Rule Sharing for Fraud Detection via Adaptation

Ariel Jarovsky¹ Tova Milo¹

Slava Novgorodov¹

¹*Tel Aviv University* ²*Recruit Institute of Technology* ¹{arielhe, milo, slavanov}@post.tau.ac.il ²wangchiew@recruit.ai

Abstract—Writing rules to capture precisely fraudulent transactions is a challenging task where domain experts spend significant effort and time. A key observation is that much of this difficulty originates from the fact that such experts typically work as "lone rangers" or in isolated groups, or work on detecting frauds in one context in isolation from frauds that occur in another context. However, in practice there is a lot of commonality in what different experts are trying to achieve. In this paper, we present the GOLDRUSH system, which facilitates knowledge sharing via effective adaptation of fraud detection rules from one context to another. GOLDRUSH abstracts the possible semantic interpretations of each of the conditions in the rules at the source context and adapts them to the target context. Efficient algorithms are used to identify the most effective rule adaptations w.r.t a given cost-benefit metric. Our extensive set of experiments, based on real-world financial datasets, demonstrate the efficiency and effectiveness of our solution, both in terms of the accuracy of the fraud detection and the actual money saved.

I. INTRODUCTION

Financial frauds are unauthorized financial operations (called transactions) that obtain money or goods illegally from accounts. Financial frauds are a billion dollar industry and financial companies (banks, credit card companies, etc.) invest significant resources to detect frauds and prevent them. Online fraud detection systems monitor incoming transactions and use models based on data mining and machine learning (ML) techniques to detect frauds [15]. A typical approach is to score each transaction and every transaction whose score is above a threshold is classified as fraudulent.

However, such approaches may still not achieve high precision and recall as legitimate transactions may be misclassified as fraudulent and, likewise, fraudulent transactions may be missed. Also, the derived threshold does not provide a semantic explanation of the underlying causes of the frauds like the ways rules do. For this reason, financial companies typically rely, in addition to ML models, on rules that are carefully specified by domain experts. The rules refine the ML scores with domain specific conditions, for automatically determining fraudulent transactions.

Writing rules to capture precisely fraudulent transactions is a challenging task where domain experts spend significant effort and time. A key observation is that much of this difficulty is due to the fact that experts typically work as "lone rangers" or in isolated groups to define the rules, or work on detecting frauds in one context in isolation from frauds in another context. However, in practice there is a lot of commonality in what different experts are trying to achieve. Very often, rules defined by experts in one given context may be useful for understanding how to detect frauds in another context. Such collaboration between experts is desired between the different branches of financial entities, and even between competitors that are willing to share knowledge for the greater cause of battling financial crime (as imposed by local regulations [24]). The goal of the GOLDRUSH system presented here is to facilitate knowledge sharing via effective adaptation of fraud detection rules from one context to another.

Wang-Chiew Tan²

We examine the conditions of the rules, and the source and target contexts, to map rules from one context to another. The possible semantic interpretations of each condition are abstracted and then instantiated to the target context. As the number of possible abstractions (and corresponding instantiations) may be large, each rule can be mapped to the target context in many different ways. The efficient algorithms underlying GOLDRUSH identify the most effective ones, in terms of improving the fraud detection accuracy in the target context. To illustrate, let us consider the following example.

Example 1.1: Consider two fraud detection experts, A and B, working in two different collaborating companies, located in the USA and Germany, respectively. Figures 1 and 2 show a portion of each company's transaction relation (for simplicity, we present here just some real-world attributes, including the transaction time, amount, type, and country). We also present in the last column the ground truth labeling for each transaction: fraud (marked as F) or legitimate (L).

$$\varphi^{A}$$
:Type = "Stock Trade" \land Amount $\geq 100K \land$
Time $\geq 16:00 \land$ Country $\in \{Dinotopia, Jamonia\}$

The above rule is designed by expert A to detect the frauds in Figure 1. In practice each rule also includes a threshold (not shown) on the transaction score as derived by a Machine Learning module, i.e., the degree of confidence that the transaction is fraudulent, as well as additional conditions on the user/settings/etc. We will omit the scores and the additional conditions for simplicity and focus only on the rules in this example.

The rule defined by expert A in her context may be useful, after some appropriate adaptation, for the context of expert B, and different interpretations will yield different target conditions. Observe that the condition over Time may refer to the local time or the time after the closing of the local Stock Exchange Market (i.e., 16:00 for the New York Stock Exchange). A time-zone adaptation should be applied for the former interpretation (+6 hours) to adapt it to the time in Germany. However, the target condition should be Time $\geq 20:00$

		•	
107K	Stock Trade	Dinotopia	L
104K	Stock Trade	Dinotopia	F
111K	Stock Trade	Jamonia	F
102K	Stock Trade	Dinotopia	F
96K	Stock Trade	Dinotopia	L
:	:	:	:
	107K 104K 111K 102K 96K :	107K Stock Trade 104K Stock Trade 111K Stock Trade 102K Stock Trade 96K Stock Trade : :	107K Stock Trade Dinotopia 104K Stock Trade Dinotopia 111K Stock Trade Jamonia 102K Stock Trade Dinotopia 96K Stock Trade Dinotopia : : :

Fig. 1: Expert A Transactions

(the Frankfurt Stock Exchange closing time) if it is the latter interpretation. Similarly, Country $\in \{Dinotopia, Jamonia\}$ may refer to all over the world attacks originating from these two specific (fictional) countries, in which case an identity mapping should be employed. Conversely, if the condition deals with a specific attack against the local market (USA in this case), we should map it by the rule target country market attackers, e.g. Orsinia. Finally, Amt $\geq 100K$ may be a condition in terms of the local currency, in which case a translation from US dollars to Euro may be applied (translates to about 95K). Or, it may be a condition that captures "exceptionally large amounts", which, considering the trade amounts distribution in context B, should be mapped to 200K Euro. A resulting possible translation for rule φ^A then may be the following:

$$\varphi^B$$
: Type = "Stock Trade" \land Amount $\ge 95K \land$
Time $\ge 20: 00 \land$ Country $\in \{Orsinia\}$

Our goal is to facilitate meaningful and effective rules mappings between different contexts. A fundamental challenge is that rule semantics is often undocumented. To overcome this, we focus on the individual rule conditions and derive a set of candidate value abstraction and concretization functions that may capture the possible mappings between the conditions in the given contexts. These include common built-in mappings (e.g. currency and distribution-based mappings) as well as semantic-aware mappings, derived by analyzing the given data instance. To choose between the possible mappings we define an intuitive cost-benefit model that reflects the improvement in fraud detection that the resulting rules may bring. Finally we provide a set of efficient algorithms to choose the best translation among the candidates set. We show that finding the optimal rule adaptation is NP-hard (intuitively, as all the combinations between the conditions adaptation candidate values need to be considered), but the problem can be modeled as an Integer Linear Programming (ILP) problem. As the performance of ILP solvers often deteriorates for large number of variables (which is the case here), we employ a dedicated data reduction technique that clusters together transactions that "behave uniformly" w.r.t the conditions in the rule, dramatically reducing the number of variables and yielding efficient performance¹.

Our contributions can be summarized as follows.

• We formulate and present a novel rule adaptation framework that employs value abstraction and concretization to map individual rule conditions from one context to

¹While the ILP solver still goes through exponentially many combinations, it gives very good performance in practice on the reduced data.

Time	Amount	Туре	Country		
19:53	140K	Stock Trade	Orsinia	L	
20:02	97K	Stock Trade	Orsinia	F	
20:03	230K	Stock Trade	Orsinia	F	
20:05	92K	Stock Trade	Orsinia	L	
20:07	206K	Stock Trade	Orsinia	F	
:	:	:	:	:	

Fig. 2: Expert B Transactions

another. An intuitive cost-benefit model measures the improvement in fraud detection that the resulting rules may bring.

- We show that the problem of identifying the best rule adaptation is NP-hard and propose an effective data reduction technique which, together with a dedicated ILP formulation of the reduced problem, yields nevertheless a practically efficient performance.
- We have implemented our solution in the GOL-DRUSH prototype system and applied it on real use cases, demonstrating the effectiveness and efficiency of our approach as a whole, as well as the individual contribution of each components. We performed experimental evaluations on five real-life datasets of financial transactions that showed how our algorithms consistently outperformed alternative baseline solutions, both in terms of the accuracy of the frauds detection and the actual money saved.

We note that while some previous work has studied the problem of adapting rules to changes in the data (e.g. [25], [11], [16]), they mainly focus on incremental maintenance of rules in a *given context* as new transactions come in. In contrast our work examines rule adaptations entailed by context switching. Closest to our work is research on transfer learning, that attempt to reuse knowledge gained while training a model in one context to save training cost in another context. In particular fraud detection rules may be viewed as a form of decision trees to be transferred between contexts. However, since their objective function is different, as we show in our experiments, not only the resulting decision trees after the ML-based transfer are less meaningful in terms of (semantics) understability to the human expert, they are also less effective in terms of fraud detection.

While we focus here on a particular application domain, our techniques are generic and may be generally utilized to adapt classification rules under similar cost-benefit objectives.

The paper is organized as follows. Section II presents the data and cost-benefit model that we employ and formalizes the rule adaptation problem. The problem complexity and the optimized algorithms to solve it are considered in Section III. The classes of value mappings that we consider are detailed in Section IV. The experiments are described in Section V. We present related work in Section VI and conclude in VII.

II. PRELIMINARIES

In the following section we present all the necessary definitions to formally define the problem. We start by explaining the data model (transactions and rules) followed by the definition of rule adaptation and the cost/benefit model that we employ.

A. Transactions

Our model consists of a transaction relation $T(A_1, A_2, \ldots, A_m) = \{t_1, \ldots, t_n\}$, which is a set of tuples (or *transactions*) over the attributes A_1, A_2, \ldots, A_m . In a financial context, a tuple t_i captures an operation (i.e. transfer, payment) made through some bank or credit card. The transaction relation is appended with more transactions over time. For brevity, when the attributes and the transactions set are clear, we will omit them and use just T. Every attribute A_i belongs to a domain $dom(A_i)$ which has a partial order that is reflexive, antisymmetric and transitive, with a greatest element \top_A and a least element \perp_A . Attributes associated with a partial order but not with a total order are called categorical attributes. The elements in such partial order are sometimes referred to as concepts. W.l.o.g. we also assume that \top_A and \perp_A do not appear in any of the tuples ². These special elements will be useful when defining rule mappings. For brevity, when an attribute name is clear from the context we omit it and simply use the notations \top and \perp .

A transaction may be classified *fraudulent* which means that the transaction was carried out illegally. Conversely, a transaction may be verified to be *legitimate* and classified accordingly. The labeling is assumed to correspond to the (known part of the) ground truth. Given a transaction relation T we denote the set of fraudulent (legitimate resp.) transactions as fraud(T) (legit(T)). In addition, each transaction is attached with a score between 0 and 1, (computed automatically using machine learning techniques) which represents the estimated probability of each transaction to be fraud. Figures 1 and 2 depict two sample transaction relations. For the simplicity of the presentation, the score is omitted from the figures. Also, we assume (for simplicity) that all contexts have the same schema. Otherwise one may apply schema matching techniques before employing our algorithms.

B. Rules

Fraud detection experts specify rules that can be applied to a transaction relation to identify the fraudulent tuples. For efficiency of execution, the rules are typically written over a single relation, which is a universal transaction relation that includes all necessary attributes (possibly aggregated or derived from many other database relations) for fraud detection. Hence, it is not necessary to consider explicit joins over different relations in the rule language. A *rule* φ in our setting is defined as a conjunction of one or more conditions over the attributes of a transaction relation T, For simplicity, each rule includes only one condition over each attribute, but multiple disjunctive conditions over the same attribute can be expressed using multiple rules. Note that the rule language that we consider, albeit simple, forms the core of common rule languages used by actual systems and real-world industry companies [16]. Indeed, as shown in the experiments, it covered all the obtained rules.

More formally, a rule is of the form $\varphi = \bigwedge_{1 \le i \le m} \alpha_i$ where α_i is a condition of the form ' $A_i \ op_i \ v_i$ ', $op_i \in \{=, \neq, <, >, \le, \le, \in, \notin\}$. We assume that the rules are well-defined: the v_i values belong to the corresponding attribute's domain and the operators' semantics is the corresponding from the domain. For readability, in our examples we show only the non-trivial conditions on attributes, namely omit conditions of the form $A_i \le \top$. The semantics of a rule φ is defined as a boolean function over the tuples of a relation T such that for a tuple $t = < A_1 : a_1, A_2 : a_2, \ldots, A_m : a_m >$,

$$\varphi(t) = \bigwedge_{1 \le i \le m} \alpha_i [a_i / v_i]$$

Here a_i/v_i denotes the replacement of v_i by a_i in α_i , and the semantics of $\alpha_i[a_i/v_i]$ is defined by the standard op_i semantics, except for the case where $a_i = \top$ for which α_i is always satisfied. Similarly, the semantics of applying a rule φ on T denote the set of all tuples that are *captured* by φ , i.e.: $\varphi(T) = \{ t \in T \mid \varphi(t) = True \}$. Let Φ denote a set of rules over T, then $\Phi(T) = \bigcup_{\varphi \in \Phi} \varphi(T)$. In other words, $\Phi(T)$ denotes the result of the union of evaluating every rule in Φ over T. Observe that for every T and for every $\varphi \in \Phi$, $\varphi(T) \subseteq \Phi(T) \subseteq T$, since every rule selects a subset of transactions from T.

As an example, the rules φ^A and φ^B from the Introduction are rules in the presented language.

C. Rule-Driven Transactions Classification

Recall that transactions may be associated with a classification label that indicates whether they are fraudulent or legitimate. Given a rule φ and a set of transactions T (typically the full transaction relation), we define the set of *fraudulent* transactions captured by φ as $F_C(\varphi, T) = fraud(T) \cap \varphi(T)$. Conversely, the set of uncaptured fraudulent transactions is defined as $F_U(\varphi, T) = fraud(T) \setminus F_C(\varphi, T)$. The set $L_C(\varphi, T)$ of legitimate transactions (wrongly) captured by φ and the set $L_U(\varphi, T)$ of uncaptured legitimate transactions are similarly defined, by replacing fraud(T) with legit(T) (and correspondingly replacing F_C with L_C in the L_U formula).

Notice that these definitions can be extended also to be used with a rule set Φ by simply replacing φ with Φ . For brevity, when the rule/rules set and the transaction relation are clear, we will omit them and simply use the notations F_C , F_U , L_C and L_U . Ideally, we aim to build a set of fraud detection rules in which F_C and L_U are as large as possible (and so F_U , L_C are as small as possible).

D. Rule Adaptation

The *context* of an expert includes her transaction relation and her current set of fraud-detection rules, as well as relevant contextual information such as the company she works for, its location, currency, language, local regulations, etc.

While there might be many ways to map rules from a source context to a target, we center our attention here on mappings

 $^{^{2}}$ If not the case, add a new special element to the domain and set it smaller/greater, in the partial order, than all other elements.

that consider the individual conditions of the rule, substituting (when needed) the values v_i used in the rule conjuncts by those that best match the target context. As we will see in Section V, such value-based mappings are extremely effective.

Which value substitutions should one consider when adapting a rule from one context to another? As illustrated in our running example from the Introduction, there are multiple semantics to consider for each condition (attribute). For instance, the value in the *Amount* attribute conjunct can represent an absolute amount in the local context currency or some regulatory value, or alternatively correspond to some percentile in the amounts distribution. The candidate mappings are thus obtained by first abstracting the value from the source context, using each of these possible relevant semantics, and then concretizing the abstracted values to the target context.

More formally, given an attribute A_i , we denote by $\Sigma(A_i)$ its given set of possible semantics. For a context s, an attribute A_i and a possible value semantics $\sigma \in \Sigma(A_i)$, the *attribute abstraction function* $\alpha_{\sigma}^s : dom(A_i) \to dom(\sigma)$ maps values from the attribute domain $dom(A_i)$ to the abstract semantic domain $dom(\sigma)$. Conversely, the *attribute concretization function*, γ_{σ}^s , is the inverse function that maps an abstract semantic value $y \in dom(\sigma)$ to its set of possible origins in $dom(A_i)$, except for the special \top element which is mapped to itself: $\gamma_{\sigma}^s(y) = \{x | \alpha_{\sigma}^s(x) = y \text{ if } y \neq \top$, else $\top\}$

To map a rule φ from a source context s to a target context t, we will consider, for each conjunct A_i op v_i in φ , its possible abstractions (under the different semantic $\Sigma(A_i)$), and then concretize each to the target domain. Formally,

Definition 2.1 (Attribute Mapping Candidates):

Given a source context s, a target context t and an attribute A_i , the mapping candidates for a value $v_i \in dom(A_i)$ are the values

$$V_i^{st}(v_i) = \bigcup_{\sigma \in \Sigma(A_i)} \gamma_{\sigma}^t(\alpha_{\sigma}^s(v_i))$$

When s, t and v_i are clear from the context, we will omit them and use just V_i .

Example 2.2: To continue with our example from the Introduction, consider the two financial institutes A and B that want to share the knowledge about fraud attacks. Let us focus on the Amount attribute of the first rule φ^A . Naturally, there are several possible mappings for this attribute. For example one of them may use the exchange rate of US Dollars (the local currency of A in USA) to Euros (the local currency of B in Germany). In this case the abstraction here maps the local amount to some agreed upon reference currency (Swiss Francs in the example); a conversion from the reference currency to Euro is then used as the concretization function (marked as CC in Figure 3). A second possible option is a distribution-based mapping that maps the value to its corresponding percentile (e.g. upper 5%) among the transaction amounts in context A. The concretization function then maps the abstract percentile to the corresponding concrete value in the target context B (marked as VP). Another possibility is using semantics that maps the amount to some local financial regulation, assuming

Semantics	Abstraction	Concretization
Identity [ID]	100K	100K
Exchange Rate [CC]	97K (CHF)	95K(EUR)
Distribution [VP]	upper 5%	200K
Regulation Limits [RL]	after hours	120K

Fig. 3: Mapping of $\{Amount \ge 100K\}$ from context A to B

that the value of the attribute in the rule is exactly the regulation limit. Here the abstraction will be the "type of regulation" (e.g. after hours trade) and the concretization is the value of the corresponding regulation limit in the target context country (marked as **RL**). In addition, there is an identity mapping that leaves the value unchanged (marked as **ID**). The summary of the different mappings discussed is depicted in Figure 3.

We are now ready to define the set of candidate rule mappings. Intuitively, it includes all the rules obtained by substituting each of the values appearing in φ by one of its corresponding possible mappings.

Definition 2.3 (Rule Mapping Candidates): Given a source context s with a rule $\varphi = \bigwedge_{1 \le i \le m} A_i \ op_i \ v_i$, and a target context t. Let $V_i \ (1 \le i \le m)$, be the set of attribute mapping candidates for v_i . The set of rule mapping candidates for φ is the set of rules

 $\Psi^{st}(\varphi) = \{\varphi[v'_1/v_1, \dots, v'_m/v_m] \mid v'_1 \in V_1, \dots, v'_m \in V_m\}$ Here again, when s and t are clear from the context we omit them and simply use $\Psi(\varphi)$.

E. Cost & Benefit model

As previously explained, to compare between the different rule candidates in $\Psi(\varphi)$ and determine which is the most suitable mapping, we need to measure the "cost & benefit" it entails. Intuitively, the gain from a new rule can be measured by the increase in the number of fraudulent transactions that are captured by adding it (i.e. the fraudulent transactions that were not captured by the existing rules), minus the number of legitimate transactions that it misclassifies (i.e. legitimate transactions that were correctly classified by previous rules).

Problem 2.4 (Best Rule Adaptation Problem): Let s be a source context with a rule φ . Let t be a target context with transaction relation T_t and an existing set of rules Φ_t . Let F_C , F_U , L_C and L_U denote the set of captured and uncaptured fraudulent and legitimate transactions in t w.r.t to T_t and Φ_t (as defined above). Let $\Psi(\varphi)$ be the rules mapping candidates set as previously defined.

The BEST RULE ADAPTATION PROBLEM is to compute a rule adaptation $\varphi' \in \Psi(\varphi)$ such that:

$$w(\varphi') = (\alpha \cdot |\varphi'(F_C)| + \beta \cdot |\varphi'(F_U)|) - (\gamma \cdot |\varphi'(L_C)| + \delta \cdot |\varphi'(L_U)|)$$
(1)

is maximized for a given $\alpha, \beta, \gamma, \delta \ge 0$.

The term $(\alpha \cdot |\varphi'(F_C)| + \beta \cdot |\varphi'(F_U)|)$ represents the benefit that can be obtained from adding the rule adaptation φ' to the existing rule set Ψ_t in terms of the fraudulent transactions captured, while the term $(\gamma \cdot |\varphi'(L_C)| + \delta \cdot |\varphi'(L_U)|)$ represents the rule costs (in terms of legitimate transactions captured). By weighting each of the components one can tune, depending on the application, the precision and recall of the selected rule. For instance, setting $\alpha = 1$, $\beta = 1$, $\gamma = 0$, $\delta = 0$ will maximize the recall of the new rule set, while the dual assignment ($\alpha = 0$, $\beta = 0$, $\gamma = 1$, $\delta = 1$) will maximize its precision. A more balanced function (which represents the common policy of financial companies) can be $\alpha = 0.5$, $\beta =$ 1, $\gamma = 0.5$, $\delta = 1$. This function aims to capture as much uncaptured fraudulent transactions as possible (especially new ones), while at the same time maintaining a low number of (new) false positives.

Additional constraints (e.g. thresholds on the precision/recall/number of captured legitimate transactions) may be added to the problem definition, to fine tune the adaptation. For simplicity we omit this here but our algorithms extend to such more refined definition.

III. FINDING THE BEST RULE ADAPTATION

Let us first assume that we have a single rule that we wish to adapt. We will discuss the case where multiple rules are available afterwards. From the exposition in the previous section it follows that there are two main challenges to overcome to yield a good rule adaptation: (1) the choice of suitable abstraction/concretization functions that will allow to build an effective yet not too large set of candidate mappings for each rule attribute, and (2) the design of efficient algorithms to identify the best, cost&benefit-wise, rule adaptation. We will first discuss (2) in this section, assuming that the set of possible mappings for every attribute value is given. Then, in the next section, we will explain which mappings are used.

A naïve algorithm to identify the best rule adaptation would iterate over all possible attribute mappings, evaluate Equation 1 for each of the mappings, and choose the one with maximal value. However, this algorithm's complexity is exponential in the number of attributes and becomes prohibitively expensive when the relations contains many attributes, each with multiple possible mappings, as it is often the case in practice (see Section V). Indeed, we can show the problem to be NP-hard.

Theorem 3.1: Testing whether a rule φ has an adaptation $\varphi' \in \Psi(\varphi)$ whose score $w(\varphi')$ exceeds a given threshold θ is NP-Hard in the number of attributes in the transaction relation, even if each attribute has only two possible mappings.

The proof is presented in the technical report [13]

However, as we will show below, the problem can be modeled as an Integer Linear Programming (ILP) problem. While this is still NP-hard, ILP solvers are known to be efficient in practice especially if the number of variables is not too large. Unfortunately, in our setting, the variables in the ILP formulation correspond to the number of transactions, which are typically in the order of millions of them. Hence, we employ a dedicated preprocessing data-reduction step that clusters together transactions that behave uniformly w.r.t the conditions in the rule, representing them as a single tuple, and thereby significantly reducing the ILP problem space. Our experimental results (shown in Section V) on real-world datasets prove the efficiency of our solution. We will start our exposition by describing the direct ILP formulation of the problem, then explain how it is optimized via data reduction. Finally, we will consider the general case which, given a set of candidate rules, computes a set of top-k rule adaptations.

A. ILP Formulation of the Best Rule Adaptation Problem

Consider a rule $\varphi = \bigwedge_{1 \le i \le m} A_i \ op_i \ v_i$, and let V_i , $i = 1 \dots m$, be the set of attribute mapping candidates for v_i . Intuitively, our ILP problem consists of a boolean variable for every value in each V_i (with the value 1 symbolizing that the corresponding mapping was chosen, and 0 that it was not). We also have a boolean variable for each of the tuples in the target transaction relation T_t (with the value 1 symbolizing that it satisfies the rule under the chosen mapping). The objective function will be analogous to the one defined in problem 2.4. We also have a constraint for every attribute A_i in the rule, ensuring that a single value is selected from V_i . Formally, we define the ILP as follows:

Notations Let $\pi(\Gamma)$ be the set of indexes of transactions in T_t belonging to the set Γ for $\Gamma \in \{F_C, F_U, L_C, L_U\}$. Let $c_i = |V_i|$ be the number of different mapping candidates for attribute A_i and let $n = |T_t|$, then:

- Boolean (0-1) Variables:
 - $x_i = 1$ iff transaction t_i was captured by the chosen rule adaptation
 - $A_{ij} = 1$ iff mapping candidate *j* of the set V_i was selected for attribute A_i

$$\alpha_{ijk} = 1$$
 if attribute A_i on transaction t_k is satisfied by

mapping
$$A_{ij}$$
, otherwise 0

Model Formulation Our goal is to maximize the following objective function under the given set of constraints.

Objective Function:

$$\sum_{i \in \pi(F_C)} \alpha \circ x_i + \sum_{i \in \pi(F_U)} \beta \circ x_i - \sum_{i \in \pi(L_C)} \gamma \circ x_i - \sum_{i \in \pi(L_U)} \delta \circ x_i$$

Constraints:

First, we want to allow just one mapping to be selected for every attribute i:

$$\bigvee_{1 \le i \le m} \sum_{j=1}^{c_i} A_{ij} = 1 \tag{2}$$

Then, we have the constraints that define if transaction k is satisfied by the selected adaptation: we require for every attribute A_i that if j was the selected mapping among the candidates in V_i , then the constant α_{ijk} must be equal to 1. For this purpose, we check the transaction satisfaction by counting the number of attributes satisfied by the chosen candidates. Thus, a count result of m is equivalent to a full record satisfaction. We use a lower bound equation and an upper bound one in order to force the variables x_k to behave as expected.

Upper bound: On the first hand, we want x_k to be 1 only in case of full satisfaction. Thus, we need an expression that will allow x_k to be 1 only if there are at least m attributes satisfied.

$$\bigvee_{1 \le k \le n} m \cdot x_k \le \sum_{i=1}^m \sum_{j=1}^{c_i} \alpha_{ijk} \cdot A_{ij} \tag{3}$$

Lower bound: On the other hand, we want x_k to be 1 in every case of full satisfaction. Thus, we need an expression that will require x_k to be at least 1 if there are m attributes satisfied.

$$\bigvee_{1 \le k \le n} m - 1 + x_k \ge \sum_{i=1}^m \sum_{j=1}^{c_i} \alpha_{ijk} \cdot A_{ij}$$
(4)

This concludes the construction.

The size of the ILP problem is quadratic in the number of transactions in the target relation T_t . More precisely, it includes $3n + m + \Sigma |V_i|$ constraints over $n + \sum |V_i|$ variables, where n, m are the number of transactions and attributes in T_t , resp. As will be shown experimentally in Section V, this presents a performance bottleneck when the number of transactions in T_t is large. To overcome this, we use a data reduction technique that allows to represent a set of "indistinguishable" transactions by a single tuple.

B. Optimization via Data Reduction

Given a rule φ , its sets V_i , $i = 1 \dots m$, of attribute mapping candidates, and the target relations T_t , we build a reduced relation T'_t which we call the *RMC (Rule Mappings Clustered)* relation. Then we define a reduced ILP problem for the compressed relation.

The reduced RMC relation: The RMC relation is built by clustering together sets of transactions that would be equally labeled by the adapted rule, independently of which specific attribute mappings are chosen.

Example 3.2: To illustrate, recall our running example from the Introduction. Consider expert B's transaction relation from Figure 2, and rule φ^A that we wish to adapt to her context. Assume that expert B has no other rules at the moment and so, all the fraudulent and legitimate transactions in the figure are currently uncaptured. Suppose that the attribute mapping candidate for adapting φ^A to context B are $V_{Amount} =$ $\{95K, 100K, 120K, 200K\}$ and $V_{Time} = \{16:00, 20:00\}$. Consider the third and fifth transactions in the Figure. One can verify that, independently of the chosen mappings, both transactions will be both captured or uncaptured by the chosen rule adaptation (since both *Amounts* are above 200*K* and both *Times* are above 20:00). Our RMC construction, described below, merges such indistinguishable tuples.

Intuitively, for every attribute A_i , we partition its domain $dom(A_i)$, using the values in V_i . We choose (as will be explained below) a representative for each partition, then replace each attribute value in the transaction relation by the representative of the partition to which it belongs. Recall that conjuncts in the rule are of the form A_i op_i v_i . The domain partitioning (and chosen representatives), and consequently the value replacements, are dictated by operator op_i used for

Time	Amount	F_C	F_U	L_C	L_U
16:00	120K	0	0	0	1
20:00	1	0	0	0	1
20:00	95K	0	1	0	0
20:00	200K	0	2	0	0
Fig. 4: Expert B RMC Transactions Relation					

value comparison. We use below $h_i^{op_i}$ to denote the value replacement function for attribute A_i .

To illustrate, let us consider the operators \geq and = (the other operators work similarly). For \geq , when the domain is totally ordered, we use the values in V_i to partition the domain into disjoint intervals [v, v'), $v, v' \in V_i \cup \{-\infty, \infty\}$, and set each partition representative value as the minimal value in the partition (v). (For the case of $v = -\infty$ we use the special \perp symbol). The value replacement functions h_i^{\geq} for attribute A_i is then defined as follows:

$$h_i^{\geq}(v) = argmax_{b \in V_i} \{ v \ge b \}$$

For =, each element in V_i forms a partition of itself whereas all the others values belong to a "complement" partition whose representative is the \perp element. The value replacement functions $h_i^{=}$ for attribute A_i is then defined as follows:

$$h_i^{=}(v) = v$$
 if $v \in V_i$ otherwise \perp

Example 3.3: Following Example 3.2, the next table presents for each *Amount* value in Figure 2, its partition and representative as induced by V_{Amount} and h_{Amt}^{\geq} .

v_i	$v_i's \ partition$	$h_{\overline{A}mt}^{\geq}(v_i)$
92K	$(-\infty, 95K)$	\perp
97K	[95K, 100K)	95K
140K	[120K, 200K)	120K
206K	$[200K,\infty)$	200K
230K	$[200K,\infty)$	200K

Suppose also that $V_{Type} = \{\text{Stock Trade, Payment}\}$. Clearly, using $h_{Type}^{=}$, all the transactions in Figure 2 will remain with the *Stock Trade* type.

Finally, we cluster indistinguishable tuples and attach, to each representative tuple, counters for the four transaction classes

 (L_C, L_U, F_C, F_U) , counting the number of tuples in the cluster belonging to the corresponding class.

Example 3.4: Continuing with Example 3.2, Figure 4 shows the *RMC table* for the 5 transactions in Figure 2 (we assume that all of them are still uncaptured). We omit the *Type* and *Country* attributes since they have the same values for all the records in the example (and thus will be indistinguishable for any possible adaptation).

The reduced ILP problem: We can now define, for the reduced RMC table, a corresponding reduced ILP Model. In this model, each of the x_i variables represents a summarized transaction in the RMC relation, and the α_{ijk} constants are defined w.r.t the x_i summarized tuples. The definition of the A_{ij} variables stays the same.

We use the same set of constraints as before and adapt just the objective function. Let χ_i^{Γ} denote the number of

original transactions of class $\Gamma \in \{F_C, F_U, L_C, L_U\}$ that were summarized into the summarized tuple x_i , then the updated objective function becomes:

$$\sum_{i \in \pi(F_C)} \alpha \circ \chi_i^{F_C} \circ x_i + \sum_{i \in \pi(F_U)} \beta \circ \chi_i^{F_U} \circ x_i - \sum_{i \in \pi(L_C)} \gamma \circ \chi_i^{L_C} \circ x_i - \sum_{i \in \pi(L_U)} \delta \circ \chi_i^{L_U} \circ x_i$$

As demonstrated in the experiments, the smaller size of the transaction relation, and consequently the smaller size of the ILP model, leads to significant performance improvement.

To conclude the discussion, observe that a large number of potential mappings allows on the one hand flexibility in rule translation but on the other had makes the data reduction task harder, as it may lead to smaller clusters. As we can see in the experiments, the compression ratio of the RMC database size, compared to the original one, indeed decreases when the number of mappings grows. However, it is still highly effective even for extremely large number (10^6) of possible mappings, and the size is reduced by 2-4 orders of magnitude. This shows that in practice, in spite of the large number of mappings, many irrelevant mapping combinations are pruned, allowing the RMC database size to be kept relatively small.

C. k-Rule Adaptation

So far we have explained how the best rule adaptation, for a single rule φ is computed. When a set of *n* rules is available (from the same or different sources), the same algorithm can be applied to each of the rules φ in the set (i.e. computing a reduced RMC relation for φ and running ILP on it). The scores of the best adaptation for each of the rules are compared, and the one with the highest score is selected. If the score is positive (or above a minimal quality threshold set by the system administrator), the translation is added to the target rules set. The process can then be iterated, to choose another qualifying rule adaptation, and so on. In the full version of the paper [13] we present an optimized version of this iterative process where "unpromising" rule candidates are pruned via a threshold-based technique. We omit this here for space constraints.

IV. GENERATING MAPPING CANDIDATES

We employ in GOLDRUSH three classes of attribute mappings. The first class consists of a standard build in set of value-based mappings. The second includes distribution based mappings. The third class is data driven and employs an ontological knowledge to determine possible abstractions/concretizations. We briefly overview each class below.

A. Mapping by Value

Our set of value-based mappings includes standard value conversions for currency, temperature, time zone, length and weight metrics. The corresponding abstraction (resp. concretization) functions here map, in each context, every value from (to) its local scale to (from) a universally agreed one.

We also include in this class three useful mappings:

Identity: This is the trivial mapping that leaves the values unchanged. Here both the abstraction and the concretization functions are simply the identity function.

Any-to-Any: This mapping allows to replace a value by any other value in the domain. Such a mapping is useful for semantic-less categorical attributes such as category code names or internal user ids, where the corresponding value in the target domain is unknown at mapping time, and allows the algorithm to examine all possible instantiations. The abstraction function here maps all values to the abstract "Any" element. The inverse concretization then facilitates all possible mappings. To avoid blowup in the number of mapping options, we apply this mapping only to categorical attributes with bounded domain size (up to 100 in our implementation).

Wildcard: This mapping allows to replace a value by the special \top element, which satisfies all conditions, thereby essentially removes the corresponding conjunct from the rule. This is particularly useful when a given rule has conditions that are relevant to the source context but not to the target. For instance conditions on the state name may be relevant to US bank branches but not for the European ones).

B. Mapping by Distribution

The second class includes distribution-based mappings such as percentile and frequency (top/bottom-k). We distinguish here between numerical and categorical attributes.

Numerical attributes: For numerical attributes there is often a correlation between the value used in the condition and the attribute's value distribution in the underlying relation. For example, assuming that only 5% of the money amounts recorded in a transaction relation are above 100K, a condition of the form of $\{Amount \ge 100K\}$ may in fact mean "in the top 5% amounts". To capture this we use an abstraction function that maps each value to its corresponding percentile in the attribute values distribution, in the source relation. The inverse concretization function then uses the target's distribution to map back to the appropriate target value. (A bottom n% mapping works in a similar manner.)

In our implementation we examined two standard methods for estimating the percentiles and their values: building a *Cumulative Histogram* and using *1-D Kernel Density Estimations* (*KDE*) [21]. As both methods yielded almost the same results we show results only for Cumulative Histograms.

Categorical attributes: The top-k (resp. bottom-k) most frequent is the analogous version of the numerical top n% (bottom n%), for categorical attributes. The abstraction function here maps every attribute value to the corresponding maximal (minimal) k. The inverse concretization function then maps each k to the appropriate value in the target context.

C. Mapping by Semantics

For categorical attributes, we use ontologies to examine the semantic relationship between the values mentioned in the attribute condition (e.g. Dinotopia and Jamonia in *country* \in {*Dinotopia*, *Jamonia*}) and each of the values mentioned in the source context attributes (e.g. USA in *location* = *USA* of



Fig. 5: A sample of geopolitical ontology expert A in the Introduction). We use in GOLDRUSH (portions of) general purpose ontologies such as YAGO [22] and DBPedia [9] as well as domain specific ontologies. We assume the ontologies to be sufficiently clean and accurate, and indeed those that we use in our implementation were built and curated by the domain experts of the participating financial institutes.

The simple, yet effective, semantic inference that we employ examines the ontology graph. It identifies, for each semantically meaningful value mentioned in the source context, simple path expressions (sequence of labeled edges) connecting the value to the values mentioned in the condition. Each such path expression captures a possible semantic abstraction. To illustrate, two possible path expressions connect the USA node to the two nodes labeled Dinotopia and Jamonia: (1) a single (yellow) edge *scams*, semantically representing the fact that both countries are the source of common scams against USA, and (2) a sequence of two blue *inContinent* edges, capturing that both countries belong to a single continent as the USA. Each of the two path expressions captures a possible semantic abstraction of the attribute value, relative to the source context. Their conjunction captures a third abstraction representing the conjunction of the two semantic properties.

For each such abstraction, the inverse concretization function retrieves, in the ontology graph, values that are accessible by the corresponding path expression(s) from the analogous target context attribute value. For instance, when adapting the rule to a context where country = Germany, the concretization of the first semantic abstraction (the scams edge) includes Orsinia, whereas the concretization w.r.t the second semantic abstraction (same continent) includes Orsania and Qurac. The concretization of the third semantic abstraction is the intersection of the two sets and includes only Orsinia, the only value connected to Germany by the two path expressions.

V. IMPLEMENTATION AND EXPERIMENTS

We have implemented the algorithms described in the previous sections in the GOLDRUSH system, using Python (backend service), PHP/JavaScript (frontend) and MySQL as the database engine. The system architecture is detailed in [13]. We next describe the datasets used in our experiments, the compared algorithms, and the experimental evaluation.

A. Datasets

We run the experiments over real-world financial transaction relations obtained from our industrial collaborators (names

DS	Monthly Trx#	Fraud Ratio	# of Rules	Precision	Recall	F1 Score
Α	45K	0.02	80 (27)	0.86	0.53	0.65
В	800K	0.002	44 (9)	0.65	0.0007	0.001
С	2.5M	0.01	55 (14)	0.52	0.1	0.17
D	4.2M	0.001	110 (34)	0.42	0.1	0.16
E	250K	0.001	62 (14)	0.6	0.02	0.05
Fig. 6: Datasets used for GOLDRUSH experiments						

omitted per companies request). Due to the sensitivity of the financial information, we received a masked version of the datasets in which personal information such as user names, account IDs and IPs were omitted, and locations were reduced to include only the city name. We obtained five datasets belonging to five different financial institutes (FIs) around the world, for the second quarter of 2016. Each dataset consists of (i) a transactions relation, of two full months activity, and (ii) the full set of fraud detection rules for that period. As mentioned in Section II, the rule language that we consider, albeit simple, forms the core of common rule languages used by actual systems and indeed covered all the obtained rules. For each FI, we used the first month data as our training set and tested our generated rules against the second month.

Transactions and Rules: The transaction relations include payment and authentication activities performed by the FI clients. Each transaction is labeled as fraudulent/legitimate. The labeling is done as part of the standard operation of the FIs (using e.g. user notifications and periodic user approval/disproval of transactions and fraud notifications) and we use these labels as the ground truth. The relations have between 30-70 attributes, numerical attributes (such as the transaction amount, number of actions in the last hour, etc.) and categorical ones (location, client type, activity type, etc.). Each transaction also includes a risk-score attribute (a value between 0 - 1000), generated by the machine learning module of our industrial collaborators, that indicates the probability of a transaction to be fraudulent and is usually used as one of the conditions in the rules. The rules contain between 3-20 conditions.

Figure 6 provides some additional statistics about the different datasets used for the experiments. The reported statistics are for the training set. The test set is quite similar. The number of transactions in the relations varies from 50K to 4.2M per month, where 0.1% to 2% of them are labeled fraudulent and the rest legitimate. The number of fraud detection rules for each FI varies from 44 to 110 (ignore for now the numbers in parentheses). The rule sets have precision between 0.42 and 0.86 and for all the datasets (except A) the recall is often below 0.1. These low recall numbers are common for fraud detection, even when performed by experienced domain experts, is often an extremely challenging task. Even small percentage of improvement in fraud capturing is considered a major success and yields great financial savings.

Ontology: In the experiments we used a geopoliticalfinancial ontology that was manually built and curated by the FIs domain experts. According to our industry collaborators such a one-time effort is cost-effective, as there are many effective rules that can benefit from the derived mappings. The ontology was built using DBPedia [9] and publicly available datasets such as CIA Factbook [7] and FATF [8] which contains lists of location-based known fraud schemes (i.e. Western Union Scams [26]), known money-laundering regions, sanctioned countries, as well as geographical properties such as each country's continent and bordering countries.

B. Algorithms

Our experiments evaluate each of the algorithms presented in the previous sections, as well as the end-to-end GOL-DRUSH system. We next list the various algorithms (and competing variants) that we examine.

The GOLDRUSH system: In our end-to-end experiments we will use the full fledged GOLDRUSH . We use the performance of the given rule sets as baseline. Note that these rules reflect the execution of state-of-the-art ML to label the transactions with risk-scores and to suggest thresholds³, along with the work of professional domain experts to refine the classification, and thus represent the strongest existing competitor. We denote this baseline by ML+E. We measure the contribution of GOLDRUSH by adding our generated rules to the FI rules set, measuring the improvement in performance, in terms of precision and recall, obtained for the extended rules set, compared to the original set. To compare to a transfer learning approach, we model the rules as decision trees and apply a state-of-the-art transfer learning algorithm from [20] (TL). (The construction details are given in the full paper [13]). Here again we add the transformed rules to the original FI rules set and examine the performance. To the best of our knowledge none of the existing decision trees TL algorithms support the incorporation of predefined mappings in the process. To nevertheless examine if/how they may improve the TL results, we also run an experiment where we first applied all possible mappings to the rules, then gave the resulting set as input to the TL algorithm (\mathbf{TL}^m) .

Best Rule Adaptation: As detailed in Section III, GOL-DRUSH solves the *Best Rule Adaptation* problem by: (1) compressing the transaction relation into an *RMC Relation*, and (2) solving an ILP problem for the reduced relation. We refer below to this two steps algorithm as **GOLD**. To examine the contribution of each of these two steps, we compare *GOLD* to the following competitors:

GOLD⁻: A restricted variant of GOLD that does not build the reduced RMC and runs the ILP on the original relation.

BF: The naïve brute-force algorithm which iterates over all possible attribute mappings, computing the score of each possible combination by issuing a corresponding SQL query over the transaction relation.

BF⁺: Similar to *BF*, except that the queries are issued on the summarized *RMC* relation.

C. Experiments

To illustrate the efficiency and effectiveness of our approach, we tested GOLDRUSH and its algorithms through three sets of experiments. We first performed an end-to-end experiment

Туре	Method	k = 1	k = 5	k = 10	k = 20
	Identity	7	53	97	163
Values	Any-to-any	5	20	43	78
values	Wildcard	4	17	28	52
	Currency	0	0	1	3
Distribution	Cum. Dist.	5	33	68	161
Distribution	Top-k Frequent	6	29	53	119
Semantic	Location Onto.	2	2	8	32

Fig. 7: Mapping Methods adoption rate

to examine GOLDRUSH as a whole, demonstrating its effectiveness for improved fraud detection and actual money saving. This also allowed us to illustrate the usefulness of each mappings class. We next examined the runtime performance of our algorithm, demonstrating its superiority compared to the competing algorithms.

1) End-to-End Benchmark: As a recommendation system, we evaluated the quality of the recommended rule-sets proposed by GOLDRUSH in terms of statistical measures such as recall, precision, F1-score, and the amount of money that the user would save if she adopts the recommended rules. For each FI, the candidate rules pool consists of the rules set from all other FIs, and we generated k rule recommendations for $k \in \{1, 5, 10, 20\}$. The presented rule-set statistics refer to the rule set consisting of the given FI rule set union those recommended by GOLDRUSH. We run the experiments for each of the FI's training and test sets. Since the results in the two cases were similar, we present both just for the F1-score measure, and only those of the test set for the other measures.

We use here the balanced ILP target score function presented in Section II-E with $\alpha = \gamma = 0.5$ and $\beta = \delta = 1$. We consider only adaptations that yield a positive cost-benefit score. In some cases, this implies that less than k rules would be recommended, and this is the reason why there are experiments for which the results for a given dataset remain unchanged from some k onwards. Figure 6 reports (in the parentheses) for each FI the number of rules, out of all rules, for which at least one such adaptation was found in the experiments. A sample of the results for GOLD vs. ML+E are depicted in Figures 8a-8c, and for GOLD vs. TL in Figure 9.

Figure 8a depicts the recall for varying k. The first bar for each dataset shows the recall of the ML+E on the test data while the others present the recall obtained with GOLD when adding the 1, 5, 10 and 20 recommended rules. The recall naturally increases as more rules are added. For 10-rule adaptation, the recall for all datasets more than doubles, and for dataset B, where the ML+E rules set recall was extremely low, the recall raises to almost 90%. Further note that for all the datasets, an important improvement is shown already at the 1-rule adaptation. This is in accordance with the observation that adequate rules are extremely hard to write and thus companies often miss important fraud patterns that are hidden in the data. The results demonstrate the great benefit that even minimal knowledge sharing may bring. Even the adaptation of a single rule (the most beneficial one, as determined by our algorithm) can greatly improve performance.

Regarding the precision, it increased in all the datasets (in some cases over 50%) except dataset A (where it has decreased

³For the companies privacy, the algorithms details cannot be disclosed



by 18% for the large k values). Notice that to avoid such decrease, a precision threshold constraint can be added to the ILP model (not done here). We have further observed in the experiments that the precision has always increased in the 1-rule adaptation, then gradually decreased for higher k values (yet generally stays higher than that of ML+E).

Figures 8b and 8c present the F1-score for ML+E and GOLD for the training and test dataset, respectively. We use this measure to highlight the trade-off between the recall improvement and the precision decrease (which, as mentioned, occurred just for a single dataset). First, note that the results in the training and test sets are fairly similar (especially in terms of the delta between the original rules set and ours), which is exactly what one would like to see. One can also notice that, compared to the original rules set, the F1-score is improved in all cases, even for dataset A with a decreased precision, and thus we see here overall trends similar to Figure 8a.



Fig. 9: GOLD vs TL evaluation

To compare also to TL, Figure 9 depicts the recall and F1-scores for ML+E, GOLD and TL for dataset A (the others show similar trends). While the recall value for TL is marginally higher than GOLD for k = 10, 20, the corresponding precision was extremely bad, as demonstrated by the fact that the TL F1-scores are significantly lower for all k. Moreover, the rules generated by TL were not meaningful to the experts in terms of semantic understability, as illustrated in the full version of the paper [13]. Interestingly, for TL^m, not

only that for each input rule an exponential number of rules (all possible mapping combinations) should be examined, the results did not improve relative to the original TL.

Finally, to highlight the monetary benefit of our approach, we calculated the overall saved money by summing up the amounts of all the uncaptured fraudulent transactions (F_U) that were captured by some rule in the recommended set. The results show that the rules could have saved between 50K to 100M(!) USD, depending on the dataset (see [13] for details).

Original Rule	Translated Rule
COUNTRY = $Dinotopia \land$	$COUNTRY = Qurac \land$
$TIME > 20: 00 \land \dots$	$TIME > 21:00 \land \dots$
$COUNTRY = Orsinia \land$	COUNTRY = $Dinotopia \land$
IP_CHANGED_TODAY > 3	IP_CHANGED_TODAY > 1
ACTIVITY = $Withdraw \land$	ACTIVITY = $Transfer \land$
AMOUNT > $5K \land \dots$	AMOUNT > $100K \land \dots$
$E = 10, T_{max}$	nalatad Dalaa

Fig. 10: Translated Rules

To conclude the experiment we examined, for each attribute mapping class, how often it was effectively used to construct the rules derived in the experiment. Interestingly all classes were useful. Figure 7 shows the results, for various mapping classes, for varying k. The most dominant mapping is the identity one, reflecting the fact that in many case values are context independent. The distribution-based is also particularly effective, as both the cumulative distribution mapping (used for numerical attributes) and the *top-k frequent* mapping (used for categorical attributes) are employed. For the semanticbased mappings, location-based mapping was very useful, chosen in over 50% of the cases for adapting the locationrelated attributes (country and city). Finally, we used the Anyto-Any mapping for semantic-less categorical attributes (such as internal code names) which generally had relatively small value domains (up to 100 values). As appreciated from the statistics, this mapping was critical for the successful mapping of rules with those attributes. Table 10 provides some examples of the resulting translated rules. We only show the conditions that were modified in the translation.

2) Best Rule Adaptation: In this set of experiments we compared the run-time performance of the *GOLD* algorithm against the three competitors presented in Section V-B. We show the results for the train datasets. The results for the test datasets are similar.

Single rule: We start by considering the time it takes to identify the best adaptation for a single rule. As for a given dataset, the performance depends on the number of attribute mapping candidates, we examine the run-time for a varying number of attribute mapping combinations. In our datasets, these ranged from 10 to almost 10^7 , depending on the attributes used in the given rule and their domain. Figure 8d summarizes the results of all runs, comparing the run-time of single rule adaptation performed with GOLD (corresponding to the x-axis), to the competitors (y-axis). Each point represents one experiment (performed on a given dataset, a given rule with a given number of possible mapping combinations, and a given competitor). The point is located in the x-y coordinates representing the corresponding running time for GOLD and the given competitor (identified by the point color and shape). As expected, the results of naïve BF were significantly worse than all other algorithms and are thus omitted. The points of the GOLD- algorithm for the dataset D are also omitted, as it took more than 3 hours even for the smallest set of mappings. By the construction of the graph, points above the y = x line (in black) are the runs where our GOLD algorithm outperforms its competitor, and as observed, most of the points are there. A closer look at the few BF⁺ points below the line showed this happens only for the smallest datasets (A and E), and with a rather small number of possible mapping combinations (less than 100), due to overhead of GOLD compared to BF. Also note that in all cases GOLD terminated in less than 10 seconds, while for BF⁺ some inputs required over 100 seconds and much more the GOLD⁻.

RMC compression: To highlight the source of efficiency of GOLD, compared to GOLD-, we examine, for the different datasets, the compression ratio between the RMC relations size and that of the original relations (Figure 8e), as well as the actual RMC size (Figure 8f). In both figures we show how they vary as a function of the number of the possible value-mapping combinations. Each experiment performed over a given dataset is represented by a dot. The dots shape and color show on which dataset the experiment was run. To highlight the compression trends in Figure 8e, we added to the figure lines that show the result of linear regression performed for each dataset. As expected, the fewer the available mappings, the greater is the compression. This is because, given a smaller set of values, the domain is split to larger intervals. Consequently more data values are represented by the same constant and can be compressed together. From the shape of the linear regression lines we can observe this trend holds for all datasets. Yet note that the compression is extremely effective even for an extremely large numbers of mappings: Even for 10⁶ possible mappings, the size is reduced by over 2 orders of magnitude for dataset A, and by 3-4 orders of magnitude for the rest. Also note that the larger the dataset is, the greater is the compression ratio (for the same number of mappings). This is because more tuples are unified and can be compressed together.

We can see the actual RMC sizes in Figure 8f. Since large data sets get compressed more, the distribution of RMC sizes of all datasets is even and we get, for all datasets, RMCs of varying sizes. But in all cases they are much smaller than the original relations. As a result, even for extremely large datasets (such as dataset D with 4.2M transactions) and extremely large number of mappings (above 10^7), the data is compressed into a practical size that allows for efficient ILP solving.

Multiple rules: As mentioned, Figure 8d depicts the running times for a single rule adaptation. Given a set of candidate rules, the same process is applied to each rule in the set and the best-score rule translation is selected. In our experiments, for each FI, the candidate rules pool consists of the rule sets from all other FIs, and the running time to identify from this pool the best rule translation for FIs A, B, C, D, and E, is 0.3, 1.8, 4.4, 7.9, and 0.6 minutes, resp. (depending as expected on the size of the customer's transactions relation). Note that even for the largest dataset the time is below 8 minutes, which is very reasonable for an offline computation. More rule translations may be naturally obtained by repeating the process (see Section III-C). This is where our k-Rule threshold-based optimization comes into play, pruning unpromising rule candidates and reducing drastically the average iteration time for the FI's to 0.03, 0.4, 0.5, 1.1 and 0.15 minutes, resp. See [13] for the full details.

VI. RELATED WORK

The identification of fraudulent transactions is essentially a classification task which is a fundamental problem in machine learning (ML) and data management [17], [19]. The relationship between ML and expert-written rules has been widely discussed in the literature [5], [1], [4] for different tasks including information extraction, conflict resolution in data integration and entity linking. Other fraud detection techniques include the use of decision tree [14], or genetic programming [3] to classify transactions into suspicious and non-suspicious ones. Our work is complementary to these lines of works and the adapted rules derived by our algorithms may be used to enhance any of these techniques.

Much of the previous research on rule-based classifiers focus on how to learn rules from the training data. Some systems, such as Chimera [23] and [16], build the rules interactively by using both ML techniques and human experts that write and refine the rules. These works however do not consider the problem of rules sharing and adaptation between different contexts. The incremental maintenance of rules in response to new incoming data has also been extensively studied in the literature (e.g. [25], [11]). In contrast our work focuses on a full context switch and employs a dedicated cost/benefit model to optimize the integration of the adapted rules to the already existing rules in the target.

Our implementation uses a geopolitical-financial ontology built with the help of our industry collaborators. Financial information about countries was also used for the task of preventing fraud in some previous work, such as [12]. Our experiments indicate that such data can indeed be effective in deriving rules with good prediction quality.

Collaboration between different parties for the improvement of fraud detection systems is sometimes enforced by local regulations (e.g. [24]). Works such as [6] present methods to encourage this collaboration and are based on centralized repositories that contain lists of fraud patterns built from fraud attacks experienced by the clients. A main drawback is that the patterns are either very general (so that they can be used by broad range of clients) and thus have low precision and/or recall, or conversely, they are specific to a certain client and require context adaptation (such as the one we provide in this paper) before they can be applicable to other clients. Combining our framework with a privacy preservation mechanism which assures that the source rule, or specific parts of it, are not discoverable from the adapted rule is an intriguing direction for future work.

Transferring knowledge from one context to another has also attracted much interest of the ML community, in areas referred to as transfer learning or domain adaptation (e.g. [18], [2]). As mentioned, closest to our work is the work on transfer learning for decision trees and random forests [20]. However, as demonstrated in our experiments, as their objective function is different the process yields rules that are less effective for the given context and often not meaningful in terms of semantics understandability.

Conceptually, the work on query reformulation (e.g., [10]) is similar to our approach. While they reformulate queries written over a source schema to a target one (under a given set of tuple/equality generating dependencies), we replace the rule conditions over a single schema using source to target condition mappings. Furthermore, our focus is on finding the best translation in terms of capturing frauds and is different from the query optimization goal which aims to minimize the size of the query.

VII. CONCLUSION

In this work we introduced GOLDRUSH, a system which facilitates knowledge sharing via effective adaptation of fraud detection rules from one context to another. Our solution employs values abstraction and concretization to map individual rule conditions from one context to another. An intuitive costbenefit model measures the improvement in fraud detection that the resulting rules bring. We employ an effective data reduction technique which, together with a dedicated ILP formulation of the reduced problem, yields a practically efficient algorithm. An extensive set of experiments on real-life datasets demonstrates the effectiveness of our approach.

There are several challenging directions for future work. Specifically, we plan to extend GOLDRUSH by supporting richer rule languages, including time/event sequences, aggregates and negation. Our mappings in this work focus on individual attributes. Considering more complex mappings that involve attribute combinations is also challenging and should take into account the dependencies among rule components. Finally, applying our rule adaptation techniques to other domains such as cyber security and medical classification is another intriguing research direction.

Acknowledgements This work has been partially funded by the European Research Council under the FP7, ERC grant MoDaS, agreement 291071, and by grants from Intel, the Blavatnik Cyber Security center, the Israel Innovation Authority and the Israel Science Foundation. Work was done while Tan was at UCSC. Tan was partially supported by NSF grant IIS-1524382 at UCSC.

REFERENCES

- B. Alexe, M. Roth, and W. Tan. Preference-aware integration of temporal data. *PVLDB*, 8(4):365–376, 2014.
- [2] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning from different domains. *Machine learning*, 79(1), 2010.
- [3] P. J. Bentley, J. Kim, G.-H. Jung, and J.-U. Choi. Fuzzy darwinian detection of credit card fraud.
- [4] D. Burdick, R. Fagin, P. G. Kolaitis, L. Popa, and W. C. Tan. A declarative framework for linking entities. In *ICDT*, 2015.
- [5] L. Chiticariu, Y. Li, and F. R. Reiss. Rule-based information extraction is dead! long live rule-based information extraction systems! In *EMNLP*, 2013.
- [6] C.-C. Chiu and C.-Y. Tsai. A web services-based collaborative scheme for credit card fraud detection. In *EEE*, pages 177–181, 2004.
- [7] CIA Factbook. https://cia.gov/library/publications/the-world-factbook/.
- [8] Financial Action Task Force. http://www.fatf-gafi.org/.
- [9] DBPedia. http://dbpedia.org.
- [10] A. Deutsch, L. Popa, and V. Tannen. Query reformulation with constraints. SIGMOD Record, 35(1):65–73, 2006.
- [11] F. J. Ferrer-Troyano, J. S. Aguilar-Ruiz, and J. C. R. Santos. Data streams classification by incremental rule learning with parameterized generalization. In SAC, pages 657–661, 2006.
- [12] E. Grace, A. Rai, E. M. Redmiles, and R. Ghani. Detecting fraud, corruption, and collusion in international development contracts: The design of a proof-of-concept automated system. In *BigData*, 2016.
- [13] A. Jarovsky, T. Milo, S. Novgorodov, and W.-C. Tan. Rule sharing for fraud detection via adaptation (technical report). http://slavanov.com/research/goldrush-tr.pdf
- [14] A. I. Kokkinaki. On atypical database transactions: Identification of probable frauds using machine learning for user profiling. *Knowledge* and Data Exchange, IEEE Workshop on, 0:107, 1997.
- [15] Y. Kou, C.-T. Lu, S. S, and Y.-P. Huang. Survey of Fraud Detection Techniques. ISNSC, 2:749–754, 2004.
- [16] T. Milo, S. Novgorodov, and W.-C. Tan. Interactive Rule Refinement for Fraud Detection. In *EDBT*, 2018.
- [17] T. M. Mitchell. *Machine learning*. McGraw Hill series in computer science. McGraw-Hill, 1997.
- [18] S. J. Pan and Q. Yang. A survey on transfer learning. *Knowledge and Data Engineering, IEEE Transactions on*, 22(10):1345–1359, 2010.
- [19] R. Ramakrishnan and J. Gehrke. *Database management systems (3rd ed.)*. McGraw-Hill, 2003.
- [20] N. Segev, M. Harel, S. Mannor, K. Crammer, and R. El-Yaniv. Learn on source, refine on target: A model transfer learning framework with random forests. *TPAMI*, 2016.
- [21] B. W. Silverman. Density estimation for statistics and data analysis, volume 26. CRC press, 1986.
- [22] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: a core of semantic knowledge unifying wordnet and wikipedia. In WWW, 2007.
- [23] C. Sun, N. Rampalli, F. Yang, and A. Doan. Chimera: Large-scale classification using machine learning, rules, and crowdsourcing. *PVLDB*, 7(13):1529–1540, 2014.
- [24] Data Sharing for the Prevention of Fraud: https://gov.uk/government/publications/data-sharing-for-the-preventionof-fraud-code-of-practice.
- [25] G. Widmer and M. Kubat. Learning in the presence of concept drift and hidden contexts. *Machine Learning*, 23(1):69–101, 1996.
- [26] Common types of western union scams. https://www.westernunion.com/us/en/fraudawareness/fraud-types.html.