

Appendices

A Features for Semantic Labeling

Below are the description of features used by SERENE semantic model. Note that features are extracted from the name and values of each attribute in the data source. Thus, all the data from a single column is reduced to some feature vector of the concatenation of values computed by the extractors below. We refer to a single value from the column as an entry. For example, in the table below, John is an entry of the column Person.

Table 2: Example table.

Person	Age	Gender
John	20	M
Jane	30	F

Table 3: Feature extractors that produce a scalar value

Feature name	Description
num-unique-vals	A simple count of unique entries of a column
prop-unique-vals	The proportion of unique entries of a column.
prop-missing-vals	The proportion of all missing/empty entries from all entries in a column.
ratio-alpha-chars	The average of the proportion of alphabetic characters from each entry.
prop-numerical-chars	The average of the proportion of numeric characters from each entry.
prop-whitespace-chars	The average of the proportion of whitespace characters from each entry.
prop-entries-with-at-sign	The proportion of entries with an '@' sign.
prop-entries-with-hyphen	The proportion of entries with a '-'.
prop-range-format	The proportion of entries which follow some numerical range format ([0-9]+)-([0-9]+)
is-discrete	A binary indicator if the column entries are discrete.
entropy-for-discrete-values	The entropy of the entries of a column.
shannon-entropy	Shannon’s entropy for the vocabulary specified in the description of the feature ‘char-dist-features’

Table 4: Feature extractors that produce a vector value. Some of these are computed together to reduce computation time. Others are values that are computed for all classes.

Feature name	Description
inferred-data-type	Produces a vector which contains the inferred data type. Only a single element is non-zero in the vector (1-of-K format). The set of inferred types are in float, integer, long, bool, date, time, datetime, string.
char-dist-features	Produces a vector for character distribution of column values; considered characters are printable characters
stats-of-text-length	Computes the mean, median, mode, min, max of string lengths from a columns entries.
stats-of-numerical-type	Computes the mean, median, mode, min, max of entries from a column of numerical type.
prop-instances-per-class-in-knearestneighbours	This finds the k columns from the training set with most similar names. It then produces a vector of class proportions from these k neighbours. Parameters: “num-neighbours”: The number of neighbours to find (k).
mean-character-cosine-similarity-from-class-examples	For some instance t , computes the character distribution from this instance’s entries, and produces a vector of the average cosine similarity of t with training instances of each class.
min-editdistance-from-class-examples	For some instance t , produces a vector of the minimum column name edit distance between t and all training instances of each class. Parameters: “cache-size”: size of cache to store edit distances of words; “max-comparisons-per-class”: max number of training instances to compare for each class.
min-wordnet-jcn-distance-from-class-examples	For some instance t , produces a vector of the minimum column name JCN WS4J distance between t and all training instances of each class. Parameters: “cache-size”: size of cache to store wordnet distances of words; “max-comparisons-per-class”: max number of training instances to compare against for each class.
min-wordnet-lin-distance-from-class-examples	For some instance t , produces a vector of the minimum column name LIN WS4J distance between t and all training instances of each class. Parameters: “cache-size”: size of cache to store wordnet distances of words; “max-comparisons-per-class”: max number of training instances to compare against for each class.

B MINIZINC Model

```
1  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2  %%                               %%
3  %% MiniZinc model to solve the Rel2Onto problem           %%
4  %% as a composition of Steiner Tree Problem             %%
5  %% and a matching problem                               %%
6  %%                               %%
7  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
8
9  include "globals.mzn";
10
11  %%%%%%%%% Steiner Tree Problem
12  int: nbV;                               % Nb of nodes in alignment graph
13  int: nbE;                               % Nb of edges in alignment graph
14  set of int: nodes = 1..nbV;
15  set of int: edges = 1..nbE;
16  set of int: cnodes;                     % Class nodes in alignment graph
17  set of int: dnodes;                     % Data nodes in alignment graph
18  array [nodes] of var bool: vs;          % Is a node of the alignment graph in the solution?
19  array [edges] of var bool: es;          % Is an edge of the alignment graph in the solution?
20  array [nodes,edges] of bool: adjacent; % adjacent[n,e] <=> node n is an endnode of edge e
21  array [edges] of nodes: heads;          % Endpoint 1 of each edge (the graph is undirected)
22  array [edges] of nodes: tails;          % Endpoint 2 of each edge
23  array [1..2,edges] of int: endpairs = array2d(1..2,edges,tails++heads); % Endpoints of each edge
24  array [edges] of int: ws;               % Weights of edges in alignment graph
25  var 0..sum(ws): w;                     % Weight of the solution corresponding to the alignment part
26  constraint steiner_tree(vs,es,adjacent,endpairs,w,ws); % Steiner Tree constraint
27  constraint forall (d in dnodes)
28      (vs[d] <-> sum(e in edges where adjacent[d,e]) (es[e]) == 1); % Only 1 edges per data node
29
30  %%%%%%%%% Matching Problem
31  int: nbA;                               % Nb of attribute nodes
32  set of int: atts = 1..nbA;
33  array [atts] of set of dnodes: attribute_domains; % Data nodes to which each attribute can attach (usually all)
34  array [atts] of var dnodes: match;        % Match between Attributes and Data nodes
35  array [atts,dnodes] of int: match_costs; % Weight of match edges
36  constraint forall (a in atts) (match[a] in (attribute_domains[a])); % Must match someone in the domain
37  constraint forall (a in atts) (vs[match[a]]); % The data node matched must be in the subgraph of the alignment
38  constraint alldifferent(match);           % Exact match
39  var 0..sum(a in atts) (sum(n in dnodes) (match_costs[a,n])): wm; % Total cost from matching part :
40  constraint wm = sum(a in atts) (match_costs[a,match[a]]);
41
42
```

Model 1: MINIZINC model (part 1 of 2)

```

1 %##### Unknown nodes
2 int: nb_unknown_nodes; % Number of unknown data nodes
3 set of int: unknowns = 1..nb_unknown_nodes;
4 array [int] of int: all_unk_nodes; % List of all unknown nodes (first two elements are class nodes, by convention)
5 int: node_root = all_unk_nodes[1]; % "Root" class node
6 int: node_unk = all_unk_nodes[2]; % "Unknown" class node
7 constraint if nb_unknown_nodes > 0 then
8     vs[node_unk] <-> vs[node_root] % If class node "Unknown" is used, class node "Root" is used as well
9 /\ vs[node_root] <-> sum(e in edges where adjacent[node_root,e]) (es[e] in {1,2})
10 else true %Ignore constraint
11 endif;
12
13 %##### Patterns
14 int: nb_patterns; % Number of patterns
15 set of int: PATTS = 1..nb_patterns;
16 array[PATTS] of int: patterns_w; % Support of patterns
17 array[PATTS] of set of int: patterns; % The actual patterns
18 array[PATTS] of var bool: c_p; % Is pattern chosen to be part of the solution?
19 constraint forall (p in PATTS) ((sum (e in patterns[p]) % All edges of a pattern are used <-> pattern is used
20 (es[e] == card(patterns[p])) <-> c_p[p]);
21 % Total prize collected from patterns :
22 var 0..sum(patterns_w): wp;
23 constraint if nb_patterns = 0 then wp = 0 else wp = sum(p in PATTS) (c_p[p] * patterns_w[p]) endif;
24
25 %##### Final objective function
26 var int: obj = w + wm - wp; % Total = (cost of semantic model) - (prize collected from patterns)
27
28 %##### Search strategy
29 array[atts,dnodes] of int: match_costs_sorted; % Match of attributes, sorted by cost.
30 ann: cheap_matches = int_search([ match[a] == match_costs_sorted[a,i] | a in atts, i in dnodes],
31 input_order,indomain_max,complete);
32
33 array[edges] of int: sorted_edges = arg_sort(ws);
34 ann: cheap_edges = int_search([es[sorted_edges[nbE + 1 - e]] | e in edges],input_order,indomain_min,complete);
35
36 ann: cheap_edges_fill_pattern; % Custom search strategy for patterns (implemented in Chuffed)
37 ann: custom(ann, array[int] of int, array[int] of var bool);
38 array[int] of ann: cefp =
39 if nb_patterns == 0 then [] % This builds a list in a format understood by Chuffed, to work on the "es" vars
40 else [custom(cheap_edges_fill_pattern,[nbE]++[nb_patterns]++ws++[ card(patterns[p]) | p in PATTS]++patterns_w++arrayld([
41 e | p in PATTS, e in patterns[p]]),es)] endif;
42
43 % The solver will use these strategies in order:
44 % 1) Try using the cheapest matches, until a complete match is found (i.e. for all attributes)
45 % 2) If there are patterns, use a specific search called "cheap_edges_fill_pattern". (Implemented in Chuffed only)
46 % It will start adding edges in increasing order of cost, while trying to fill patterns that are "worth" filling.
47 % E.g. We already set es[x] = true for x in {e1,e2,e3}; if adding ws[e5] = 4 and ws[e7] = 6, but
48 % {e1,e2,e3,e7} form a pattern p of patterns_w[p] = 3, then ws[e7] - patterns_w[p] < ws[e5], so e7 is selected first
49 % If there are no patterns, this strategy is ignored.
50 % 3) Try to eliminate edges by decreasing cost: it will eliminate the more expensive edges first, until
51 % the strictly necessary cheap edges are added by the propagators. This is used only when there are no
52 % patterns, otherwise (2) already selects the edges.
53
54 solve ::seq_search([cheap_matches]++cefp++[cheap_edges])
55 minimize obj;

```

Model 2: MINIZINC model (part 2 of 2)