

# Schema Matching on Streams with Accuracy Guarantees

Szymon Jaroszewicz<sup>1</sup>, Lenka Ivantysynova<sup>2</sup>, and Tobias Scheffer<sup>2</sup>

<sup>1</sup>Technical University of Szczecin, Poland, sj@cs.umb.edu

<sup>2</sup>Humboldt-Universität zu Berlin, Germany

{ivantysy, scheffer}@informatik.hu-berlin.de

## Abstract

*We address the problem of matching imperfectly documented schemas of data streams and large databases. Instance-level schema matching algorithms identify likely correspondences between attributes by quantifying the similarity of their corresponding values. However, exact calculation of these similarities requires processing of all database records—which is infeasible for data streams. We devise a fast matching algorithm that uses only a small sample of records, and is yet guaranteed to match the most similar attributes with high probability. The method can be applied to any given (combination of) similarity metrics that can be estimated from a sample with bounded error; we apply the algorithm to several metrics. We give a rigorous proof of the method’s correctness and report on experiments using large databases.*

## 1. Introduction

In many practical situations, data mining processes are preceded by a preprocessing step in which multiple heterogeneous transaction databases are integrated into a single data warehouse. Integration requires the identification of semantic correspondences between attributes across multiple transaction database schemas [2]. This schema matching process is not only a preprocessing step to other analysis steps, but constitutes a knowledge discovery challenge in itself. In practical applications that often involve large databases, schema matching is a tedious and complex task, often further hindered by lack of documentation.

Research revolves around the question how schema matching can be supported, or even automated, effectively. Schema matching algorithms generally match elements of the distinct schemas such that some similarity criterion is maximized. Schema-level matchers utilize a similarity criterion that refers to schema information such as attribute names, descriptions, or the schemas’ structure. Unfortunately, the available schema information is often insuffi-

cient. Cases arise in which attribute names are opaque and clues on their semantics can only be found by inspecting the values. Instance-level schema matchers quantify the similarity of attributes by comparing properties of their values.

In order to guarantee that the produced mapping in fact maximizes the chosen similarity function, instance-level matchers need to exercise at least one entire pass over the databases. A complete pass is unreasonable for transactional databases that record high-speed data streams as they occur, for instance, in retail chains, and telecommunication and banking applications and can easily grow into the order of hundreds of terabytes. The ad-hoc solution of processing only a small sample of transactions results in a loss of any guarantee on the optimality of the produced match. Known methods that utilize sampling do not guarantee any properties of the retrieved mappings.

We formalize the matching problem in a way that is both mathematically rigorous and closely tied to practical applications. We devise a schema matching algorithm, which is *guaranteed* to produce a “near-optimal” match; we detail the term “near-optimal” by means of probabilistic bounds on the similarity of matched attributes. The guarantee holds even though the algorithm relies on sampling.

A useful similarity metric for the matching problem at hand is a base requirement for any sampling algorithm, and designing such a problem-specific metric clearly is a challenging issue in itself. Therefore, we design our solution to be generic with respect to the similarity criterion; weighted combinations of several schema-level and instance-level criteria can be applied. We review a selection of criteria, many more can be inserted into the algorithm.

The rest of this paper is organized as follows: after reviewing related work in Section 2, we formalize the problem in Section 3. We phrase the schema matching problem in a way that pays tribute to typical practical application scenarios, and is sufficiently rigid to allow us to state and rigorously prove the correctness of our solutions. We detail our solution in Section 4, discuss attribute similarity metrics in Section 5 and report on experimental results on real-world data in Section 6. Section 7 concludes.

## 2. Related Work

Bernstein and Rahm [16] provide an overview of the schema matching problem [10, 13] and a taxonomy of schema matching algorithms. Matchers can be dichotomized into schema-level and instance-level methods. Schema-level matchers maximize some similarity function that refers to attribute names and other structural information (e.g., [15, 12, 6, 4]) whereas the similarity metrics employed by instance-level approaches (e.g., [10, 1, 14]) refer to the instance data; that is, to the attributes' values.

The instance-level approach can even be applied in the complete absence of useful schema-level information [9] and it allows to even identify complex matching relations – e.g., a factor-equivalence relation between an income attribute in Euros and a corresponding salary in Yen [3].

Most practical systems – such as Clio [19], SemInt [11], and LSD [5] – combine schema- and instance-level similarity metrics. Assembling and balancing the similarity metric requires experience and domain knowledge; machine learning [5] and experiments on synthetic matching problems can guide the parameter optimization [17]. The algorithm that we present supports the possibility of combining schema- and instance-level similarity metrics. However, accessing the database records is only required for calculation of instance-level metrics.

Our algorithm relies on sampling and on data-dependent confidence bounds. Sampling strategies play an important role for a range of data mining tasks. They have been utilized to find the approximately most interesting association rules [7, 18], or the most significant differences between graphical models and corresponding databases [8].

## 3. Approximately Optimal Schema Matching

We focus on the following *schema matching problem setting*. We are given two schemas  $R$  and  $S$  over attributes  $\langle r_1, \dots, r_n \rangle$  and  $\langle s_1, \dots, s_m \rangle$ , respectively. The schemas contain imperfect information about the type and meaning of the attributes. The goal of schema matching is to identify pairs of attributes  $(r_i, s_j)$  that are semantically identical.

In order to approach this problem computationally, a similarity metric  $f_{R,S}(r_i, s_j)$  quantifies the similarity of attributes  $r_i$  and  $s_j$ . The similarity can exploit schema-level information such as attribute names if such information is available, and instance-level information. Instance-level similarity metrics refer to the values that occur in  $R$  and  $S$ . Naturally, instance-level information is always available—unless the databases are empty or inaccessible.

The contribution of this paper is relevant when the similarity metric  $f_{R,S}(r_i, s_j)$  involves instance-level information. In order to evaluate an instance-level metric exactly, it

is generally necessary to exercise a pass over the databases  $R$  and  $S$ ; the main challenge that we address is that this is impossible for streams. The similarity criterion  $f_{R,S}(r_i, s_j)$  is problem-specific and given in advance. We will only have to assume that it is possible to obtain an estimate of  $f_{R,S}$  that lies within a bounded confidence interval.

In practice, the schema matching process is highly interactive. A user typically seeks to find out which attributes in  $R$  have counterparts in  $S$  and identify these counterparts. We phrase this problem setting as follows. *We seek to construct an algorithm that finds, for every attribute in  $R$ , a set of the  $k$  best matching candidates in  $S$ , sorted by their similarity to the attribute from  $R$ , and vice versa for all attributes in  $S$ . In addition, the algorithm has to order all attributes in  $R$  and  $S$  by their likelihood of having a counterpart in the other schema.*

Our formulation of the problem setting builds upon  $\varepsilon$ -correct orderings of sequences. In an  $\varepsilon$ -correct ordering, values are sorted in decreasing order, but in case of near-ties between adjacent values (similarity differences of no more than  $\varepsilon$ ), any ordering is considered correct. The special case of  $\varepsilon = 0$  corresponds to a regular ordering.

**Definition 1 ( $\varepsilon$ -Correct Ordering)** *A sequence of values  $(x_1, \dots, x_k)$  is in  $\varepsilon$ -correct ordering if, for all  $i$  between 1 and  $k$  and all  $j$  between 1 and  $i - 1$ , the inequality  $x_j \geq x_i - \varepsilon$  is satisfied.*

In order to formalize the *approximately optimal matching problem* we refer to stochastic approximations with confidence bounds. Intuitively, the parameters  $\varepsilon$  and  $\delta$  that are involved have the following meaning. When the algorithm is restarted multiple times, then the resulting match is “ $\varepsilon$ -close” to the optimal match in a fraction of at least  $1 - \delta$  of all runs.

**Definition 2 (Approximately Optimal Matching)** *Given schemas  $R$  and  $S$  over attributes  $\langle r_1, \dots, r_n \rangle$  and  $\langle s_1, \dots, s_m \rangle$  respectively, a desired number of matches  $k$ , approximation and confidence parameters  $\varepsilon$  and  $\delta$ , find a sequence of  $k$  attributes  $(s_{i_1}, \dots, s_{i_k})$  for each  $r_i$  in  $R$ , a sequence of  $k$  attributes  $(r_{j_1}, \dots, r_{j_k})$  for each  $s_j$  in  $S$ , and permutations  $\pi^r$  and  $\pi^s$  such that, with confidence  $1 - \delta$ ,*

1. *for all attributes  $r_i$ , the retrieved sequence  $(s_{i_1}, \dots, s_{i_k})$  contains the best matching attributes (accurately up to  $\varepsilon$ ) and, analogously, the sequence  $(r_{j_1}, \dots, r_{j_k})$  contains the best matching attributes for each  $s_j$ ; i.e., there is no offending  $s_{i'} \notin (s_{i_1}, \dots, s_{i_k})$  such that  $f_{R,S}(r_i, s_{i'}) > \min_{i'' \in \{i_1, \dots, i_k\}} f_{R,S}(r_i, s_{i''}) + \varepsilon$  and no  $r_{j'} \notin (r_{j_1}, \dots, r_{j_k})$  such that  $f_{R,S}(r_{j'}, s_j) > \min_{j'' \in \{j_1, \dots, j_k\}} f_{R,S}(r_{j''}, s_j) + \varepsilon$ ;*

2. for all attributes  $r_i$ , the corresponding  $(s_{i_1}, \dots, s_{i_k})$  are ordered by their similarity to  $r_i$  and for all  $s_j$ , the corresponding  $(r_{j_1}, \dots, r_{j_k})$  are ordered by their similarity to  $s_j$ ; i.e.,  $(f_{R,S}(r_i, s_{i_1}), \dots, f_{R,S}(r_i, s_{i_k}))$  and  $(f_{R,S}(r_{j_1}, s_j), \dots, f_{R,S}(r_{j_k}, s_j))$  are  $\varepsilon$ -correct orderings;
3. the permutations  $\pi^r$  and  $\pi^s$  sort the attributes of  $R$  and  $S$ , respectively, according to their likelihood of having a matching partner in the other schema; i.e., the sequences  $(\max_{j'} f_{R,S}(r_{\pi^r(1)}, s_{j'}), \dots, \max_{j'} f_{R,S}(r_{\pi^r(n)}, s_{j'}))$  and  $(\max_{i'} f_{R,S}(r_{i'}, s_{\pi^s(1)}), \dots, \max_{i'} f_{R,S}(r_{i'}, s_{\pi^s(m)}))$  are  $\varepsilon$ -correct orderings.

## 4. Schema Matching Algorithms

We will present two algorithms that differ from each other in the way they derive the sample size required to solve the approximately optimal matching problem.

The first algorithm, called FSM, determines a “worst-case” sample size without accessing the data; the sample size only depends on the similarity metric and parameters  $\varepsilon$  and  $\delta$ . The sample size determined suffices to assure the quality guarantee for any possible stream; it may be pessimistically large for a given, particular stream.

The second algorithm is called progressive FSM. Its sample size depends on the database at hand. The sample-dependent bounds are technically more involved than the worst-case bounds; therefore, FSM is a natural baseline for progressive FSM and one of the questions that we will investigate empirically is whether progressive FSM outperforms the simpler FSM algorithm.

Both algorithms use *similarity confidence intervals* which bound the similarity  $f_{R,S}$ . These confidence intervals are of salient importance to our solution because they bound the inaccuracy incurred by the estimation process. A similarity confidence interval lays out a range of values such that, with a probability of at least  $1 - \delta$ , the true similarity lies within that range. The benefit is that a confidence interval can be calculated based on a random sample of any size whereas a full database pass is required to determine the true similarity value. The confidence interval is wide for small samples and tightens for increasingly large samples. Similarity confidence intervals can be constructed for a wide range of instance-level similarity functions.

**Definition 3 (Similarity Confidence Interval)** Let  $f_{R,S}$  be a similarity function. Let  $R', S'$  be random samples from  $R$  and  $S$ ; Then,  $\lfloor f_{R',S'} \rfloor_\delta$  and  $\lceil f_{R',S'} \rceil_\delta$  form a similarity confidence interval if, with probability at least  $1 - \delta$ , there is no pair of any  $r_i \in R$  and  $s_j \in S$  that violates

$$\lfloor f_{R',S'}(r_i, s_j) \rfloor_\delta \leq f_{R,S}(r_i, s_j) \leq \lceil f_{R',S'}(r_i, s_j) \rceil_\delta. \quad (1)$$

**Table 1. FSM: Fast Schema Matching.**

**Input:** Schemas  $R$  and  $S$ , respectively; desired number of candidates  $k$ ; approximation and confidence parameters  $\varepsilon$  and  $\delta$ .

1. **Let**  $N$  be the smallest number such that  $\lceil f_{R',S'}(\cdot, \cdot) \rceil_\delta - \lfloor f_{R',S'}(\cdot, \cdot) \rfloor_\delta \leq \varepsilon$  when  $R'$  and  $S'$  contain  $N$  records. (Where we parameterize the bound with “ $\cdot$ ”, this means “for any argument”.)
2. Draw  $N$  records at random from  $R$  and  $S$ , let them be  $R'$  and  $S'$ .
3. **For** all  $(r_i, s_j)$ , determine  $\hat{f}_{R',S'}(r_i, s_j) = \frac{1}{2}(\lceil f_{R',S'}(r_i, s_j) \rceil_\delta + \lfloor f_{R',S'}(r_i, s_j) \rfloor_\delta)$  based on the samples.
4. **For** all attributes  $r_i$ : **Let**  $H_{r_i}^*$  be the  $k$  elements  $s_j$  that maximize  $\hat{f}_{R',S'}(r_i, s_j)$ ; **For** all attributes  $s_j$ : **Let**  $H_{s_j}^*$  be the  $k$  elements  $r_i$  that maximize  $\hat{f}_{R',S'}(r_i, s_j)$ .
5. **Let**  $\pi^r$  sort the attributes in  $R$  by  $\max_{s_{i'} \in H_{s_{i'}}^*} \hat{f}_{R',S'}(r_i, s_{i'})$ ; **let**  $\pi^s$  sort the attributes in  $S$  by  $\max_{r_{j'} \in H_{r_{j'}}^*} \hat{f}_{R',S'}(r_{j'}, s_j)$ .
6. **Return**  $\pi^r, \pi^s$ , and for all attributes  $p$  in  $R \cup S$ : return the sequence of elements in  $H_p^*$  sorted by  $\hat{f}_{R',S'}(p, q)$  as  $k$  best matches.

With probability  $1 - \delta$ , the true similarity of *all* the pairs  $(r_i, s_j)$  are within their bounds. The probability of any violation by at least one pair is bounded by  $\delta$ . We will detail intervals for selected similarity functions in Section 5.

The first algorithm, “FSM”, is detailed in Table 1. In Step 1, it uses a binary search to determine the smallest sample size  $N$  that suffices to guarantee that the similarity confidence intervals are tight up to  $\varepsilon$ .  $N$  is calculated without taking any samples from the database. For a given similarity metric  $f_{R,S}$ , the corresponding confidence interval that extends from  $\lceil f_{R',S'}(\cdot, \cdot) \rceil_\delta$  to  $\lfloor f_{R',S'}(\cdot, \cdot) \rfloor_\delta$  is inserted into the algorithm. (Parameterization with “ $\cdot$ ” indicates “for any pair of attributes”.)

By contrast, the progressive FSM algorithm described in Table 2 uses sample-dependent bounds. It starts to draw batches of example records from the databases and then calculates confidence intervals in each step that depend on the sample processed so far. The progressive algorithm can terminate earlier than the sample-independent bound would have suggested, depending on characteristics of the database. The progressive algorithm terminates early if it is “easy” to confidently identify the matching partners for the attributes. Identifying matching partners is easy if each at-

**Table 2. Progressive FSM Algorithm.**

**Input:** Schemas  $R$  and  $S$ , respectively; desired number of candidates  $k$ ; approximation and confidence parameters  $\varepsilon$  and  $\delta$ .

**Let**  $R' = S' = \emptyset$ . **Let** the batch size parameter be 10,000, by default. **Let**  $t = 1$ .

**Let**  $N$  be the smallest number such that  $\lceil f_{R',S'}(\cdot, \cdot) \rceil_{\frac{\delta}{2}} - \lfloor f_{R',S'}(\cdot, \cdot) \rfloor_{\frac{\delta}{2}} \leq \varepsilon$  when  $R'$  and  $S'$  contain  $N$  records. **Let**  $T = N/\text{batch size}$ .

**For**  $t = 1 \dots T$  **Repeat:**

1. **Let**  $\delta_t = \frac{\delta}{2t(t+1)}$ .
2. Draw batches of records at random from  $R$  and  $S$ , and add them to  $R'$  and  $S'$ , respectively.
3. **For** all  $(r_i, s_j)$ , calculate  $\lfloor f_{R',S'}(r_i, s_j) \rfloor_{\delta_t}$ ,  $\lceil f_{R',S'}(r_i, s_j) \rceil_{\delta_t}$ , and  $\hat{f}_{R',S'}(r_i, s_j) = \frac{1}{2}(\lceil f_{R',S'}(r_i, s_j) \rceil_{\delta_t} + \lfloor f_{R',S'}(r_i, s_j) \rfloor_{\delta_t})$ .
4. **If** all pairs  $(r_i, s_j)$  satisfy that  $\lceil f_{R',S'}(r_i, s_j) \rceil_{\delta_t} - \lfloor f_{R',S'}(r_i, s_j) \rfloor_{\delta_t} \leq \varepsilon$  **then break**.

**For** all attributes  $r_i$ : **Let**  $H_{r_i}^*$  be the  $k$  elements  $s_j$  that maximize  $\hat{f}_{R',S'}(r_i, s_j)$ ; **For** all attributes  $s_j$ : **Let**  $H_{s_j}^*$  be the  $k$  elements  $r_i$  that maximize  $\hat{f}_{R',S'}(r_i, s_j)$ .

**Let**  $\pi^r$  sort the attributes in  $R$  by  $\max_{s_i^* \in H_{r_i}^*} \hat{f}_{R',S'}(r_i, s_i^*)$ ; **let**  $\pi^s$  sort the attributes in  $S$  by  $\max_{r_j^* \in H_{s_j}^*} \hat{f}_{R',S'}(r_j^*, s_j)$ .

**Return**  $\pi^r$ ,  $\pi^s$ , and for all attributes  $p$  in  $R \cup S$ : return the sequence of elements in  $H_p^*$  sorted by  $\hat{f}_{R',S'}(p, q)$  as  $k$  best matches.

tribute in  $R$  is similar to one, or very few, attributes in  $S$  and dissimilar from most others. Our main result is that both, the FSM and progressive FSM algorithm solve the approximately optimal matching problem.

**Theorem 1** *Let  $\lceil f_{R',S'}(\cdot, \cdot) \rceil_{\delta}$  and  $\lfloor f_{R',S'}(\cdot, \cdot) \rfloor_{\delta}$  be similarity confidence bounds that satisfy Definition 3. Then both, the FSM (Table 1) and progressive FSM algorithm (Table 2) find, for any pair of schemas  $R$  and  $S$  and any input parameters  $0 < \delta \leq 1$ ,  $0 < \varepsilon$ , and  $k$ , an approximately optimal matching as detailed in Definition 2.*

The proof of Theorem 1 is given in Appendix A. Copies of the implementation can be obtained for research purposes from the authors.

## 5. Attribute Similarity Metrics

We will now discuss three exemplary similarity measures for attributes and show how confidence bounds can be derived. We will sketch the individual bounds as lemmas; the main result of this section is the bound for the weighted sum of all presented similarity measures. Any similarity metric can be integrated into the algorithm, if it can be estimated with bounded error based on a sample.

Various similarity measures require measuring the similarity of two probability distributions. For this purpose we will use a measure  $E$  based on Euclidean distance. Let  $P = (p_1, \dots, p_n)$  and  $Q = (q_1, \dots, q_n)$  be probability distributions.  $E$  is defined as  $E(P, Q) = 2 - \sum_{i=1}^n (p_i - q_i)^2$ . We subtract the sum of squares from 2 in order to guarantee that  $E$  has high values for similar distributions, while remaining nonnegative.

Similarity measures will take two attributes as arguments and will be denoted with lowercase letter  $f$  with different superscripts. When deriving the bounds for our similarity metrics, the following (well-known) observations will be helpful. The following are upper and lower bounds on an underlying probability  $p$ , given a corresponding observed frequency (a count of events divided by the total number of observations)  $\hat{p}$ .

$$\lfloor p \rfloor_{\delta} = \max(\hat{p} - e, 0), \quad \lceil p \rceil_{\delta} = \min(\hat{p} + e, 1),$$

where  $\hat{p}$  is the value of  $p$  estimated from data. The term  $e$  can be expressed in terms of the well-known Hoeffding inequality or in terms of the inverse normal distribution.

Based on Hoeffding's inequality,  $e = \sqrt{\frac{1}{2N} \log \frac{2}{\delta}}$ ; and using the Normal distribution,  $e = z_{1-\frac{\delta}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{N}}$ , where  $N$  denotes the sample size and  $z_{1-\delta}$  is the  $1 - \delta$  quantile of the normal distribution that can be looked up in a table.

### 5.1. Euclidean Distance

Given two attributes  $r$  and  $s$  with identical domains, a natural measure of their similarity is given by

$$f_{R,S}^E(r, s) = E(P_r, Q_s), \quad (2)$$

where  $P_r$  and  $Q_s$  are probability distributions of  $r$  and  $s$  respectively.

**Lemma 1** *Let  $f_{R,S}(r_i, s_j)$  be defined as in Equation 2. Then for given schemas  $R$ ,  $S$  and all pairs of attributes  $r \in R$ ,  $s \in S$ , the following similarity confidence intervals satisfy Definition 3.*

$$\begin{aligned} & \lfloor f_{R',S'}^E(r, s) \rfloor_{\delta} \\ &= \left\{ 2 - \sum_{i=1}^{n_r} \left( \left( \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \right)^2 + \left( \lceil q_{s,i} \rceil_{\frac{\delta}{M}} \right)^2 \right) \right\} \end{aligned}$$

$$-2\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \lfloor q_{s,i} \rfloor_{\frac{\delta}{M}} \Big\}, (3)$$

$$\begin{aligned} & \left[ f_{R',S'}^E(r,s) \right]_{\delta} \\ &= \left\{ 2 - \sum_{i=1}^{n_r} \left( \left( \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \right)^2 + \left( \lfloor q_{s,i} \rfloor_{\frac{\delta}{M}} \right)^2 \right. \right. \\ & \quad \left. \left. - 2\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \lfloor q_{s,i} \rfloor_{\frac{\delta}{M}} \right) \right\}, (4) \end{aligned}$$

where  $R', S'$  are samples taken from  $R$  and  $S$  respectively,  $n_r$  and  $n_s$  are the size of the domains of  $r$  and  $s$ ,  $(p_{r,1}, \dots, p_{r,n_r})$  and  $(q_{s,1}, \dots, q_{s,n_s})$  are probability distributions of  $r$  and  $s$  respectively, and  $M = \sum_{r \in R} n_r + \sum_{s \in S} n_s$ , over categorical attributes.

The proof of Theorem 2 and its lemmas is given in Appendix B.

## 5.2. Permuted Euclidean Distance

When comparing two categorical attributes, it is not clear which values correspond to each other. Therefore, all possible permutations of values in the domains of attributes have to be considered.

Let  $\pi$  be a permutation of  $\{1, \dots, n\}$  and let  $r$  and  $s$  be categorical attributes such that  $n_r = n_s = n$ . Denote by  $\pi Q_s$  the result of applying permutation  $\pi$  to the distribution  $Q_s = (q_{s,1}, \dots, q_{s,n})$  of  $s$ , i.e.,  $(q_{s,\pi(1)}, \dots, q_{s,\pi(n)})$ . The similarity measure is now defined as

$$f_{R,S}^{\pi E}(r,s) = \max_{\pi \in \Pi(n)} E(P_r, \pi Q_s), \quad (5)$$

where  $\Pi(n)$  is the set of all permutations of  $\{1, \dots, n\}$ .

**Lemma 2** Let  $f_{R,S}(r_i, s_j)$  be defined as in Equation 5. Then for given schemas  $R, S$  and all pairs of attributes  $r \in R, s \in S$ , the following similarity confidence intervals satisfy Definition 3.

$$\left[ f_{R',S'}^{\pi E}(r,s) \right]_{\delta} = \max_{\pi \in \Pi} \left[ f_{R',S'}^E(r, \pi s) \right]_{\delta}, \quad (6)$$

$$\left[ f_{R',S'}^E(r,s) \right]_{\delta} = \max_{\pi \in \Pi} \left[ f_{R',S'}^{\pi E}(r, \pi s) \right]_{\delta}, \quad (7)$$

where  $R', S'$  are samples taken from  $R$  and  $S$  respectively,  $(p_{r,1}, \dots, p_{r,n_r})$  and  $(q_{s,1}, \dots, q_{s,n_s})$  are probability distributions of  $r$  and  $s$  respectively, and  $M = \sum_{r \in R} n_r + \sum_{s \in S} n_s$ , over categorical attributes.

The proof of Theorem 2 and its lemmas is given in Appendix B.

## 5.3. Bigrams

This measure is based on the distribution of bigrams (i.e., two character subsequences) in attribute values.

Let  $B$  be a small integer (10 in our case). For an attribute  $r$  create a synthetic attribute  $r'$  computed as follows: for a given record  $t$ , pick a random bigram from the value  $t[r]$  and convert it to an integer based on ASCII values of the two characters. Take this number modulo  $B$ . The obtained number is the value of  $r'$  in the record  $t$ .

It is easy to see that  $r'$  is a categorical attribute with  $B$  values. Such attributes can easily be compared using  $f_{R,S}^E$  giving rise to the following similarity measure

$$f_{R,S}^{Bg}(r,s) = f_{R,S}^E(r',s'). \quad (8)$$

Drawing a random bigram assures that  $r'$  and  $s'$  are sampled iid (independently and from identical distributions). Using the hashing scheme with  $B$  buckets reduces the domain of the random variables compared to using the set of all bigrams which, in turn, narrows the confidence bounds.

Empirically, we find that matchings obtained using this method are similar those found without hashing; however, hashing allows for tight confidence intervals.

**Lemma 3** Let  $f_{R,S}(r_i, s_j)$  be defined as in Equation 8. Then for given schemas  $R, S$ , samples  $R', S'$  taken from those schemas and all pairs of attributes  $r \in R, s \in S$ , the following similarity confidence intervals satisfy Definition 3.

$$\begin{aligned} & \left[ f_{R',S'}^{Bg}(r,s) \right]_{\delta} \\ &= 2 - \sum_{i=1}^{n_{r'}} \left[ \left( \left\lfloor p_{r',i} \right\rfloor_{\frac{\delta}{M}} \right)^2 + \left( \left\lfloor q_{s',i} \right\rfloor_{\frac{\delta}{M}} \right)^2 \right. \\ & \quad \left. - 2\left\lfloor p_{r',i} \right\rfloor_{\frac{\delta}{M}} \left\lfloor q_{s',i} \right\rfloor_{\frac{\delta}{M}} \right] \quad (9) \end{aligned}$$

$$\begin{aligned} & \left[ f_{R',S'}^E(r,s) \right]_{\delta} \\ &= 2 - \sum_{i=1}^{n_{r'}} \left[ \left( \left\lfloor p_{r',i} \right\rfloor_{\frac{\delta}{M}} \right)^2 + \left( \left\lfloor q_{s',i} \right\rfloor_{\frac{\delta}{M}} \right)^2 \right. \\ & \quad \left. - 2\left\lfloor p_{r',i} \right\rfloor_{\frac{\delta}{M}} \left\lfloor q_{s',i} \right\rfloor_{\frac{\delta}{M}} \right] \quad (10) \end{aligned}$$

The proof of Theorem 2 and its lemmas is given in Appendix B.

## 5.4. Character Proportions

This measure compares the proportions of characters that fall into predefined subsets. We use four subsets: lowercase letters, uppercase letters, digits, and an additional set for all remaining printable characters. For instance, Clio utilizes this idea [19].

Let  $r$  be an attribute. For each record  $t$  we compute the proportions  $p_l, p_u, p_d, p_p$  of characters of  $t[r]$  belonging to each subset.

After taking the averages of  $p_l, p_u, p_d, p_p$  over all records we obtain a vector  $P^r = (P_l^r, P_u^r, P_d^r, P_p^r)$ . Elements of

this vector lie in the interval  $[0, 1]$  and sum up to 1; hence, such vectors can be compared like probability distributions. The similarity measure can now be defined as

$$f_{R,S}^c(r, s) = E(P^r, P^s).$$

Bounds for  $f_{R,S}^c$  are defined analogously to Lemma 3 but in terms of vectors  $P^r$  and  $P^s$  defined above.

### 5.5. Weighted Sum

A weighted sum of all similarity measures forms the most general case. Any weighted subset of measures can be obtained by setting weights to either zero or desired nonzero values.

$$f_{R,S}(r, s) = \frac{w_{\pi E} f_{R,S}^{\pi E}(r, s) + w_{Bg} f_{R,S}^{Bg}(r, s) + w_c f_{R,S}^c(r, s)}{w_{\pi E} + w_{Bg} + w_c} \quad (11)$$

**Theorem 2** *Let  $f_{R,S}(r_i, s_j)$  be defined as in Equation 11. Then for given schemas  $R, S$ , samples  $R', S'$  taken from those schemas, and all pairs of attributes  $r \in R, s \in S$ , the following similarity confidence intervals satisfy Definition 3.*

$$\begin{aligned} & \lfloor f_{R',S'}(r, s) \rfloor_{\delta} \quad (12) \\ &= \frac{w_{\pi E} \lfloor f_{R',S'}^{\pi E}(r, s) \rfloor_{\delta} + w_{Bg} \lfloor f_{R',S'}^{Bg}(r, s) \rfloor_{\delta} + w_c \lfloor f_{R',S'}^c(r, s) \rfloor_{\delta}}{w_{\pi E} + w_{Bg} + w_c} \end{aligned}$$

$$\begin{aligned} & \lceil f_{R',S'}(r, s) \rceil_{\delta} \quad (13) \\ &= \frac{w_{\pi E} \lceil f_{R',S'}^{\pi E}(r, s) \rceil_{\delta} + w_{Bg} \lceil f_{R',S'}^{Bg}(r, s) \rceil_{\delta} + w_c \lceil f_{R',S'}^c(r, s) \rceil_{\delta}}{w_{\pi E} + w_{Bg} + w_c} \end{aligned}$$

The proof is given in Appendix B.

### 5.6. Combining Schema- and Instance-Level

When schema-level information is available – for instance in the form of attribute names or descriptions – it is advisable to utilize this information. The similarity function can combine schema- and instance-level components:

$$f_{R,S}(r, s) = w_S f^{Schema}(r, s) + w_I f_{R,S}^{Instance}(r, s).$$

For instance,  $f^{Schema}$  can quantify the similarity of attributes' names and descriptions. Since  $f^{Schema}$  is independent of the database, it follows immediately that the bounds are

$$\lfloor f_{R',S'}(r, s) \rfloor = w_S f^{Schema} + w_I \lfloor f_{R',S'}^{Instance}(r, s) \rfloor, \text{ and}$$

$$\lceil f_{R',S'}(r, s) \rceil = w_S f^{Schema} + w_I \lceil f_{R',S'}^{Instance}(r, s) \rceil.$$

## 6. Experiments

We want to investigate empirically (a) whether the FSM and progressive FSM algorithms are practically applicable for large databases. We want to (b) compare the performances of FSM (using Hoeffding's bounds and normal bounds), progressive FSM (using normal bounds), and a baseline instance-based matcher that uses the same similarity function but processes the entire database. The baseline is not practical for streams and large databases, it always retrieves the matches that maximize the similarity on the database. Theorem 1 *guarantees* that FSM and progressive FSM return approximately optimal matches; nevertheless, we will (c) empirically study the chance of the algorithms finding the correct matches for attributes.

The similarity function used for all experiments is the weighted average of all similarity functions studied in Section 5, with all weights fixed to 1.

### 6.1. Benchmark Problems

In order to assure the reproducibility of our results, we use three publicly available databases. The KDDCup 1998 customer relationship management database contains 481 attributes and over 190,000 records; the majority of attributes are numerical but some categorical attributes are included. The KDDCup 1999 database is another standard benchmark to be audited; it includes a wide variety of intrusions simulated in a military network environment. The database has in total 42 attributes, 17 of which are almost always zero; another 14 almost always assume the value "0.00". This database contains nearly 5 million records.

The census database is an extract from the 1994 and 1995 current population surveys conducted by the U.S. Census Bureau. The database has 42 attributes in total, the majority of them contain text. There are close to 300,000 records. The execution time of both FSM and progressive FSM is constant in the desired number of solutions  $k$  (except for the final step of printing the output); the experimental results are identical for any possible value of  $k$ .

In order to conduct controlled experiments, we randomly split each of the databases into two parts containing half of the records – for the KDDCup 1998 database, we use the split that is provided with the data. We then use the schema matching algorithms to match the two halves of the databases. When the algorithm matches an attribute with itself, this is counted as a true positive. Based on the number of attributes that are matched with itself, we determine precision and recall, and F-measure.

Figure 1 shows the number of database records that FSM and progressive FSM draw from the database and process before finding an approximately optimal match. For the census and KDD Cup 1998 data, the early stopping criterion

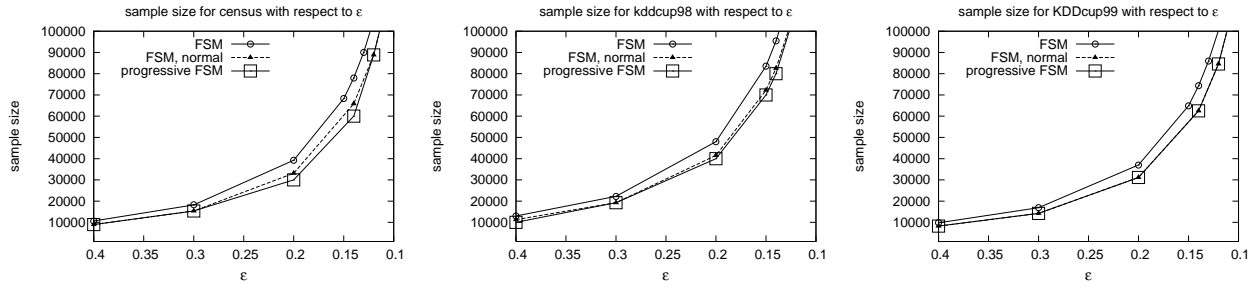


Figure 1. Sample size against  $\epsilon$  for schema matching algorithms.

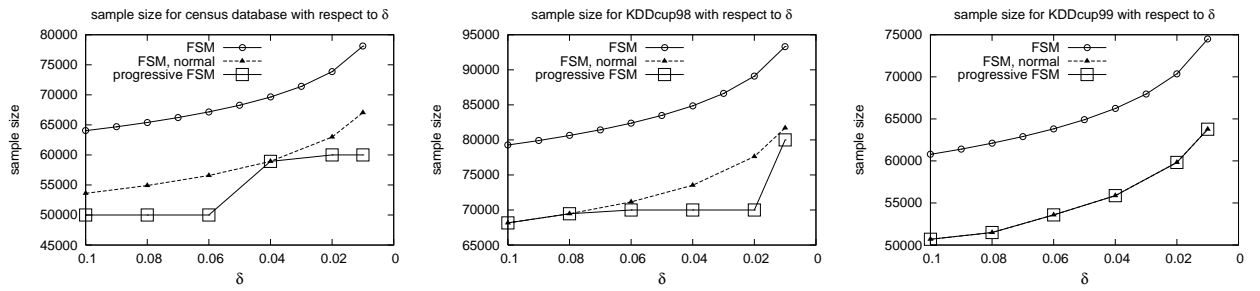


Figure 2. Sample size against  $\delta$  for schema matching algorithms.

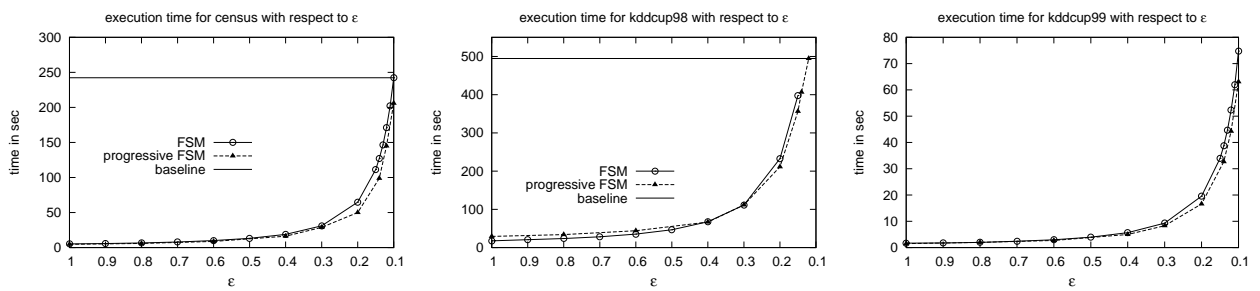


Figure 3. Execution time against  $\epsilon$  for schema matching algorithms.

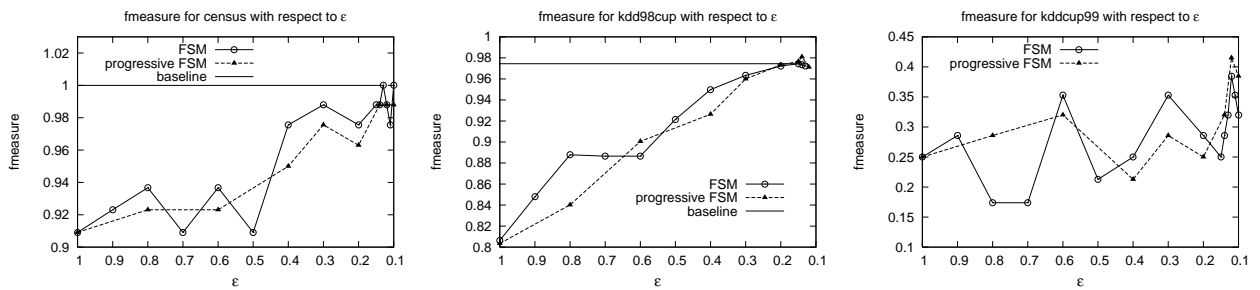


Figure 4. F-measure against  $\epsilon$  for schema matching algorithms.

(Step 4) of progressive FSM occasionally kicks in, reducing the number of samples on average compared to FSM. For the KDD Cup 1999 database, early stopping in Step 4 is never exercised. The reason is that due to the sparseness of the data the matching problem is very hard. Many attributes assume zeros most of the time which makes many potential matching partners appear equally good. For all problems, we observe that the normal bounds reduce the required sample size over the Hoeffding bounds.

Figure 2 compares the sample size required by FSM and progressive FSM for varying values of  $\delta$  (using  $\varepsilon = 0.15$ ). For small values of  $\delta$ , we can observe the benefit of the progressive version of the algorithm more clearly. In particular, while the sample-independent bounds of FSM and FSM with normal bounds follow a smooth curve, we observe that the point of termination of progressive FSM depends on the data and progresses in steps of 10,000; progressive FSM can terminate significantly earlier, depending on the data and on the value chosen for  $\delta$ .

Figure 3 compares the execution time of FSM, progressive FSM, and the baseline algorithm that executes a pass over the entire database for varying values of  $\varepsilon$  (using  $\delta = 0.1$ ). Again, progressive FSM is the fastest method, followed by FSM. The baseline algorithm that calculates the similarities based on the entire database is substantially slower, depending on the database size. For the KDD Cup 1999 database, we are unable to obtain results for the baseline algorithm which here exceeds our patience.

Finally, Figure 4 compares the F-measure of FSM, progressive FSM and the baseline algorithm that processes the entire database. For the census and KDD Cup 1998 databases, the F-measures are above 0.9 and 0.8, respectively, even for  $\varepsilon = 1$ ! For decreasing values of  $\varepsilon$ , the F-measure grows further. The KDD Cup 1999 database, by contrast, is very difficult. Since many attributes almost always assume a value of zero, most similarity values are very close and it is difficult to identify the best match; the F-measure is lower, accordingly.

Note that the low F-measure for the KDD Cup 1999 database is not in contradiction to the guarantee provided by FSM: The F-measure quantifies how frequently attributes are matched to their actually semantically equivalent partner (in our experimental setting, to themselves). By contrast, Theorem 1 guarantees that the retrieved match for each attribute is “nearly as good as the optimal match”. The baseline algorithm processes the entire database in the first place and therefore obtains a constantly high F-measure for the census and KDD Cup 1998 databases. For the KDD Cup 1999 problem, the baseline algorithm exceeds reasonable execution time, no result can be provided. For the FSM and progressive FSM algorithms, execution time and sample size as well as the quality of the solution depend on the parameters  $\varepsilon$  and  $\delta$  rather than on the size of the database.

## 6.2. Weather Data Case Study

In this set of experiments, we explore the behavior of the progressive FSM algorithm on schemas of distinct weather databases that share only a common subset of attributes. Semantically identical attributes of the two distinct databases vary greatly, and the data contain many missing values which makes this problem a challenging benchmark.

We choose a snapshot of 275,000 records in 2005 from the University of Washington and Arlington Airport weather streams. Excerpts are visualized in Figure 5. We use the progressive FSM with normal bounds. For  $\delta = .1$  and  $\varepsilon = 0.2$ , progressive FSM terminates after seeing  $N = 20,000$  records; for  $\delta = 0.1$  and  $\varepsilon = 0.15$  after 40,000 transactions. The retrieved match does not improve after processing the entire 275,000 records. Therefore, the early stopping does not come at the price of a poorer match.

Figure 5 shows the attributes of both schemas sorted according to the ordering determined by progressive FSM; *i.e.*, the algorithm considers the leftmost attributes most likely (the rightmost attributes least likely) to have corresponding attributes in the other database. The arrows indicate the single best match for each attribute; to avoid cluttering the diagram, we draw only arrows between attributes with a particularly high similarity value ( $\geq 2$ ). The rightmost attributes therefore have no displayed arrows.

The algorithm finds exact matches for time, date, Julian date, and  $T_{air}$ . These attributes can be matched based on character bigrams that occur in their values. Attribute  $T_{dew}$  is erroneously matched to  $T_{air}$ . Attribute Cht (all missing values) is matched to avg (also all missing values). The Arlington stream contains attributes with almost only missing values that are not present in the Washington stream; they are matched to vis (also almost only missing values).

It is inevitable that an instance-based schema matcher will map attributes with all missing values to other attributes with all missing values, when no additional schema-level clues are available. For  $T_{dew}$ , the diverging ranges of values prevent the involved similarity metrics from identifying these attributes as equivalent.

These results underscore that (a) the sampling-based approximation and stopping criterion are feasible, scalable, and do not deteriorate the quality of the results compared to processing all data. (b) The results emphasize that instance-information can in many cases be used to match equivalent attributes; however, the similarity metric has to be designed to be appropriate for the problem at hand, and attributes with all missing values and attributes whose values diverge too strongly between databases limit the applicability of instance-level metrics.



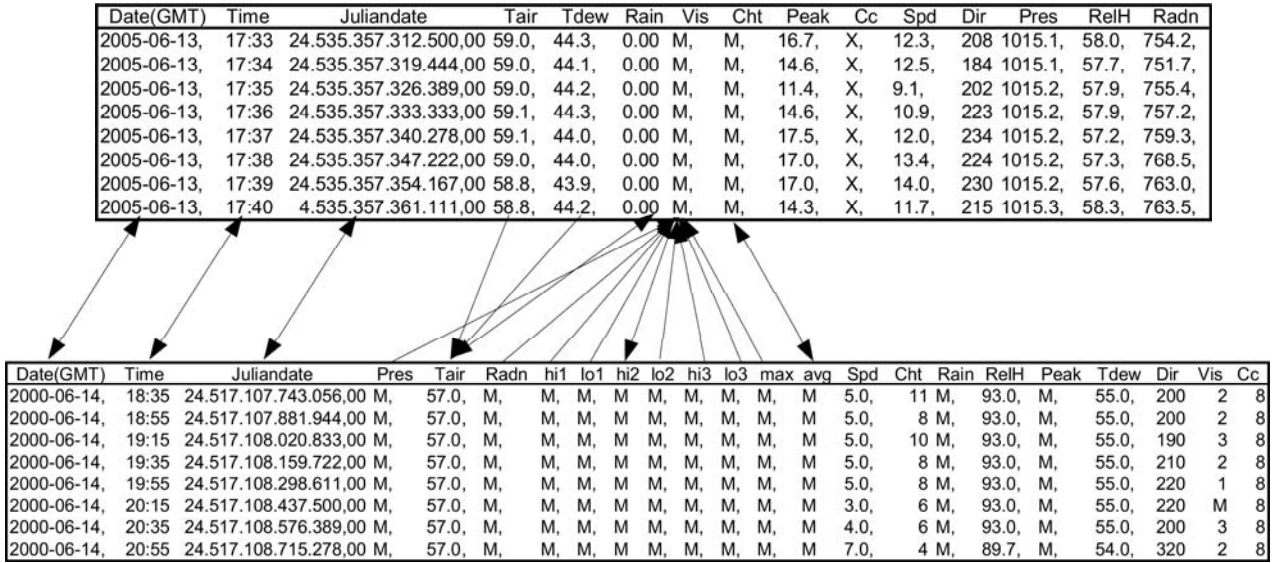


Figure 5. Uni Washington (top) matched against Arlington Airport (bottom) weather stream.

## 7. Conclusion

Our formulation of the approximately optimal schema matching problem is closely tied to practical applications and at the same time mathematically rigorous. The FSM and progressive FSM algorithm provably solve this problem; that is, for each attribute of either schema, they find the  $k$  approximately best matching partners, and approximately order them according to their similarity. Finally, the attributes in either schema are approximately sorted according to their likelihood of having a partner in the other schema. The term “approximately” here means “up to user-specified approximation parameters  $\epsilon$  and  $\delta$ ”.

The required sample size and execution time of the algorithms depends on these parameters, but are independent of the amount of data available; therefore, the algorithms apply to infinite streams. For the progressive FSM algorithm, sample size and execution time depend on the stream at hand. If there are some similar but many dissimilar potential matches, then progressive FSM can identify the similar matches faster and terminate early.

Our experiments lead to a number of conclusions. (a) FSM and progressive FSM are feasible and practically applicable for streams and very large databases. They can be applied to a range of similarity metrics; we specified bounds for a selection of measures that can easily be extended. (b) FSM with normal bounds and progressive FSM are equally fast for difficult matching problems with many very similar attributes. Progressive FSM is faster than FSM otherwise.

The normal bounds outperform the Hoeffding bounds. (c) While theoretical results guarantee an approximately optimal match, we observe empirically that attributes are often actually matched to semantically equivalent attributes. For the census and KDD Cup 1998 databases, high F-measures are observed even for large values of  $\epsilon$ .

## Acknowledgment

Tobias Scheffer is supported by the German Science Foundation DFG under Grant SCHE540/10-2.

## References

- [1] J. Berlin and M. Motro. Autoplex: automated discovery of content for virtual databases. In *Proceedings of the International Conference on Cooperative Information Systems*, 2001.
- [2] P. Bernstein and E. Rahm. Data warehouse scenarios for model management. In *Proceedings of the International Conference on Entity Relationship Modeling*, 2000.
- [3] R. Dhamankar, Y. Lee, A. Doan, A. Halevy, and P. Domingos. imap: Discovering complex semantic matches between database schemas. In *Proceedings of the ACM SIGMOD Conference*, 2004.
- [4] H.-H. Do and E. Rahm. Coma - a system for flexible combination of schema matching approaches. In *Proceedings of the International Conference on Very Large Databases*, 2002.

- [5] A. Doan, P. Domingos, and A. Halevy. Reconciling schemas of disparate data sources: a machine-learning approach. In *Proceedings of the ACM SIGMOD Conference*, 2001.
- [6] A. Doan, P. Domingos, and A. Levy. Learning source descriptions for data integration. In *Proceedings of the WebDB Workshop*, 2000.
- [7] C. Domingo, R. Gavaldà, and Osamu Watanabe. Adaptive sampling methods for scaling up knowledge discovery algorithms. *Data Mining and Knowledge Discovery*, 6(2):131–152, 2002.
- [8] S. Jaroszewicz and T. Scheffer. Fast discovery of unexpected patterns in data, relative to a bayesian network. In *Proceedings of the ACM SIGKDD Conference*, 2005.
- [9] J. Kang and J. Naughton. On schema matching with opaque column names and data values. In *Proceedings of the ACM SIGMOD Conference*, 2003.
- [10] W. Li and C. Clifton. Semantic integration in heterogeneous databases using neural networks. In *Proceedings of the International Conference on Very Large Databases*, 1994.
- [11] W. Li and C. Clifton. Semint: a tool for identifying attribute correspondences in heterogeneous databases using neural network. *Data and Knowledge Engineering*, 33(1):49–84, 2000.
- [12] J. Madhavan, P. Bernstein, and E. Rahm. Generic schema matching with Cupid. In *Proceedings of the International Conference on Very Large Databases*, 2001.
- [13] R. Miller, Y. Ioannides, and R. Ramakrishnan. Schema equivalence in heterogeneous systems: bridging theory and practice. *Information Systems*, 19(1):3–31, 1994.
- [14] F. Naumann, C. Ho, X. Tian, L. Haas, and N. Megiddo. Attribute classification using feature analysis. In *Proceedings of the International Conference on Data Engineering*, page 271, 2002.
- [15] L. Palopoli, D. Sacca, and D. Ursino. An automatic technique for detecting type conflicts in database schemas. In *Proceedings of the International Conference On Information and Knowledge Management*, 1998.
- [16] E. Rahm and P. Bernstein. A survey of approaches to automatic schema mapping. *VLDB Journal*, 10:334–350, 2001.
- [17] M. Sayyadian, Y. Lee, A. Doan, and Arnon Rosenthal. Tuning schema matching software using synthetic scenarios. In *Proceedings of the VLDB Conference*, 2005.
- [18] T. Scheffer and S. Wrobel. Finding the most interesting patterns in a database quickly by using sequential sampling. *Journal of Machine Learning Research*, 3:833–862, 2002.
- [19] L. Yan, R. Miller, L. Haas, and R. Fagin. Data driven understanding and refinement of schema mappings. In *Proceedings of the ACM SIGMOD Conference*, 2001.

## A. Proof of Theorem 1

We will prove Theorem 1 separately for FSM and progressive FSM.

### A.1. Theorem 1 for FSM

From Definition 3 (definition of a similarity confidence interval) it follows that, after Step 1 has been executed, with probability  $1 - \delta$  there is no pair of attributes  $r_i$  in  $R$  and  $s_j$  in  $S$  such that  $\lfloor f_{R',S'}(r_i, s_j) \rfloor \delta \leq f_{R,S}(r_i, s_j) \leq \lceil f_{R',S'}(r_i, s_j) \rceil \delta$ . We will now show that the nonexistence of such an offending pair  $(r_i, s_j)$  implies that the returned solution satisfies the three conditions of Definition 2. Since the risk of existence of an offending pair is bounded by  $\delta$ , this suffices to prove the Theorem.

Before proving the three properties that characterize an approximately optimal match, let us show a useful property of the algorithm. Let  $(r, s)$  and  $(r', s')$  be two pairs of attributes whose similarity confidence bounds have been determined in Step 3 of the algorithm. Then,

$$\hat{f}_{R',S'}(r, s) \geq \hat{f}_{R',S'}(r', s') \Rightarrow f_{R,S}(r, s) \geq f_{R,S}(r', s') - \varepsilon \quad (14)$$

Implication 14 is proven as follows. Equation 15 follows from the definition of  $\hat{f}_{R',S'}$  in Step 3 of the algorithm. This implies Equation 16 because, in Step 1,  $N$  is chosen such that for all pairs  $\lceil f_{R',S'}(\cdot, \cdot) \rceil - \lfloor f_{R',S'}(\cdot, \cdot) \rfloor \leq \varepsilon$ . In addition, from Definition 3 and the first sentence of this proof we know that, for any pair,  $\lceil f_{R',S'}(\cdot, \cdot) \rceil \geq f_{R,S}(\cdot, \cdot) \geq \lfloor f_{R',S'}(\cdot, \cdot) \rfloor$ .

$$\begin{aligned} \hat{f}_{R',S'}(r, s) &\geq \hat{f}_{R',S'}(r', s') \\ &\Leftrightarrow \frac{1}{2}(\lceil f_{R',S'}(r, s) \rceil + \lfloor f_{R',S'}(r, s) \rfloor) \\ &\geq \frac{1}{2}(\lceil f_{R',S'}(r', s') \rceil + \lfloor f_{R',S'}(r', s') \rfloor) \quad (15) \\ &\Rightarrow f_{R,S}(r, s) + \frac{\varepsilon}{2} \geq f_{R,S}(r', s') - \frac{\varepsilon}{2} \quad (16) \\ &\Leftrightarrow f_{R,S}(r, s) \geq f_{R,S}(r', s') - \varepsilon \quad (17) \end{aligned}$$

1. In step 4, the algorithm selects  $H_{r_i}^*$  such that  $\hat{f}_{R',S'}(r_i, s_j)$  is maximized. That is, for every element  $s_j$  inside  $H_{r_i}^*$ , including  $s_{i_{min}} = \operatorname{argmin}_{s \in H_{r_i}^*} \hat{f}_{R',S'}(r_i, s)$  and every element  $s'$  outside  $H_{r_i}^*$ , we have that  $\hat{f}_{R',S'}(r_i, s_{i_{min}}) \geq \hat{f}_{R',S'}(r_i, s')$ . Then Equation 14 implies  $f_{R,S}(r_i, s_{i_{min}}) \geq f_{R,S}(r_i, s') - \varepsilon$ , and so the first condition of an approximately optimal match is satisfied.
2. For all attributes  $r_i$ , the retrieved sequence  $(s_{i_1}, \dots, s_{i_k})$  is sorted according to  $\hat{f}_{R',S'}(r_i, s_{i'})$ . Therefore, if some element  $s_{i'}$  occurs prior to  $s_{i''}$  it must be true that  $\hat{f}_{R',S'}(r_i, s_{i'}) \geq \hat{f}_{R',S'}(r_i, s_{i''})$ . Equation 14 then implies that  $f_{R,S}(r_i, s_{i'}) \geq f_{R,S}(r_i, s_{i''}) - \varepsilon$  which satisfies the definition of an  $\varepsilon$ -correct ordering (Definition 1). The same argument applies to the ordering of  $(r_{j_1}, \dots, r_{j_k})$ .

3. In Step 5, the permutation  $\pi_r$  sorts the attributes in  $R$  according to  $\max_{s_{i'} \in H_{r_i}^*} \hat{f}_{R',S'}(r_i, s_{i'})$ . When  $\pi_r$  sorts some attribute  $r$  prior to  $r'$ , then  $\max_{s_i \in H_r^*} \hat{f}_{R',S'}(r, s_i) \geq \max_{s_{i'} \in H_{r'}^*} \hat{f}_{R',S'}(r', s_{i'})$ . Equation 14 now implies  $\max_{s_{i'} \in H_{r'}^*} f_{R,S}(r, s_{i'}) \geq \max_{s_{i'} \in H_{r'}^*} f_{R,S}(r', s_{i'}) - \varepsilon$ ; that is,  $(\max_{j'} f_{R,S}(r_{\pi^r(1)}, s_{j'}), \dots, \max_{j'} f_{R,S}(r_{\pi^r(n)}, s_{j'}))$  is an  $\varepsilon$ -correct ordering. The same argument applies to the ordering  $(\max_{i'} f_{R,S}(r_{i'}, s_{\pi^s(1)}), \dots, \max_{i'} f_{R,S}(r_{i'}, s_{\pi^s(m)}))$ .

This proves Theorem 1.  $\blacksquare$

## A.2. Theorem 1 for Progressive FSM

The main loop of the algorithm in Table 2 can terminate for two reasons. Either  $N$  examples have been processed and from the calculation of  $N$  and the definition of the similarity confidence interval (Definition 3) we know that in this case, with probability  $1 - \frac{\delta}{2}$ , there is no pair  $(r_i, s_j)$  that violates  $\lfloor f_{R',S'}(r_i, s_j) \rfloor_{\frac{\delta}{2}} \leq f_{R,S}(r_i, s_j) \leq \lceil f_{R',S'}(r_i, s_j) \rceil_{\frac{\delta}{2}}$ . Alternatively, the algorithm can terminate in step 4 when all pairs  $(r_i, s_j)$  satisfy  $\lceil f_{R',S'}(r_i, s_j) \rceil_{\delta_t} - \lfloor f_{R',S'}(r_i, s_j) \rfloor_{\delta_t} \leq \varepsilon$ . In this case, Definition 3 says that with probability  $1 - \delta_t$ , there is no pair  $(r_i, s_j)$  that violates  $\lfloor f_{R',S'}(r_i, s_j) \rfloor_{\delta_t} \leq f_{R,S}(r_i, s_j) \leq \lceil f_{R',S'}(r_i, s_j) \rceil_{\delta_t}$ . After terminating in either of these steps, the algorithm returns results analogously to the algorithm in Table 1.

Assuming that all similarities lie within their corresponding similarity confidence intervals, Equation 14 is satisfied again and the proof in Section A.1 shows that the returned solution satisfies the three conditions of an approximately optimal matching characterized in Definition 2. However, we have assumed that all the similarity values throughout the execution of the algorithm lie within their corresponding similarity confidence interval. We now have to bound the risk that during the course of execution of the algorithm there is any pair  $(r_i, s_j)$  whose true similarity lies outside its similarity confidence interval.

When  $N$  examples have been processed, then the risk that any pair  $(r_i, s_j)$  violates

$$\lfloor f_{R',S'}(r_i, s_j) \rfloor_{\frac{\delta}{2}} \leq f_{R,S}(r_i, s_j) \leq \lceil f_{R',S'}(r_i, s_j) \rceil_{\frac{\delta}{2}} \quad (18)$$

is at most  $\frac{\delta}{2}$ ; this follows immediately from Definition 2. Furthermore, during the  $t$ -th iteration of the main loop the algorithm incurs an additional risk of  $\delta_t$  that there is at least one pair whose similarity lies outside its confidence interval. We use the union bound to assess this risk as follows.

$$\begin{aligned} Pr[\text{During execution, there is a } (r_i, s_j) \text{ that violates (18)}] \\ \leq \frac{\delta}{2} + Pr[\text{Step 4, iteration } t \text{ some } (r_i, s_j) \text{ violates (18)}] \end{aligned}$$

$$\leq \frac{\delta}{2} + \sum_{t=1}^{\infty} \delta_t = \frac{\delta}{2} + \sum_{t=1}^{\infty} \frac{\delta}{2t(t+1)} \quad (19)$$

$$\leq \frac{\delta}{2} + \frac{\delta}{2} \sum_{t=1}^{\infty} \frac{1}{t(t+1)} = \delta \quad (20)$$

This completes the proof.  $\blacksquare$

## B. Proof of correctness of confidence intervals

We will now prove Theorem 2; during the course of the proof, we will also prove Lemmas 2 and 3. Recall that  $f_{R,S}$  is a weighted sum of three other measures. The following Lemma makes it possible to derive intervals for  $f_{R,S}$  from the intervals for its three constituents.

**Lemma 4** *If  $\llbracket x_i \rrbracket_{\delta}, \lceil x_i \rceil_{\delta}$  for  $i = 1, \dots, m$  are simultaneous confidence intervals for  $x_1, \dots, x_m$  respectively, with significance level  $\delta$ , and  $w_i \geq 0$  for  $i = 1, \dots, m$ , then*

$$\lfloor w_1 x_1 + \dots + w_m x_m \rfloor_{\delta} = \sum_{i=1}^m w_i \lfloor x_i \rfloor_{\delta} \quad (21)$$

$$\lceil w_1 x_1 + \dots + w_m x_m \rceil_{\delta} = \sum_{i=1}^m w_i \lceil x_i \rceil_{\delta} \quad (22)$$

is a confidence interval for  $w_1 x_1 + \dots + w_m x_m$  with significance level  $\delta$ .

**Proof.**

$$Pr \left\{ w_1 x_1 + \dots + w_m x_m \notin \left[ \sum_{i=1}^m w_i \lfloor x_i \rfloor_{\delta}, \sum_{i=1}^m w_i \lceil x_i \rceil_{\delta} \right] \right\} \quad (23)$$

$$\leq Pr \{ \exists i : w_i x_i \notin [w_i \lfloor x_i \rfloor_{\delta}, w_i \lceil x_i \rceil_{\delta}] \} \quad (24)$$

$$= Pr \{ \exists i : x_i \notin [\lfloor x_i \rfloor_{\delta}, \lceil x_i \rceil_{\delta}] \} \leq \delta. \quad (25)$$

It is easy to see from the definitions in Section 5 that we only need to prove confidence intervals for  $f_{R,S}^E$  and  $f_{R,S}^{\pi E}$ . We will therefore now prove Lemmas 1 and 2.

**Proof.** Using basic algebraic transformations it is easy to see that

$$\begin{aligned} (\forall r \in R, \forall i \in \{1 \dots n_r\} : p_{r,i} \in [\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}}, \lceil p_{r,i} \rceil_{\frac{\delta}{M}}]) \\ \text{and } (\forall s \in S, \forall j \in \{1 \dots n_s\} : q_{s,j} \in [\lfloor q_{s,j} \rfloor_{\frac{\delta}{M}}, \lceil q_{s,j} \rceil_{\frac{\delta}{M}}]) \end{aligned}$$

implies

$$\begin{aligned} \forall r \in R, \forall i \in \{1 \dots n_r\}, \forall \{s \in S : n_s = n_r\}, \\ \forall j \in \{1 \dots n_s\} : (p_{r,i} - q_{s,j})^2 \in \\ [(\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}})^2 + (\lfloor q_{s,j} \rfloor_{\frac{\delta}{M}})^2 - 2 \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \lfloor q_{s,j} \rfloor_{\frac{\delta}{M}}, \\ (\lceil p_{r,i} \rceil_{\frac{\delta}{M}})^2 + (\lceil q_{s,j} \rceil_{\frac{\delta}{M}})^2 - 2 \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \lceil q_{s,j} \rceil_{\frac{\delta}{M}}], \end{aligned}$$

and after replacing  $j$  with an  $i$ -th index of a permutation of  $\{1, \dots, n\}$ :

$$\begin{aligned} & \forall r \in R, \forall \{s \in S : n_s = n_r\}, \forall i \in \{1 \dots n_r\}, \\ & \forall \pi \in \Pi(n_r) : (p_{r,i} - q_{s,\pi(i)})^2 \in \\ & \left[ \left( \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \right)^2 + \left( \lfloor q_{s,\pi(i)} \rfloor_{\frac{\delta}{M}} \right)^2 - 2 \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \lfloor q_{s,\pi(i)} \rfloor_{\frac{\delta}{M}}, \right. \\ & \left. \left( \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \right)^2 + \left( \lceil q_{s,\pi(i)} \rceil_{\frac{\delta}{M}} \right)^2 - 2 \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \lceil q_{s,\pi(i)} \rceil_{\frac{\delta}{M}} \right]. \end{aligned}$$

After summing over all  $i = 1, \dots, n_r$  we get

$$\begin{aligned} & \forall r \in R, \forall \{s \in S : n_s = n_r\}, \forall \pi \in \Pi(n_r) : \\ & \sum_{i=1}^{n_r} (p_{r,i} - q_{s,\pi(i)})^2 \in \\ & \left[ \sum_{i=1}^{n_r} \left( \left( \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \right)^2 + \left( \lfloor q_{s,\pi(i)} \rfloor_{\frac{\delta}{M}} \right)^2 - 2 \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \lfloor q_{s,\pi(i)} \rfloor_{\frac{\delta}{M}} \right), \right. \\ & \left. \sum_{i=1}^{n_r} \left( \left( \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \right)^2 + \left( \lceil q_{s,\pi(i)} \rceil_{\frac{\delta}{M}} \right)^2 - 2 \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \lceil q_{s,\pi(i)} \rceil_{\frac{\delta}{M}} \right) \right]. \end{aligned}$$

Subtracting from 2 we get

$$\begin{aligned} & \forall r \in R, \forall \{s \in S : n_s = n_r\}, \forall \pi \in \Pi(n_r) : \\ & 2 - \sum_{i=1}^{n_r} (p_{r,i} - q_{s,\pi(i)})^2 \in \\ & \left[ 2 - \sum_{i=1}^{n_r} \left( \left( \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \right)^2 + \left( \lfloor q_{s,\pi(i)} \rfloor_{\frac{\delta}{M}} \right)^2 - 2 \lfloor p_{r,i} \rfloor_{\frac{\delta}{M}} \lfloor q_{s,\pi(i)} \rfloor_{\frac{\delta}{M}} \right), \right. \\ & \left. 2 - \sum_{i=1}^{n_r} \left( \left( \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \right)^2 + \left( \lceil q_{s,\pi(i)} \rceil_{\frac{\delta}{M}} \right)^2 - 2 \lceil p_{r,i} \rceil_{\frac{\delta}{M}} \lceil q_{s,\pi(i)} \rceil_{\frac{\delta}{M}} \right) \right]. \end{aligned}$$

After substituting confidence interval for  $f_{R,S}^E$

$$\begin{aligned} & \forall r \in R, \forall \{s \in S : n_s = n_r\}, \forall \pi \in \Pi(n_r) : \\ & f_{R,S}^E(r, \pi s) \in [\lfloor f_{R',S'}^E(r, \pi s) \rfloor_{\delta}, \lceil f_{R',S'}^E(r, \pi s) \rceil_{\delta}], \end{aligned}$$

and finally

$$\begin{aligned} & \forall r \in R, \forall \{s \in S : n_s = n_r\} : \max_{\pi \in \Pi(n_r)} f_{R,S}^E(r, \pi s) \in \\ & \left[ \max_{\pi \in \Pi(n_r)} \lfloor f_{R',S'}^E(r, \pi s) \rfloor_{\delta}, \max_{\pi \in \Pi(n_r)} \lceil f_{R',S'}^E(r, \pi s) \rceil_{\delta} \right]. \end{aligned}$$

Below ' $\Rightarrow$ ' denotes logical implication. By a similar (and simpler) argument we can derive

$$\begin{aligned} & (\forall r \in R, \forall i \in \{1 \dots n_r\} : p_{r,i} \in [\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}}, \lceil p_{r,i} \rceil_{\frac{\delta}{M}}]) \\ & \text{and } (\forall s \in S, \forall j \in \{1 \dots n_s\} : q_{s,j} \in [\lfloor q_{s,j} \rfloor_{\frac{\delta}{M}}, \lceil q_{s,j} \rceil_{\frac{\delta}{M}}]) \\ & \Rightarrow \forall r \in R, \forall \{s \in S : n_s = n_r\} : f_{R,S}^E(r, s) \in \\ & \left[ \lfloor f_{R',S'}^E(r, s) \rfloor_{\delta}, \lceil f_{R',S'}^E(r, s) \rceil_{\delta} \right]. \end{aligned}$$

Combining the two above implications and using modus tollens we get

$$\exists r \in R, \exists s \in S : n_s = n_r, f_{R',S'}^E(r, s) \notin$$

$$\begin{aligned} & [\lfloor f_{R',S'}^E(r, s) \rfloor_{\delta}, \lceil f_{R',S'}^E(r, s) \rceil_{\delta}] \\ & \vee \exists r \in R, \exists s \in S : n_s = n_r, f_{R,S}^{\pi E}(r, s) \notin \\ & \left[ \lfloor f_{R',S'}^{\pi E}(r, \pi s) \rfloor_{\delta}, \lceil f_{R',S'}^{\pi E}(r, \pi s) \rceil_{\delta} \right] \quad (26) \\ & \Rightarrow \exists r \in R, i \in \{1 \dots n_r\} : p_{r,i} \notin [\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}}, \lceil p_{r,i} \rceil_{\frac{\delta}{M}}] \\ & \vee \exists s \in S, j \in \{1 \dots n_s\} : q_{s,j} \notin [\lfloor q_{s,j} \rfloor_{\frac{\delta}{M}}, \lceil q_{s,j} \rceil_{\frac{\delta}{M}}]. \end{aligned}$$

Let  $A$  denote an event that one of the similarity measures considered is outside of its confidence interval, i.e. Equation 26. Since if for two events  $A$  and  $B$  we have: if  $A \Rightarrow B$  then  $\Pr\{A\} \leq \Pr\{B\}$ , it follows that

$$\begin{aligned} & \Pr\{A\} \\ & \leq \Pr\{\exists r \in R, i \in \{1 \dots n_r\} : p_{r,i} \notin [\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}}, \lceil p_{r,i} \rceil_{\frac{\delta}{M}}] \\ & \quad \vee \exists s \in S, j \in \{1 \dots n_s\} : q_{s,j} \notin [\lfloor q_{s,j} \rfloor_{\frac{\delta}{M}}, \lceil q_{s,j} \rceil_{\frac{\delta}{M}}]\} \\ & \leq \sum_{r \in R} \sum_{i=1}^{n_r} \Pr\{p_{r,i} \notin [\lfloor p_{r,i} \rfloor_{\frac{\delta}{M}}, \lceil p_{r,i} \rceil_{\frac{\delta}{M}}]\} \\ & \quad + \sum_{s \in S} \sum_{j=1}^{n_s} \Pr\{q_{s,j} \notin [\lfloor q_{s,j} \rfloor_{\frac{\delta}{M}}, \lceil q_{s,j} \rceil_{\frac{\delta}{M}}]\} \\ & \leq \sum_{r \in R} \sum_{i=1}^{n_r} \frac{\delta}{M} + \sum_{s \in S} \sum_{j=1}^{n_s} \frac{\delta}{M} = \delta. \end{aligned}$$

Let  $P = (p_1, \dots, p_n)$  be a probability distribution. We will now give confidence intervals for a single probability value  $p_i$ .

**Proof of Theorem 2.** Lemma 1 and the introductory notes of Section 5 prove bounds for  $f_{R,S}^{\pi E}$  and  $f_{R,S}^{Bg}$ . For the case of  $f_{R,S}^c(r, s)$ , note that the elements of vectors  $P^r$  and  $P^s$  are averages of numbers from the interval  $[0,1]$ , so they are, in fact, random variables on that interval, and Hoeffding and normal bounds (through Central Limit Theorem) apply to them. Now, Lemma 1 applies giving confidence intervals for  $f_{R,S}^c$ . Applying Lemma 4 to confidence intervals for  $f_{R,S}^{\pi E}$ ,  $f_{R,S}^{Bg}$  and  $f_{R,S}^c$  completes the proof.  $\blacksquare$