

# Entity-Attribute Resolution

Decision-theoretic  
symbol matching 102



## Relevant Work

### **A Distance Based Approach to Entity Reconciliation in Heterogeneous Databases**

- D. Dey and S. Sarkar (200109)
- Uses decision theory and user-solicited distances to determine yes/no entities in disparate databases are same.
- Assumes that semantic-heterogeneity of attributes has already been resolved.

# Relevant Work

## **The Field Matching Problem: Algorithms and Applications**

- A. Monge and C. Elkan (199608)
- Uses text matching algorithms to determine contents of attributes are same
- Does not work with symbolic data
- No subjective measure of success

## Relevant Work

### Automatic Ontology Mapping for Agent Communication

- F. Wiesman, N. Roos, P. Vogt (200207)
- Establishes mappings using a *joint attention set* based on word co-occurrence
- Cannot resolve symbolic data
- No metric for deciding whether translation was feasible

# Relevant Work

## Algorithms for Ontological Mediation

- E. Campbell and S. Shapiro (199803)
- “Common words have common meanings”
- Uses WordNet lexical ontology to translate between two ontologies.
- Does not operate on symbolic data

# Formulation – Records

Defining a record (of an event)...

- Let  $\mathbf{G}$  be a finite, discrete set of “global keys”
- Let  $\mathbf{K}$  be a finite, discrete set of “prediction keys”
- Let  $\mathbf{S}$  be a finite, discrete set of “symbols” that encode an “attribute-of-interest”
- A record  $\mathbf{x}$  is a triple of the form  $\langle \mathbf{g}_i, \mathbf{k}_i, \mathbf{s}_i \rangle$  where  $\mathbf{g}_i \in \mathbf{G}$ ,  $\mathbf{k}_i \in \mathbf{K}$  and  $\mathbf{s}_i \in \mathbf{S}$

# Formulation - Metadatabase

Define a metadatabase as:

- Let  $\mathbf{P}(\mathbf{K})$  denote the distribution of  $\mathbf{K}$  over the domain of discourse
- Let  $\mathbf{P}(\mathbf{S})$  denote the distribution of  $\mathbf{S}$  over the domain of discourse.
- Let  $\mathbf{C}$  be a set of conditional probabilities such that  $\exists \mathbf{P}(\mathbf{s}_i | \mathbf{k}_j)$  such that  $\mathbf{P}(\mathbf{s}_i | \mathbf{k}_j) \neq \mathbf{P}(\mathbf{s}_i)$
- A metadatabase **META** is the tuple...  
 $\langle \mathbf{G}, \mathbf{K}, \mathbf{P}(\mathbf{K}), \mathbf{S}, \mathbf{P}(\mathbf{S}), \mathbf{C} \rangle$

## Formulation – Global Keys

Elements of  $G$  are globally unique identifiers.

- For any two records  $r_i$  and  $r_j$  in a set of records,  $R$ :  $g_i \neq g_j$



# Formulation - Database

A database is then defined as follows...

- A set of records  $\mathbf{R} = \{r_1, \dots, r_n\}$
- A database  $\mathbf{D}$  is the tuple  $\langle \mathbf{M}, \mathbf{R} \rangle$  where  $\mathbf{M}$  is a metadatabase describing the the domain of discourse and  $\mathbf{R}$  is a set of recorded observations of events in the domain of discourse.

# Formulation - Agent

An agent is an entity that...

- Has knowledge about its environment
- is **active**
- is **autonomous**
- seeks **new information resources**
- **acts on this information** to increase utility

# Formulation – Agent – Observation

We allow for the possibility that an agent's observations are not perfect

- This is represented using conditional probabilities of the form  $\mathbf{P}(\mathbf{s}_o | \mathbf{s}_a)$
- Let **DISCRIMINATOR** be an  $\mathbf{S} \times \mathbf{S}$  matrix encoding the conditional probabilities.

# Formulation – Agent – Actions

An agent takes action based upon observation of objects or events in the domain-of-discourse.

- Let the actions (strategies) available to the agent be the set **ACTIONS** =  $\{a_1, \dots, a_n\}$

# Formulation – Agent – Outcomes

Upon the execution of some action, an agent may experience any one of a set of finite outcomes.

- Let **OUTCOMES** be the finite, discrete set of possible results from the execution of some action.

# Formulation – Agent – Payoffs

A payoff is what an agent experiences when executing a particular action on the stimulus of some event or object.

- Let **PAYOFFS** be a matrix of **S** × **OUTCOMES**
- A payoff **PAYOFF** (**a** | **s**) is some distribution over **OUTCOMES**.

# Formulation – Agent – Information

An agent may have one or more information sources.

- Let **INFO** be a set of tuples of the form  $\langle \mathbf{t}_i, \mathbf{R}_i \rangle$  where  $\mathbf{t}_i$  is a translator that maps the attribute of interest of records in  $\mathbf{R}_i$  into the symbol set **S** of the agent's metadatabase.

# Formulation - Agent

An agent is thus formulated as the tuple of...

**A = <META,  
DISCRIMINATOR,  
ACTIONS,  
OUTCOMES,  
PAYOFFS,  
INFO>**



# Formulation – Agent - Operations

We will need these functions for later computation:

- **choice** ( $\mathbf{A}, \mathbf{s}$ ) – Agent  $\mathbf{A}$ 's optimal strategy for acting on observation of object of class  $\mathbf{s}$  with certainty.
- **action** ( $\mathbf{A}, \mathbf{s}_o, \mathbf{s}_a$ ) – The distribution of outcomes of agent  $\mathbf{A}$  executing **choice** ( $\mathbf{A}, \mathbf{s}_o$ ) on an event/object of class  $\mathbf{s}_a$ .

## Formulation - Translation

- A translation function  $\tau$  is a mapping defined over  $S_r \times S_1: \tau(S_r) \rightarrow S_1$
- Let  $|S_r| = m$
- Let  $|S_1| = n$
- The set of all translation functions is denoted  $T_{S_r S_1}$  and contains  $m^n$  translations.
- Translation is also defined on records where the attribute of interest is translated through  $\tau$ .

# Formulation - Translation

Some useful functions for later...

- **default** ( $S_r$ ) – Constructs a default translator by mapping all elements of  $S_r$  to the unknown symbol.
- **update** ( $t, s_r, s_1$ ) – Produces a new translator by associating the remote symbol  $s_r$  with the local symbol  $s_1$

# Statement

- Given an agent under local control  $\mathbf{A}_1$
- Given a remote set of records  $\mathbf{R}_r$  with  $\mathbf{K}_1 \cap \mathbf{K}_r \neq \{\}$
- Find translation function  $\mathbf{t}$  in  $\mathbf{T}_{s_1 s_r}$  such that  $\Sigma \mathbf{E}[\mathbf{action}(\mathbf{A}_1, \mathbf{t}(\mathbf{r} \in \mathbf{R}_r))]$  is maximal.

## Solution – The Unknown

We begin the solution by augmenting  $\mathbf{s}_r$  with a new symbol (?) that indicates a complete lack of knowledge.

Likewise, the **ACTIONS** set of  $\mathbf{a}_1$  must also be augmented with a null strategy which is executed only on observation of the unknown.

## Solution – E[Value] of a Record

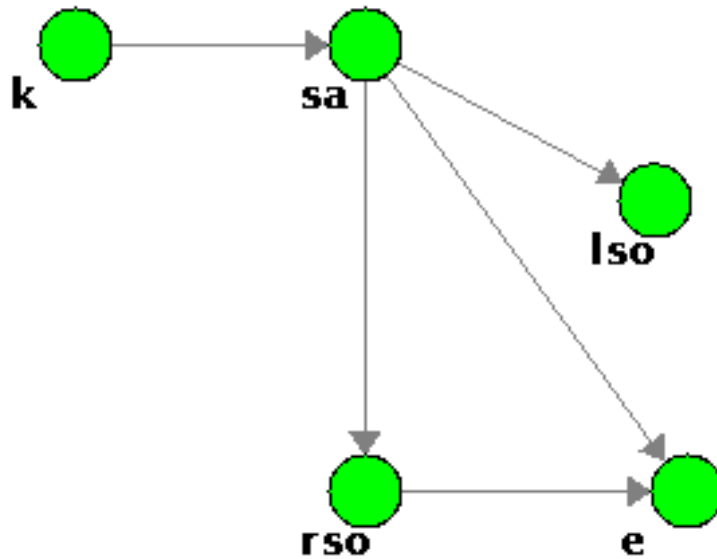
To compute the aggregate value of a translation, we first develop a function to determine the expected value of acting on an observation.

## Solution – E[Value] of a Record

The process will consist of:

1. Apply translation:  $\mathbf{r}_t \leftarrow \mathbf{t}(\mathbf{r})$
2. Select the best strategy:  
**choice(A, class( $\mathbf{r}_t$ ))**
3. Use a BBN to compute the distribution of **class( $\mathbf{r}_t$ )** over **S**.
4. Compute the expected value of executing this strategy.

# Solution – E[Value] of a Record



BBN for computing distribution of class  $s_a$ .



## Solution – E[Value] of a Record

- $k$  is `predictionKey(r)`
- $s_a$  is the *actual* class of  $r$ , rather than the observed class
- $rso$  is the observed class of the record after translation from  $S_r$  to  $S_1$  and accounting for possible error via the discrimination matrix  $DISCRIMINATOR_r$

## Solution – Evidence - Joint Sets

The BBN can incorporate evidence from objects in a *joint set*.

The joint set,  $\text{joint}(\mathbf{R}_1, \mathbf{R}_2)$  between any two record sets  $\mathbf{R}_1$  and  $\mathbf{R}_2$  is the set of all objects with a global key that is in both  $\mathbf{R}_1$  and  $\mathbf{R}_2$

For a record  $r$  that is known by AI, we can instantiate the `iso` node to add evidence to the translation.

## Solution – Evidence - Experiments

Additional evidence can be accumulated by taking action on a translated record.

The observed outcome of the experiment is instantiated on the  $e$  node of the BBN.

# Solution – Class Compression

Under certain conditions, we can compress the space of the translation and reduce computational efforts by pruning.

- Let  $s_1$  and  $s_2$  be elements of  $S_1$ .
- $s_1$  subsumes  $s_2$  if...  
$$\text{action}(A_1, s_1, s_2) \geq \text{action}(A_1, s_1, s_1)$$
- Let  $s_1' = s_1 \cup s_2$

## Solution – Class Compression

Subsumption is transitive, allowing a hierarchy to be constructed over  $S_1$  by bottom-up association.

The compression of classes by subsumption is a function of an agent's perceptions, making the constructed hierarchy independent of any remote entities.

# Solution – Class Compression

We introduce the operations...

- **basis** ( $\mathbf{A}_1$ ) – The classes of  $\mathbf{S}_1$  that cannot be subsumed.
- **next** ( $\mathbf{S}_{1s}, \mathbf{s}_1$ ) – The immediate successors of  $\mathbf{s}_1$  in the subsumption  $\mathbf{S}_{1s}$  hierarchy of  $\mathbf{A}_1$ .

## Solution – Algorithm – $t$ value

Using what has so far been developed we can construct a function to determine the expected utility gain of applying translator  $t$  to remote data set.

Let the function  $tvalue(t, R_r, A_1, S_{1s})$  be defined as follows...

# Solution – Algorithm – t value

```
value ← 0
for r in Rr:
    k ← predictionKey(r)
    rso ← class(r)
    rsot ← t(rso)
    d ← BBN(A1, k, rsot)
    for s1 in S1s:
        value ← value + E[action(A1, rsot, s1)]
```



# Solution – Algorithm – t value

Where...

- $t$  is a translator
- $R_r$  is the remote record set
- $A_1$  is the local agent
- $S_{1s}$  is some compression (including no compression) of  $S_1$ .

# Solution – Algorithm

```
t ← default(Sr)
maxv ← tvalue(t, Rr, A1, S1)
for sr in Sr:
  q ← basis(A1)
  while q not empty:
    s1 ← pop(q)
    t' ← update(t, sr, s1)
    tv ← tvalue(t, Rr, A1, S1s)
```

# Solution - Algorithm

```
if tv > maxvalue:  
    maxvalue ← tv  
    t ← t'  
    q ← next (S1s', s1)
```

At the end of the algorithm,  $\mathbf{t}$  holds the best possible translator of the set  $\mathbf{T}_{s_1s_r}$

## Future work – Better Evidence

Currently, the evidence gathered by experimentation is underutilized.

Evidence is only propagated to a single record when it should enhance confidence in the entire translation.

## Future work – EOG Graphs

Once attribute resolution has been concluded, EOG relationships between objects should be found.

This extension will allow taxonomies/ontologies to be resolved for a vastly more complete method for knowledge sharing.

## Future work – Compound Classes

An agent's actions naturally partition a schema into classes. These classes can be constructed along multiple attributes. Eg. An agent takes action a on  $\langle M,+ \rangle$  and b on  $\langle F,0 \rangle$

Simultaneous resolution of multiple attributes is much more complex than single attribute resolution, yet much more useful.

# Conclusions

- The decision theoretic method presented here complements the linguistic and textual methods used more commonly in other research efforts.
- By explicitly incorporating a set of strategies and payoffs, we allow an agent to act with high autonomy when such actions are risky.

# Conclusions

- The method for computing a translation's expected value has been improved by accounting for joint sets and experimental evidence.
- Class compression can lead to pruning of the search space with no loss of optimality. The value of this improvement will be more noticeable when resolving compound attributes.