

# ARMS: Automated rules management system for fraud detection

David Aparício  
david.aparicio@feedzai.com  
Feedzai

Ricardo Barata  
ricardo.barata@feedzai.com  
Feedzai

João Bravo  
joao.bravo@feedzai.com  
Feedzai

João Tiago Ascensão  
joao.ascensao@feedzai.com  
Feedzai

Pedro Bizarro  
pedro.bizarro@feedzai.com  
Feedzai

## ABSTRACT

Fraud detection is essential in financial services, with the potential of greatly reducing criminal activities and saving considerable resources for businesses and customers. We address online fraud detection, which consists of classifying incoming transactions as either legitimate or fraudulent in real-time. Modern fraud detection systems consist of a machine learning model and rules defined by human experts. Often, the rules performance degrades over time due to concept drift, especially of adversarial nature. Furthermore, they can be costly to maintain, either because they are computationally expensive or because they send transactions for manual review. We propose ARMS, an automated rules management system that evaluates the contribution of individual rules and optimizes the set of active rules using heuristic search and a user-defined loss-function. It complies with critical domain-specific requirements, such as handling different actions (e.g., accept, alert, and decline), priorities, blacklists, and large datasets (i.e., hundreds of rules and millions of transactions). We use ARMS to optimize the rule-based systems of two real-world clients. Results show that it can maintain the original systems' performance (e.g., recall, or false-positive rate) using only a fraction of the original rules ( $\approx 50\%$  in one case, and  $\approx 20\%$  in the other).

## CCS CONCEPTS

• **Theory of computation** → **Optimization with randomized search heuristics**; • **Software and its engineering** → **Genetic programming**; • **Applied computing** → *Online banking*; *Online shopping*; *Secure online transactions*;

## KEYWORDS

fraud detection; genetic programming; evolutionary algorithms; greedy algorithms; randomized search

### ACM Reference Format:

David Aparício, Ricardo Barata, João Bravo, João Tiago Ascensão, and Pedro Bizarro. 2020. ARMS: Automated rules management system for fraud detection. In *Proceedings of KDD '20 (submitted)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

*KDD '20 (submitted)*, August 22–27, San Diego, CA, USA

© 2020 Copyright held by the owner/author(s).

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 INTRODUCTION

Financial institutions, merchants, customers, and government agencies alike suffer fraud-related losses daily, including credit card theft and other scams. Financial fraud consists of someone inappropriately obtaining the details of a payment card (e.g., a credit/debit card) and using it to make unauthorized transactions. Frequently, the cardholder detects such illicit usage and initiates a dispute with the bank to be reimbursed (a *chargeback*), at the expense of the merchant or bank that accepted the transaction. An over-conservative decision-maker might block all suspicious activity. However, this is far from optimal, as fraud patterns are not trivial, and it prevents legitimate economic activity. Therefore, it is essential to adjust automated fraud detection systems to the risk profile of the client.

Modern automated fraud detection systems consist of a machine learning (ML) model followed by a rule-based system. The model *scores* the transaction. The rule-based system uses the score and triggers of manually defined rules to decide an action (i.e., accept, alert, or decline the transaction). Rule-based systems with many rules are complex, hard to maintain, and frequently computationally expensive. An ideal system has only a minimum set of rules that ensure performance while preserving requirements and alerts low.

Our main contributions are the following:

- (1) Identifying a new problem: how to properly evaluate a complex rules system (taking into account overlapping rule triggers with different rule priorities and blacklists)? (Section 2).
- (2) Proposing ARMS (Figure 1), a framework which handles all bookkeeping necessary to correctly evaluate such rules systems (Sections 3.1–3.4).
- (3) Exploring optimization methods (namely random search, greedy expansion, and genetic programming) to improve the original system according to user-defined criteria (Sections 3.5–3.8).
- (4) Evaluating our proposed solutions on both synthetic and real data, demonstrating improvements to existing rules systems deployed at Feedzai (Section 4).

Evaluating the performance of the whole fraud detection system is simple: given the fraud labels (i.e., the chargebacks) and the historical decisions, we compute performance metrics (e.g., recall at a given false positive rate or FPR). However, it is not enough to analyze the performance of each rule by itself. We need to consider how it contributes to the entire system as its triggers may overlap with other rules with different decisions and priorities. Blacklists are another source of dependencies. Blacklisting rules, when triggered due to fraudulent behavior, blacklist the user (or email, or card) so that their future transactions are promptly declined. Deactivating blacklisting rules has side effects on the blacklists themselves and, therefore, in triggering or not rules that verify them.

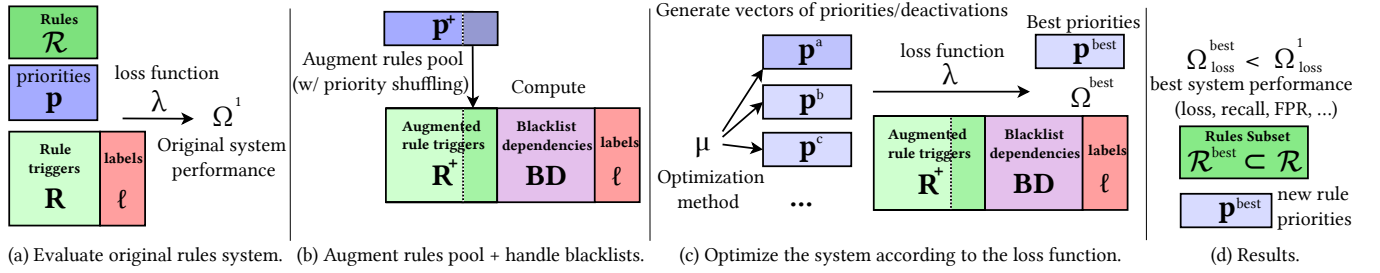


Figure 1: ARMS components: handling blacklists, priority shuffling, and optimizing a user-defined loss function.

In this work, we study the use case of a system with a pre-existing set of rules and priorities to optimize according to a user-defined objective function. As far as we know, we are the first to address the proper evaluation and optimization of such complex rules systems. A suitable goal is to minimize the number of rules and alerts while keeping the original system’s performance (e.g., recall). We explore three different methods (random, greedy, and genetic algorithms), using synthetic data and data sets from real-world online merchants.

Our results show that ARMS can significantly reduce the number of rules while maintaining the system’s performance. We stress that rules can depend on expensive aggregations (e.g., the average amount of the user’s transactions in the last month). Thus, ARMS brings meaningful gains in practical fraud detection settings.

We organize the remainder of the paper as follows. Section 2 gives an overview of fraud detection systems and discusses related work. Section 3 presents ARMS main components: handling blacklists, rules system evaluation, priority shuffling, and rules system optimization. Section 4 presents our results in synthetic data and real-world clients. Finally, we discuss our conclusions in Section 5.

## 2 BACKGROUND

### 2.1 Fraud detection

We focus on fraud detection in online payments, where a fraudster makes unauthorized transactions online. Fraud detection can be formulated as a binary classification task: each transaction is represented as a feature vector,  $\mathbf{z}$ , and labeled as either *fraudulent* (positive class,  $y = 1$ ) or *legitimate* (negative class,  $y = 0$ ). Other approaches frame it as an outlier detection problem [9] that treats fraudulent transactions as anomalies. Typical outlier detection is unsupervised, and often results in much lower performance.

We consider fraud detection as a two-step process. First, when a transaction occurs, a feature engineering step,  $g(\mathbf{z})$ , is applied to the raw features  $\mathbf{z}$ , resulting in processed features,  $\mathbf{x}$ . An example of a processed feature (a *profile*) is the number of transactions for a card in the last hour. Secondly, the automated fraud detection system evaluates the transactions and decides between three actions: to *accept* the transaction, to *decline* it, or to *alert* it to be manually reviewed (so that specialized fraud analysts investigate it and produce a final decision). Reviews are complicated (i.e., subject to human error) and expensive, as they require specialized knowledge and introduce unnecessary friction for legitimate transactions.

### 2.2 Automated fraud detection system

We consider an automated fraud detection system consisting of a machine learning model followed by a rule-based system.

**2.2.1 Machine learning model.** The supervised machine learning model trains offline using historical data. When evaluating a transaction, the model then produces a score,  $\hat{y} \in [0, 1]$ , that is typically the probability of fraud given the features,  $P(y = 1 | \mathbf{x})$ .

**2.2.2 Rule-based system.** Rules consist of conditions and corresponding actions. Depending on the action, the rules can be *accept*, *alert*, or *decline* rules. Rules may depend on the model score (e.g., if  $\hat{y} < 0.5$  then accept the transaction), and the features (e.g., if the transaction is above a *risky amount*, then alert/decline it). Since a transaction might trigger multiple rules with contradictory actions, priorities are necessary. Finally, rules can be switched on and off at any time. The rules system encapsulates all rules, their state (active or not), and priorities. Generally, the rule system is a function,  $f(\mathbf{x}, \hat{y})$ , that evaluates a list of rules and returns an action.

## 2.3 System evaluation

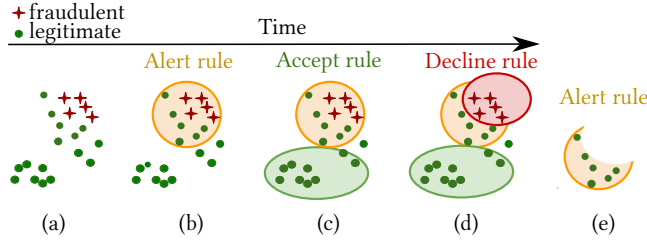
To assess system performance, we compare the system’s decisions with the labels coming from chargebacks or the fraud analysts’ decisions. Then, we compute the relevant performance metrics.

## 2.4 Rule evaluation

The rules system,  $f(\mathbf{x}, \hat{y})$ , receives the processed features and the model score and returns a decision to accept, alert, or decline. It comprises a set of rules,  $\mathcal{R} = (R_1, R_2, \dots, R_k)$ , applied *individually* on incoming transactions. Hence, transactions may trigger none, one, some, or all of these rules.

We aim to measure the contribution of individual rules to the system. Typically, at the time of the deployment of the system (i.e., after training with the latest data), rules and priorities perform well. However, as time goes by, fraud patterns change, and performance degrades. This degradation is acute in fraud detection, given the adversarial context (fraudsters often change their strategies). Whereas some rules remain beneficial, others may become redundant or even degrade the performance of the system. Figure 2 illustrates how an initially good rule can degrade over time. As rule-based systems remain in production for a long time, it is essential to monitor how individual rules are impacting the system, namely their fraud detection and computational performances (rules can be heavy to compute, e.g., if they depend on profiles).

One naive approach is to evaluate each rule independently by measuring how well its decisions match the labels (intuitively, accept rules should find legitimate transactions, while alert and decline rules should find fraudulent transactions). Then, if the rule’s performance is inadequate, it is discarded. Notwithstanding, this



**Figure 2: Rule degradation: (a) transactions in the feature space, (b-d) an alert, an accept, and a decline rule are progressively added, and trigger for some transactions. The alert rule has a good ratio of correct alerts when added, but by the end it is alerting only legitimate transactions (e).**

approach is problematic and insufficient because it disregards interactions between rules. Consider the following examples:

- Low-priority rules can perform outstandingly when no high-priority rules are triggered (e.g., in specific corner cases), but perform very poorly if used individually.
- Turning off high-priority rules allows lower-priorities rules to act; this can lead to different decisions by the system.

Instead of this naive approach, we build a rules management system that takes into account the interactions between rules with different actions and different priorities when evaluating them.

## 2.5 State-of-the-art

We review current work on the optimization of rule-based systems using search heuristics. Table 1 shows an overview of the methods.

Ishibuchi et al. propose a method to maximize correctly classified instances, while reducing the number of rules, using genetic programming [7]. This approach is not sufficiently flexible for the fraud detection use-case, as the client (e.g., a merchant or bank) may want to optimize for other metrics (e.g., recall, or a combination of metrics). Moreover, their method neglects priorities, and it is not clear if it scales up well for fraud detection data sets with millions of transactions (they used the Iris data set [5], which consists of 150 records). In a later study, they compare multiple heuristics, namely greedy search and genetic programming, in four small data sets [8].

Some approaches target specific use-cases, namely financial trading [1] or opinion mining [10]. Besides the domain, another crucial difference between our research and the work by Allen et al. is that, instead of learning new rules, we optimize existing rule systems.

Rosset et al. describe a method that learns and selects rules for telecommunication fraud detection [11]. Like us, the authors stress the importance of choosing a *good* set of rules, instead of a set of *good* rules. However, we target online transaction fraud and optimize a more complex system with priorities and blacklists.

**Table 1: Comparison of rules management systems.**

	[8]	[1]	[10]	[11]	[3]	[6]	ARMS
various rule actions	X	X	X	X	X	X	✓
rule priorities	X	X	X	X	*	X	✓
> million instances	X	X	X	X	X	X	✓
user-defined loss	X	X	X	X	X	X	✓
blacklists	X	X	X	X	X	X	✓
rule learning	X	✓	X	✓	X	X	X

\* optimizes rule weights instead

**Table 2: Notation.**

Features	$\mathcal{X} = (X_1, X_2, \dots, X_m)$
Rules	$\mathcal{R} = (R_1, R_2, \dots, R_k)$
Priority space	$\mathbb{P} = \{p \in \mathbb{Z} \mid p \geq -1\}$
Rule priority	$p_i \in \mathbb{P}$
Rule active condition	$p_i > -1$
Rules priority vector	$\mathbf{p} = (p_1, p_2, \dots, p_k)$
Priority-action map	$a : p_i \rightarrow \{\text{accept}, \text{alert}, \text{decline}\}$
Transaction feature vector	$\mathbf{x} = (x_1, x_2, \dots, x_m)$
Transactions	$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$
Transaction rules vector	$\mathbf{r} = (\{p_1, -1\}, \{p_2, -1\}, \dots, \{p_k, -1\})$
Rules triggers matrix	$\mathbf{R} = [\mathbf{r}_{x_1}, \mathbf{r}_{x_2}, \dots, \mathbf{r}_{x_n}]^T$
Labels vector	$\boldsymbol{\ell} = [\ell_{x_1}, \ell_{x_2}, \dots, \ell_{x_n}]^T$
Blacklist updater rules	$\mathcal{B}^u \subset \mathcal{R}$
Blacklist checker rules	$\mathcal{B}^c \subset \mathcal{R}$
Loss function	$\lambda$
Performance	$\Omega$ (contains $\Omega_{\text{loss}}, \Omega_{\text{recall}}$ , etc.)

Duman et al. propose a system combining genetic programming and scatter search to optimize rule *weights* and other parameters [3]. Similar to our case, the rules are based on expert knowledge and suffer from concept drift. Each rule has a weight corresponding to its contribution to a fraud score, unlike our work, which considers priorities to activate a single rule. Furthermore, Duman et al. do not consider blacklists, different rule actions (e.g., accept, alert, and decline), and uses a predefined fitness function that minimizes the money loss. Additionally, the most substantial data set considered contains only  $\approx 250$  thousand transactions and 43 rules and parameters. They report money savings of 212% at the cost of a 67% increase in false positives, and, after manual tuning, they settled for a system with savings of 189% and a 35% increase in false positives.

Gianini et al. optimize a system of 51 rules using a game theory approach [6]. They measure rule importance using Shapley values [13] as a measure of contribution to the system. They propose two strategies: (1) select the  $n$  rules with highest Shapley values (and deactivate the others) and (2) greedy expansion of the set of rules using the Shapley values of the rules. Both strategies performed identically and were able to reduce the number of rules down to 30 while maintaining the original system’s F-score. Like Duman et al. [3], this approach disregards essential constraints of the fraud detection system we are considering: rule priorities, rule actions, blacklists, and support for a user-defined loss function.

## 3 ARMS

We start this section with an overview of ARMS. Then, we describe in detail each of its main components: handling blacklisting rules, the evaluation of the rule-based system, rule priority shuffle, and, finally, the optimization strategies to select rules.

### 3.1 System overview

Algorithm 1 gives a general view of ARMS. We refer the reader to Table 2 for the notation used throughout this work.

ARMS receives the following information as **inputs**:

- **Features.** A vector of features,  $\mathcal{X}$  (e.g., username, email).
- **Transactions.** A matrix  $\mathbf{X}^{n \times m}$  containing the values of the  $m$  features for each of the  $n$  transactions. It is needed to compute

blacklists (i.e., to know blacklisted values for each feature, e.g., username = *fraudster91*).

- **Triggers or activations.** A matrix  $\mathbf{R}^{n \times k}$  containing the rule triggers of the  $k$  rules for each of the  $n$  transactions. Each cell  $R_{ij} = -1$  if rule  $R_j$  did not trigger for transaction  $\mathbf{x}_i$  or  $R_{ij} = p_j$  (i.e., the rule's priority) if it did.
- **Labels.** A vector with the label for each transaction,  $\ell$ .
- **Priorities.** A vector with the priority of each rule,  $\mathbf{p}$ .
- **Actions.** A map,  $a$ , mapping rule priorities to actions (i.e., accept, alert or decline).
- **Blacklisting rules.** A set of blacklisting rules containing rules that update the blacklist,  $\mathcal{B}^u$ , and rules that check it,  $\mathcal{B}^c$ .
- **Method.** Optimization strategy,  $\mu$  (i.e., random search, greedy expansion, or genetic programming).
- **Loss function.** A loss function,  $\lambda$ , defined by the user.
- **Priority shuffle.** A boolean,  $arp$ , specifying whether to augment the rules pool  $\mathbf{R}$  by cloning rules with different priorities.
- **Optimization parameters.** Set of parameters,  $\theta$ , which are specific to the optimization strategy (e.g., population size or mutation probability for the genetic algorithm or the number of evaluations for the random search).

ARMS starts by addressing the blacklist dependencies (line 1 of Algorithm 1; details in Section 3.2). Then, ARMS evaluates the original system's performance,  $\Omega^1$  (line 2 of Algorithm 1; Section 3.3). This evaluation runs *before* optimization because the loss function often depends on the original performance (e.g., optimize the FPR, while maintaining recall). Afterwards, ARMS augments the rules pool, if the user so desires (lines 3-4 of Algorithm 1; Section 3.4). This adds new rules with the same triggers as existing rules, but with different priorities. The rationale is that changing priorities might improve the system. Finally, ARMS optimizes the rules system (line 5 of Algorithm 1; Section 3.5). In essence, ARMS turns off rules and changes their priorities, obtaining a new priority vector,  $\mathbf{p}^{best}$ , to reduce the loss of the system,  $\Omega^{best}$ .

---

#### Algorithm 1 ARMS: Automated Rules Management System.

---

**Input:** Vector  $\mathcal{X}$ , Matrix  $\mathbf{X}$ , Matrix  $\mathbf{R}$ , vector  $\ell$ , vector  $\mathbf{p}$ , map  $a$ , set  $\mathcal{B}$ , loss function  $\lambda$ , method  $\mu$ , parameter  $arp \in \{0, 1\}$ , parameters  $\theta$

**Output:** Vector  $\mathbf{p}^{best}$ , performance  $\Omega^{best}$

- ```

1:  $\mathbf{BD} \leftarrow \text{COMPUTEBLACKLISTDEPENDENCIES}(\mathbf{R}, \mathbf{X}, \mathcal{X}, \mathcal{B})$ 
2:  $\Omega^1 \leftarrow \text{EVALUATE}(\mathbf{X}, \mathbf{R}, \ell, \mathbf{p}, a, \mathcal{B}, \mathbf{BD}, \lambda)$ 
3: if  $arp = 1$  then
4:    $\mathbf{R} \leftarrow \text{AUGMENTRULESPool}(\mathbf{R}, \mathbf{p}, a)$ 
5:  $(\mathbf{p}^{best}, \Omega^{best}) \leftarrow \mu.\text{OPTIMIZE}(\mathbf{X}, \mathbf{R}, \ell, \mathbf{p}, a, \mathbf{BD}, \lambda, \Omega^1, \theta)$ 

```
- 

### 3.2 Handling blacklists

Both analysts and rules can blacklist entities. If an analyst finds transaction  $\mathbf{x}$  to be fraudulent, they can blacklist some of its entities (e.g., in the future, always decline transactions from the email used in transaction  $\mathbf{x}$ ). Similarly, *blacklist updater rules* add entities to the blacklist when they trigger. Other rules, called *blacklist checker rules*, trigger when a transaction contains a blacklisted entity.

Therefore, blacklist rules have side effects. Deactivating blacklist updater rules can lead to blacklist checker rules not triggering, and affect the system's performance. Thus, we need to take this

into account when evaluating the system. For this purpose, ARMS keeps a state of the blacklists and manages them according to the interaction between blacklist updater and blacklist checker rules (for a detailed description, we refer to Supplementary Algorithm S1).

### 3.3 Rules system evaluation

ARMS evaluates (Algorithm 2) the original system and the configurations produced by the optimization strategies (Section 3.5). It creates an empty confusion matrix,  $\mathbf{V}$  (line 2), to be updated by traversing each transaction,  $\mathbf{x} \in \mathbf{X}$ , alongside its rule triggers,  $\mathbf{r}_i \in \mathbf{R}$ , and its label,  $\ell_i \in \ell$  (lines 3-9). For each transaction:

- (1) ARMS computes the activations  $\mathbf{r}'_i$  (i.e., what rules are active and with what priority), using priority vector  $\mathbf{p}$  (line 4). When ARMS is evaluating the original system,  $\mathbf{p}$  contains the original rules' priorities; however, rule priority shuffling and optimization strategies generate variations of  $\mathbf{p}$ .
- (2) ARMS checks whether to turn off any blacklist checker rules as a side-effect and stores that in  $\mathbf{r}''_i$  (line 5).
- (3) ARMS obtains the final decision,  $o_i$ , from  $\mathbf{r}''_i$ , i.e., to accept, alert, or decline (line 6). It is the action of the highest priority rule triggered for the transaction that is active.
- (4) ARMS evaluates the decision,  $o_i$ , against the label  $\ell_i$ , storing it in  $v_i$  (line 7). Accepting a legitimate transaction is a true negative. Declining/alerting a legitimate transaction is a false positive. Declining/alerting a fraudulent transaction is a true positive. Accepting a fraudulent transaction is a false negative. The confusion matrix,  $\mathbf{V}$ , is updated with  $v_i$  (line 8).

Finally, ARMS uses the confusion matrix  $\mathbf{V}$  to compute the rule configuration's performance,  $\Omega$ , based on a user-defined loss function,  $\lambda$  (line 9). The loss function allows optimizing (e.g., minimize the number of active rules, maximize recall) and satisfying metrics or constraints (e.g., keep the original system's FPR). We discuss loss functions used in synthetic data and real-world clients in Section 4.

---

#### Algorithm 2 Rules system evaluation.

---

- ```

1: function EVALUATE( $\mathbf{X}, \mathbf{R}, \ell, \mathbf{p}, a, \mathcal{B}, \mathbf{BD}, \lambda$ )
2:    $\mathbf{V} \leftarrow \text{INITCONFUSIONMATRIX}()$ 
3:   for all  $\mathbf{x} \in \mathbf{X}, \mathbf{r}_i \in \mathbf{R}, \ell_i \in \ell$  do
4:      $\mathbf{r}'_i \leftarrow \text{MASK}(\mathbf{r}_i, \mathbf{p})$ 
5:      $\mathbf{r}''_i \leftarrow \text{HANDLEBD}(\mathbf{r}'_i, \mathcal{B}, \mathbf{BD}[\mathbf{x}])$ 
6:      $o_i \leftarrow a(\text{MAX}(\mathbf{r}''_i))$ 
7:      $v_i \leftarrow \text{GETTRUTHVALUE}(o_i, \ell_i)$ 
8:      $\mathbf{V} \leftarrow \text{UPDATECONFUSIONMATRIX}(\mathbf{V}, v_i)$ 
9:    $\Omega \leftarrow \lambda(\mathbf{V})$ 
10:  return  $\Omega$ 

```
- 

### 3.4 Priority shuffle

Initial rule priorities require expert knowledge and are defined by clients or fraud analysts. Over time, however, the system requires adjusted priorities to deal with concept drift and incorporate emerging knowledge (e.g., new rules). In this section we discuss how ARMS addresses priority shuffling for optimization (Section 3.5). First, we discuss how ARMS changes the priority of individual rules. Then, we discuss how ARMS can augment the initial rules pool by cloning existing rules and assigning them alternative priorities.

**3.4.1 Random priority shuffle.** Since the system might have many rules and many possible priorities, the search space of all possible rule priorities can be gigantic. A more efficient alternative for such cases is to use random priority shuffle. For a given rule  $r_i$  with priority  $p_i$ , ARMS changes its priority to  $p_j \neq p_i$  with the same action, i.e.,  $a_i = a_j$ . The new rule priority is sampled considering uniform probabilities. Consider the illustrative example with three types of accept rules: *weak accept* with priority 1, *strong accept* with priority 3, and *whitelist accept* with priority 5. Random priority shuffle can, for example, change the priority of a strong accept to either 1 (weak accept) or 5 (whitelist accept).

**3.4.2 Augment rules pool.** Another option is to augment the initial pool by cloning existing rules (i.e., same triggers), but assigning them different priorities. Starting from the existing priorities,  $\mathbf{p}$ , we create variants for each  $p_i \in \mathbf{p}$ , with all possible alternative priorities with the same action,  $E$ . Then, for each  $p_j \in E$ , ARMS adds a new vector (a new "rule") with the same triggers as the original rule and the new priority,  $p_j$ , to the rules triggers matrix,  $\mathbf{R}$ .

### 3.5 Optimization strategies

ARMS uses two fundamental mechanisms to optimize a rule-based system: deactivate underperforming rules and change priorities. It is unfeasible to test all possible combinations. Instead, we employ three heuristics (methods): random search (Section 3.6), greedy expansion (Section 3.7), and genetic programming (Section 3.8).

First, we give an overview of ARMS optimization (Algorithm 3), as methods share a similar structure. The original system (i.e., with rule priorities  $\mathbf{p}$  and performance  $\Omega^1$ ) is the one to beat (line 1). Until meeting a predefined stopping criteria (line 3), ARMS generates new priority vectors,  $\mathbf{p}'$ , which are variations of the original  $\mathbf{p}$  (line 4). The criteria can be to stop after  $k$  hours, after computing  $n$  variations, or when the loss between consecutive iterations does not improve above a threshold  $\epsilon$ . ARMS saves the variation with lowest loss it finds,  $\mathbf{p}^{\text{best}}$ , alongside its performance  $\Omega^{\text{best}}$ , and returns them to the user (lines 5-8). The fundamental difference between methods is how they generate the variations,  $\mathbf{p}'$ .

---

#### Algorithm 3 ARMS optimization.

---

$\theta$ : parameters of the method

```

1: function  $\mu.\text{OPTIMIZE}(X, \mathbf{R}, \ell, \mathbf{p}, a, \text{BD}, \lambda, \Omega^1, \theta)$ 
2:    $(\mathbf{p}^{\text{best}}, \Omega^{\text{best}}) \leftarrow (\mathbf{p}, \Omega^1)$ 
3:   while  $\text{STOPPINGCRITERIANOTMET}()$  do
4:     generate a new  $\mathbf{p}'$  from  $\mathbf{p}$ 
5:     if  $\mathbf{p}'$  is the best so far then
6:       save it as  $\mathbf{p}^{\text{best}}$ 
7:       save its performance as  $\Omega^{\text{best}}$ 
8:   return  $(\mathbf{p}^{\text{best}}, \Omega^{\text{best}})$ 

```

---

### 3.6 Random search

A straightforward approach is to generate random rules priority vectors,  $\mathbf{p}'$ , and evaluate them against the original  $\mathbf{p}$ , saving the best rule configuration  $\mathbf{p}'$  that it found. While this approach seems naive, it is a natural baseline that can be better and less expensive than grid or manual searches [2]. Random search has two parameters:

- *Rule shutoff probability*,  $\rho$ . Percentage of rules to deactivate, e.g., if  $\rho = 50\%$ , then ARMS turns off  $\approx 50\%$  of the rules.
- *Rule priority shuffle probability*,  $\gamma$ . Percentage of rules with priorities changed, e.g., if  $\gamma = 50\%$ , then ARMS generates new priority vectors for  $\approx 50\%$  of the rules.

For more detail, we refer to Supplementary Algorithm S2.

### 3.7 Greedy expansion

ARMS contains a greedy expansion module, that starts from a set of inactive rules and greedily turns on rules, one at the time. Greedy solutions are not guaranteed to find the global optimum. Consider the following example, where we want to optimize recall and rules  $R_1$ ,  $R_2$ , and  $R_3$  have recall 70%, 69%, and 20%, respectively. A greedy solution would pick  $R_1$  first. Now, imagine that rules  $R_2$  and  $R_3$  are detrimental to  $R_1$ , i.e., the system becomes worse if we combine  $R_1$  with either  $R_2$  or  $R_3$ . Hence, the final solution is a system with only  $R_1$ . Imagine, however, that  $R_2$  and  $R_3$  are somewhat complementary, and that, when combined, the system's recall is  $> 70\%$ . Then, the global optimum is  $> 70\%$ , and the greedy solution is not optimal. Nevertheless, greedy heuristics can find useful solutions in a reasonable time.

For more detail, we refer to Supplementary Algorithm S3.

### 3.8 Genetic programming

Genetic programming is standard in classification tasks [4], such as fraud detection. It continuously improves a population of solutions by combining them using crossovers and random mutations, while keeping a fraction of the best solutions for the next iteration.

In our case, we build a population of random rule configurations and improve them with genetic programming. The algorithm has three parameters:

- *Population size*,  $\psi$ . Number of configurations per iteration, e.g., if  $\psi = 100$ , ARMS evaluates 100 different rule configurations per iteration.
- *Survivors fraction*,  $\alpha$ . Fraction of the top configurations that survive for the next iteration, e.g., if  $\psi = 100$  and  $\alpha = 20\%$ , only the 20 best solutions *survive* for the next iteration. If  $\alpha$  is high, then we might achieve higher variability but get stuck trying to improve bad solutions. If  $\alpha$  is low, then the lack of variability might prevent the system from reaching a good solution.
- *Mutation probability*,  $\rho$ . The percentage of rules subject to random mutation, e.g., if  $\rho = 20\%$ , then 20% of the rules are randomly mutated (i.e., the *child* rule configuration mutates the *parents* rules configuration). If  $\rho$  is high, we leave little room for genetic optimization and are essentially doing a random search. If  $\rho$  is low, we are more dependent on finding good parent configurations.

For more detail, we refer to Supplementary Algorithm S4.

## 4 EXPERIMENTS AND RESULTS

We test the following hypotheses: **(h1)** ARMS turns off rules and, at least, maintains system performance, **(h2)** ARMS changes the priority of rules and improves system performance, **(h3)** results are stable (i.e., similar across folds).

#### 4.1 Synthetic data

Since we can not find public data sets similar to our own, we use synthetic data to test hypotheses (**h1–h2**). Later, we also test (**h1–h3**) in real datasets. We generate 225k labels with a fraud rate of 5% (i.e., 11250 positive labels) and simulate accept, alert, and decline rules from the labels. The support of a rule corresponds to how many times it triggers. An accept rule with negative predictive value (NPV) of  $k\%$  is correct  $k\%$  of the times that it triggers (i.e., out of all triggers,  $k\%$  will be true negatives). The same goes for the precision (PPV) of an alert or decline rule (i.e., out of all triggers,  $k\%$  will be true positives). We sample the support, NPV and precision from Gaussian distributions and use 10 different priority levels (for details see Supplementary Section A.1) and divide the data set into three splits: train, validation, and test, with 75k "transactions" each.

#### 4.2 Methodology

We run ARMS on the train set and do parameter tuning in the validation set. We detail the parameter space in Supplementary Section A.2. We ensure that results are comparable between random search and genetic programming by keeping the number of rule configuration evaluations fixed (i.e.,  $n = 300k$ ).

We optimize the loss function from Equation 1 with  $\alpha = 0.1$ ,  $\beta = 0.5$ , and  $\gamma = 0.4$ . Note that  $\Omega^1$  and  $\Omega'$  are the performance of the original system and of a configuration found by ARMS, respectively;  $\Omega_{rules\%}$  is the percentage of active rules,  $\Omega_{recall\%}$  is the recall, and  $\Omega_{alert\%}$  is the alert rate.

$$\lambda(\mathbf{R}') = \alpha * \Omega'_{rules\%} - \beta * \Omega'_{recall} + \gamma * \Omega'_{alerts\%} \quad (1)$$

Finally, we evaluate the four final methods in the test set: the original system, and the best rule system configuration found by random search, greedy expansion, and genetic programming.

#### 4.3 Results on synthetic data

After running parameter tuning (note that the greedy expansion method does not have any parameter), we find that the following parameters were the best:

- **Random search:**  $\rho = 40\%$ .
- **Genetic programming:**  $\rho = 10\%$ ,  $\psi = 30$ ,  $\alpha = 5\%$ .

For brevity, we omit results for other parameters; we do a more thorough analysis of the parameters in real data sets (Section 4.6). We observe that all methods improved upon the original system, and that genetic programming was the one with highest performance (Table 3). When we ran augmented rules pool (ARP) before optimization, results consistently improved. Thus, we verify hypothesis (**h1**) and (**h2**).

Table 3: Performance of ARMS in synthetic data.

	recall	alerts %	rules off	loss
original	13.11%	0.779%	none	0.0376
random	79.53%	1.013%	38 (38.8%)	-0.1837
greedy	54.42%	1.746%	34 (34.7%)	-0.1998
genetic	52.82%	1.067%	45 (45.9%)	-0.2058
greedy w/ arp	53.30%	1.107%	43 (43.9%)	-0.2060
genetic w/ arp	53.09%	0.97%	45 (45.9%)	<b>-0.2075</b>

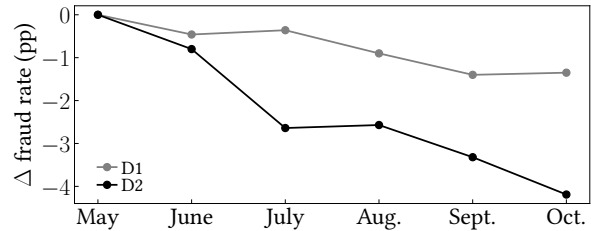


Figure 3: Datasets fraud rate evolution (concept drift).

#### 4.4 Real-world data sets

We evaluate ARMS on representative samples of real-world data sets of two online merchants. In both cases, an automated fraud detection system actively scores transactions in production. We collected the rule triggers, model decisions, and blacklists. The data sets comprise dozens of rules, with different actions (i.e., accept, alert, and decline) and multiple priorities (more details in Supplementary Section A.3). For privacy compliance, we refer to the data sets simply as **D1** and **D2**.

The data covers six months of transactions. We divide each data set in four sequential and overlapping folds of three months each (for temporal cross-validation, detailed in Section 4.5.3) and split each fold into three sequential sets (train, validation, and test) of one month each.

Unless explicitly stated, when we mention *fraud*, we are referring to validated fraud (i.e., chargebacks or fraud confirmed by analysts, not transactions declined by the automated fraud detection system). Due to the adversarial setting and other factors, we observe concept drift in both data sets. Figure 3 shows the evolution of the fraud rate in **D1** and **D2** (with May 2018 as reference), highlighting the system's ability to reduce fraud over time.

While both clients are online merchants, they have three important differences:

- (1) **D1** has more non-verified declined transactions. It has  $\approx 14x$  more auto-declined transactions than confirmed frauds, due to the specific requirements of the client. Using automatically declined transactions for training is dangerous as it creates a feedback loop. Thus, we disregard them in training and validation but use them in testing so that results are comparable to a production setting. Moreover, for this dataset, ARMS does not optimize decline rules.
- (2) Only **D2** uses blacklists.
- (3) The active rules in **D2** changed multiple times during the period under study, while the rules in **D1** never changed.

#### 4.5 Methodology

**4.5.1 Optimization metrics (loss functions).** Online merchants are required to keep the fraud-to-gross rate (FTG) under a certain threshold, or else they face fines. Thus, a sensible approach is to minimize the FPR and ensure that recall is within the legal requirements. The system should be able to pick up all the necessary fraud (ideally, all of it) without declining legitimate transactions. Additionally, reducing the number of rules and alerts decreases the overall cost of the system. We use different loss functions for each data set, showing ARMS' ability to fit diverse use-cases:



- In **D1**, the FPR is artificially high due to the many transactions declined by the automated fraud detection system. Therefore, our focus is to remove rules,  $\Omega'_{rules\%}$ , and reduce alerts,  $\Omega'_{alert\%}$ , while maintaining approximately the same recall,  $\Omega'_{recall}$ , as the original rule-based system,  $\Omega^1_{recall}$ , (Equation 2). We use  $\alpha = \beta = \frac{1}{2}$ , thus giving equal importance to both objectives.
- In **D2**, the objective is to remove rules,  $\Omega'_{rules\%}$ , but also to improve recall,  $\Omega'_{recall}$ , while maintaining approximately the same FPR,  $\Omega'_{fpr}$ , as the original system,  $\Omega^1_{fpr}$ , (Equation 3). We use  $\alpha = 0.05$  and  $\beta = 0.95$ , thus attributing more importance to improving recall than to reducing the number of rules.

$$\lambda(\mathbf{R}') = \begin{cases} \alpha * \Omega'_{rules\%} + \beta * \Omega'_{alert\%} & \text{if } \Omega'_{recall} \geq 0.95 * \Omega^1_{recall} \\ \alpha + \beta + (\Omega^1_{recall} - \Omega'_{recall}) & \text{otherwise} \end{cases} \quad (2)$$

$$\lambda(\mathbf{R}') = \begin{cases} \alpha * \Omega'_{rules\%} - \beta * \Omega'_{recall} & \text{if } \Omega'_{fpr} \leq \Omega^1_{fpr} \\ \alpha + (\Omega^1_{fpr} - \Omega'_{fpr}) & \text{otherwise} \end{cases} \quad (3)$$

4.5.2 *Baselines*. We compare ARMS optimized rule systems against three baselines:

- (1) **Original system (All on)**: system with all rules and original priorities.
- (2) **Mandatory system (All off)**: system with no rules except for the ones that cannot be deactivated due to business reasons, with the original priorities.
- (3) **Random search**: generate  $r$  independent rule configurations, using different values of  $\rho$  (Section 3.6).

If ARMS finds rule systems better than baselines 1 and 2 by turning off rules, we successfully address **(h1)**. If it further improves its performance by also tuning rule priorities, we address **(h2)**.

4.5.3 *Temporal cross validation (TCV)*. We use TCV to verify **(h3)**. For each data set, we create four folds composed of three sets (i.e., train, validation, and test) of one month each. We train ARMS with different search heuristics and parameters on each train data set, evaluate the resulting configurations on the validation data set, and identify the best one for each heuristic. Finally, we evaluate the winners and the baselines on the test set.

4.5.4 *Optimization strategies*. We run ARMS with two different optimization strategies: greedy expansion (Section 3.7) and genetic programming (Section 3.8). Results for both are shown in Sections 4.6.2 and 4.6.3, respectively.

## 4.6 Results on real data

Unless stated otherwise, the results refer to rule configurations obtained in the train data of each fold and evaluated in the respective validation set. Results shown are always relative to the original system baseline and show the gains relative to the current system in production, i.e.,  $\Delta loss$  is the difference between the loss of the system being evaluated and the original one.

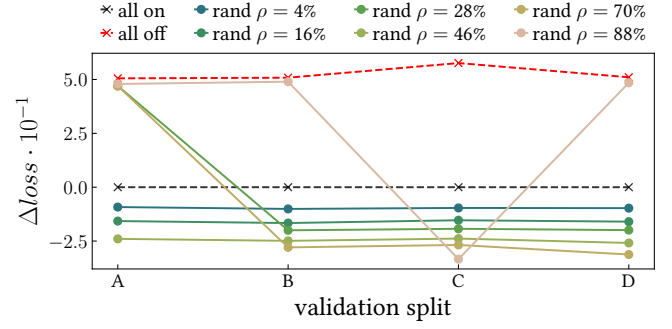


Figure 4: Baselines comparison in D1.

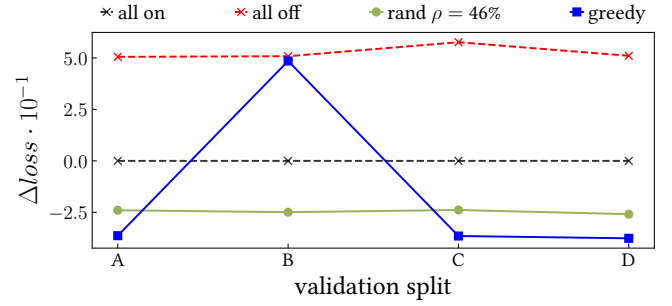


Figure 5: Greedy expansion against baselines in D1.

4.6.1 *Baselines comparison*. We compare the original system (*all on*) against the mandatory system (*all off*) and against random search, with  $n = 10000$  and  $\rho$  as a tunable parameter with values spaced out in 4% intervals (Figure 4 for **D1**). We observe that the mandatory system has a higher loss than the other systems, as it fails to meet the recall constraint from Equation 2. We also observe that random search is almost always superior to the original system, regardless of  $\rho$ . In a few cases, the random search is worse than the original system because it does not meet the recall constraint, namely with aggressive configurations (e.g.,  $\rho = 88\%$ ). On the other hand, aggressive random search (higher  $\rho$ ) can decrease the loss significantly, so there is a trade-off between being able to meet the recall constraints and lowering the loss. We observe similar behavior for **D2**, and thus omit results for brevity. Nevertheless, we show metrics besides the loss for **D2** in Supplementary Figure S1. From these results, we decide to use random search with  $\rho = 46\%$  for **D1**, and  $\rho = 58\%$  for **D2**, for the baselines, alongside the original system and the mandatory system for both data sets.

4.6.2 *Greedy expansion results*. We test greedy expansion with and without ARP. We find that ARP did not improve the system in **D1** or **D2**. One possible explanation is that greedy expansion yields simple systems with few rules, so it did not benefit from ARP. Another possibility is that the original priorities are already well-tuned for both data sets as they correspond to mature systems.

When compared against the baselines, the outcomes vary. For **D1**, the greedy expansion was superior to the baselines except for the second fold, where it failed to meet the constraints (Figure 5). In the other three folds, greedy expansion was able to remove  $\approx 75\%$  of the rules and reduce alerts. For **D2**, however, the greedy expansion

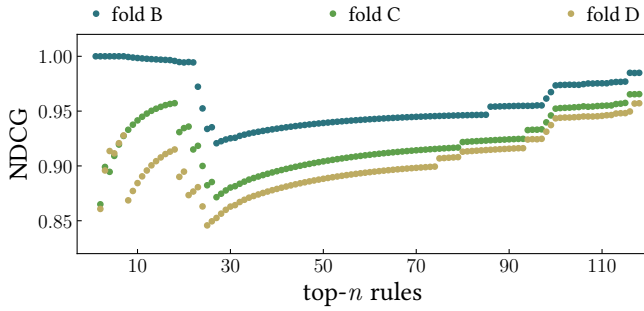


Figure 6: Greedy expansion rule order consistency in D1.

was worse than the baselines in two of the four folds, as it did not respect the constraints (Supplementary Figure S2).

We also evaluate the consistency between rules across folds. Recall that the greedy expansion obtains an ordered list of rules sequentially by importance. We compare the ordered lists across folds and compute their normalized discounted cumulative gain (NDCG) in Figure 6. We show results of the first fold of D1 compared with the other folds. We observe that rules are consistent across folds (NDCG values are consistently  $> 0.7$ ), but the NDCG line drops (e.g., important rules in fold A are more similar to important rules in fold B than in fold C). We observe similar behavior in D2 (omitted for brevity).

**4.6.3 Genetic programming results.** We evaluate how the genetic programming method (Section 3.8) improves fraud detection. Since our datasets are very big, we can not perform a grid search on all parameters. Thus, we have a three phase process.

First, we find a good set of default parameters. For this purpose, we set  $\psi = 100$  and do grid search on  $\alpha$  and  $\rho$ . We do  $n = 10000$  evaluations by default, i.e., for  $\psi = 100$ , then  $r = 100$  runs. We perform a grid search on  $\alpha \in [2\%, 5\%, 10\%, 20\%]$  and  $\rho \in [0\%, 2\%, 5\%, 10\%]$ . For D1, we find that  $\rho = 10\%$  outperforms the baselines across datasplits and that random search takes longer to achieve similar losses (e.g., for fold A; Supplementary Figure S3). The overall best parameters were found to be  $\rho = 10\%$ ,  $\alpha = 5\%$ , and  $\rho = 5\%$ .

Secondly, we study how each parameter influences the loss. For this purpose, we vary only one parameter at a time and keep the others at the default values. Since parameters  $r, \psi, \rho, \alpha$  are ordinal, we try 10 different values for each and see how increasing each parameter individually influences the loss. Figure 7 shows results for fold A of D1. We observe that, in general, increasing  $\rho$  and  $\alpha$  makes the performance worse; however the best  $\alpha$  is 10%, thus, keeping some of the best individual configurations is important. We also observe that  $r$  influences the loss much more than  $\psi$  (e.g., both  $(r = 100, \psi = 10)$  and  $(r = 10, \psi = 100)$  perform 1000 rule evaluations, but the first one leads to lower losses). Typically, the loss improves as you increase  $r$  and  $\psi$ , but it plateaus relatively quickly for both (i.e.,  $r$  at 300,  $\psi$  at 400). Similar conclusions hold for D2 (omitted for space concerns). We did not observe gains in changing rule priorities during genetic optimization.

Finally, we measure ARMS performance on the test sets. We compare ARMS using genetic programming against the baselines and ARMS using greedy search. To do this, we evaluate the rules deactivations suggested by ARMS (trained on the train sets and evaluated in the validation sets) on the respective test sets of each

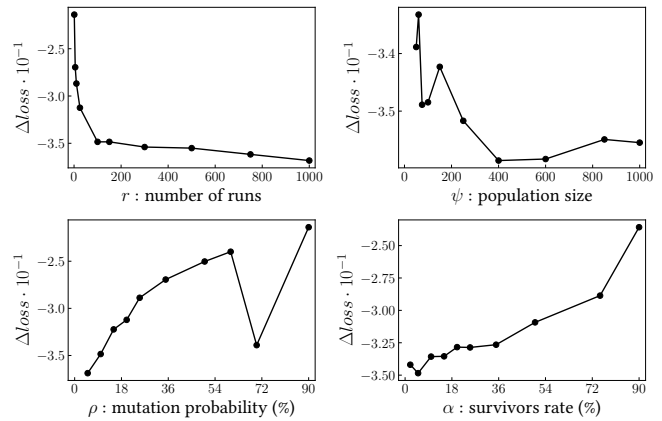


Figure 7: Genetic: influence of  $r, \psi, \rho, \alpha$  on fold A of D1.

fold. For D1, we evaluate the best rule configuration found by ARMS using  $r = 1000, \psi = 250, \alpha = 5\%, \rho = 5\%$ , and no priority shuffling. For D2, we evaluate the best rule configuration found by ARMS using  $r = 1000, \psi = 150, \alpha = 20\%, \rho = 5\%$ , and no priority shuffling.

For D1, we observe that greedy and genetic optimization performed similarly and better than random search with  $\rho = 46\%$  (Figure 8). For D2, we observe that random search and the genetic programming approaches perform similarly; the greedy method fails to comply to the constraints in two of the four folds (Figure 9).

In order to check the consistency of ARMS across data folds, we measure the Jaccard similarity [12] of the deactivated rules suggested by ARMS in different splits. We see that the Jaccard is

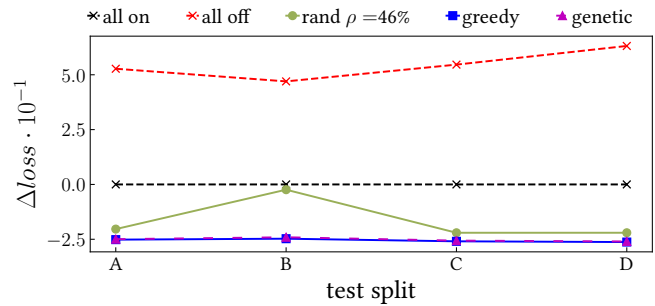


Figure 8: Performance of ARMS on the test sets of D1.

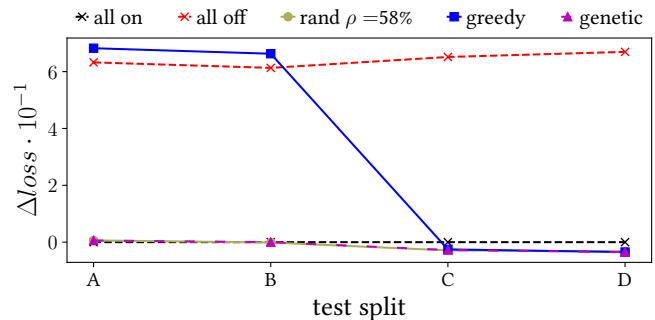


Figure 9: Performance of ARMS on the test sets of D2.



**Table 4: ARMS consistency results (i.e., across folds). We highlight in bold the lowest loss for each fold.**

	A	B	C	D		A	B	C	D		A	B	C	D		A	B	C	D
<b>A</b>	1	0.930	0.902	0.826	<b>A</b>	1	0.789	0.696	0.636	<b>A</b>	<b>0.275</b>	<b>0.344</b>	0.274	0.273	<b>A</b>	<b>-0.626</b>	<b>-0.613</b>	-0.651	-0.662
<b>B</b>	-	1	0.950	0.820	<b>B</b>	-	1	0.773	0.565	<b>B</b>	-	0.348	0.277	0.275	<b>B</b>	-	-0.612	-0.651	-0.662
<b>C</b>	-	-	1	0.829	<b>C</b>	-	-	1	0.708	<b>C</b>	-	-	<b>0.268</b>	0.267	<b>C</b>	-	-	<b>-0.678</b>	-0.696
<b>D</b>	-	-	-	1	<b>D</b>	-	-	-	1	<b>D</b>	-	-	-	<b>0.264</b>	<b>D</b>	-	-	-	<b>-0.704</b>

(a) Jaccard of removed rules (D1). (b) Jaccard of removed rules (D2).

(c) Loss on future folds (D1).

(d) Loss on future folds (D2).

higher for **D1** than **D2** (Table 4 (a)-(b)). The fact that **D2** rule set changes across folds obviously leads to intrinsically lower values (i.e., regardless of what ARMS deactivates). We also evaluate systems trained on a given fold in more recent folds (e.g., we train ARMS on fold A and evaluate it the test set of A, B, C and D). We observe that systems trained on older folds have good performance on more recent test sets (Table 4 (c)-(d)).

**4.6.4 Summary.** We evaluated ARMS on two big online merchants. For **D1**, ARMS using genetic programming (or greedy expansion) was able to remove  $\approx 50\%$  of the original 193 rules, while maintaining the original system performance (i.e., keeping 95% of the original recall). Thus, ARMS was able to improve the original system (**h1**). We also saw that results are stable across data-splits (**h3**). We did not see gains of using priority shuffling (**h2**). For **D2**, we observed that ARMS was able to remove  $\approx 80\%$  of the system rules while maintaining the original system performance (i.e., keeping a low FPR). Thus, ARMS improved the original system (**h1**). Similar to **D1**, we found evidence supporting (**h3**) but not (**h2**).

**4.6.5 Discussion.** Real-world transaction data sets for fraud detection pose several challenges. Auto-declines lead to unreliable labels, and thus we cannot verify if a system positive is a true positive, meaning that decline rules cannot be evaluated unless an analyst verifies auto-declines. In practice this is difficult because fraud analysts' time is a very limited resource. The two systems that we chose are also particularly hard to optimize since they have been in production for years and have been manually tuned by data scientists. Finally, we evaluated ARMS' performance on past transactions and did not measure its performance in production. We think that putting ARMS in production and continuously optimizing the rules system could lead to better results.

## 5 CONCLUSION

We have proposed ARMS, a framework that optimizes rules systems using search heuristics, namely random search, greedy expansion, and genetic programming. To the best of our knowledge, ARMS is the first to (1) handle different rule priorities and actions, (2) address blacklists side effects, and (3) optimize user-defined functions. These components are essential in real-world fraud detection systems. Our results in real-world clients demonstrate that ARMS is capable of maintaining the original system's performance while greatly reducing the number of rules (between 50% and 80%, in our experiments) and minimizing other metrics (e.g., alert rate).

Currently we are adding a rules suggestions module to ARMS, which is beyond the scope of this paper. In the future we also plan to incorporate a module to simultaneously tune the rules and the machine learning model threshold.

## ACKNOWLEDGMENTS

We want to thank the other members of Feedzai's research team, who always gave insightful suggestions. In particular, we want to give special thanks to Marco Sampaio, for reviewing the paper internally, and Patrícia Rodrigues, for starting ARMS.

## Note on reproducibility

We make available a binary of ARMS, the synthetic data described in Section 4.1 (as well as the script used to generate it), and all the necessary steps to reproduce our results from Section 4.3 at <https://github.com/feedzai/research-arms>. For privacy compliance, we can not share our clients data sets.

## REFERENCES

- [1] Franklin Allen and Risto Karjalainen. 1999. Using genetic algorithms to find technical trading rules. *Journal of financial Economics* 51, 2 (1999), 245–271.
- [2] James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *Journal of Machine Learning Research* 13, Feb (2012), 281–305.
- [3] Ekrem Duman and M Hamdi Ozelcik. 2011. Detecting credit card fraud by genetic algorithm and scatter search. *Expert Systems with Applications* 38, 10 (2011), 13057–13063.
- [4] Pedro G Espejo, Sebastián Ventura, and Francisco Herrera. 2009. A survey on the application of genetic programming to classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 40, 2 (2009), 121–144.
- [5] Ronald A Fisher. 1936. The use of multiple measurements in taxonomic problems. *Annals of eugenics* 7, 2 (1936), 179–188.
- [6] Gabriele Gianini, Leopold Ghemmogne Fossi, Corrado Mio, Olivier Caelen, Lionel Brunie, and Ernesto Damiani. 2020. Managing a pool of rules for credit card fraud detection by a Game Theory based approach. *Future Generation Computer Systems* 102 (2020), 549–561.
- [7] Hisao Ishibuchi, Ken Nozaki, Naohisa Yamamoto, and Hideo Tanaka. 1995. Selecting fuzzy if-then rules for classification problems using genetic algorithms. *IEEE Transactions on fuzzy systems* 3, 3 (1995), 260–270.
- [8] Hisao Ishibuchi and Takashi Yamamoto. 2004. Comparison of heuristic criteria for fuzzy rule selection in classification problems. *Fuzzy Optimization and Decision Making* 3, 2 (2004), 119–139.
- [9] Yufeng Kou, Chang-Tien Lu, Sirirat Sirwongwattana, and Yo-Ping Huang. 2004. Survey of fraud detection techniques. In *IEEE International Conference on Networking, Sensing and Control, 2004*, Vol. 2. IEEE, IEEE, 749–754.
- [10] Qian Liu, Zhiqiang Gao, Bing Liu, and Yuanlin Zhang. 2015. Automated rule selection for aspect extraction in opinion mining. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.
- [11] Saharon Rosset, Uzi Murad, Einat Neumann, Yizhak Idan, and Gadi Pinkas. 1999. Discovery of fraud rules for telecommunications challenges and solutions. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 409–413.
- [12] Cesare Baroni Urbani. 1980. A statistical table for the degree of coexistence between two species. *Oecologia* (1980), 287–289.
- [13] Eyal Winter. 2002. The shapley value. *Handbook of game theory with economic applications* 3 (2002), 2025–2054.

## A SUPPLEMENTARY MATERIALS

### A.1 Synthetic data

The set of rules comprises 8 accept rules, 30 review rules, and 60 decline rules. The support of the accept rules was sampled from a Gaussian distribution  $\mathcal{N}(45000, 22500^2)$ , while the support of the review and decline rules was sampled from  $\mathcal{N}(22.5, 225.0^2)$ . The NPV of accept rules was sampled from  $\mathcal{N}(0.75, 0.20^2)$ , while the precision of the alert and decline rules was sampled from  $\mathcal{N}(0.17, 0.05^2)$ .

Rules have ten possible priorities. Accept rules have priority  $p_a \in \{0, 1, 5, 6, 10\}$ , alert rules have priority  $p_l \in \{2, 4, 7, 9\}$ , and decline rules have priority  $p_d \in \{3, 8\}$ .

### A.2 Synthetic data parameter tuning

*A.2.1 Random search.* We use 16 mutation probabilities, i.e.,  $\rho \in [4\%, 94\%]$ , in intervals of 4%, i.e.,  $\rho = 4\%$ ,  $\rho = 8\%$ , ...

*A.2.2 Genetic programming.* We use three mutation probabilities, i.e.,  $\rho = 10\%$ ,  $\rho = 20\%$ , and  $\rho = 30\%$ . We use two population sizes, i.e.,  $\psi = 20$  and  $\psi = 30$ . Finally, we use two survivors fractions, i.e.,  $\alpha = 2\%$  and  $\alpha = 5\%$ .

### A.3 Real-world datasets

*A.3.1 D1.* The client has 198 rules, with one of three possible actions: accept, alert, and decline. Out of the 198 rules, 30 of them are accept rules, 89 are alert rules, and 79 are decline rules. Accept rules have four different priority levels  $p_a \in \{1, 8, 10, 15\}$ , alert rules have two  $p_a \in \{5, 11\}$ , and decline rules have three  $p_d \in \{6, 9, 12\}$ . If no rules are triggered, the default action is to accept the transaction.

The dataset contains few validated fraud, i.e., of the declined (by the model/rules) and fraudulent population of transactions, only a small portion was validated by analysts or via chargeback.

We note that decline rules and auto-declined transactions are ignored in the train and validation datasets. We make this choice because decline rules can not be validated. However, when we measure performance in the test set, decline rules are included in order to make results directly comparable to the results obtained in production.

We do temporal cross validation (TCV) with four folds and each set has one month of data.

*A.3.2 D2.* Unlike **D1**, which has the same activated rules for the whole period, in D2 the rules changed. During the seven months period a total of 13 rules were added, while some were removed, increasing the number of rules in the set from the original 77 to 90.

Rules have one of three outcomes: accept, alert, and alert&decline (this means that most auto-declined are verified, unlike in **D1**). From those, 6 are accept rules, 48 are alert rules, and 36 are alert&decline rules. Accept rules have priority  $p_a \in \{0, 5, 10\}$ , alert rules have priority  $p_l = 1$ , and alert&decline rules have priority  $p_d \in \{2, 4, 8\}$ . Three of the decline rules are blacklist checker rules, and all 36 alert&decline rules are blacklist updater rules.

Since **D2** has a high ratio of validated fraud, all rules are optimized by ARMS, however the auto-decline transactions are not used during the training process, but are present in the test set in order to make results directly comparable to the results obtained in production.

We do temporal cross validation (TCV) with four folds and each set has one month of data.

### A.4 Supplementary Algorithms and Figures

---

#### Algorithm S1 Blacklist propagation.

---

```

1: function COMPUTEBLACKLISTDEPENDENCIES( $R, X, \mathcal{X}, \mathcal{B}$ )
2:    $BL \leftarrow \{\}$ 
3:    $BD \leftarrow \{\}$ 
4:   for all  $x \in X$  do
5:     for all  $R_j \in \mathcal{B}^u$  do
6:       if  $r_j \neq -1$  then
7:         for all  $X_l \in \mathcal{X}$  that  $R_j$  blacklists do
8:            $BL[(R_j, X_l : x_l)].APPEND([x.time, +\infty])$ 
9:           for all  $R_q \in \mathcal{B}^c$  that checks  $X_l$  do
10:             $BD[x].ADD(R_j < R_q)$ 
11:         if  $r_j = -1$  then
12:           for all  $X_l \in \mathcal{X}$  that  $R_j$  can blacklist do
13:             if  $x.time$  is in any  $BL[(R_j, X_l : x_l)]$  then
14:                $r_j \leftarrow p_j$ 
15:         for all  $\{x_l \in x \mid x_l \text{ is in any active blacklist}\}$  do
16:           if  $\nexists R_q \in \mathcal{B}^c \mid p_q \neq -1$  then
17:             for all  $\{R_j \in \mathcal{B}^u \mid (R_j, X_l : x_l) \in BL\}$  do
18:                $BL[(R_j, X_l : x_l)].LAST() \leftarrow [x.time]$ 
19:         for all  $R_q \in \mathcal{B}^c$  do
20:           if  $r_q \neq -1$  and  $|(R_i < R_q) \in BD[x] \mid R_i \in \mathcal{B}^u| = 0$  then
21:              $BD[x].ADD(R_q < R_q)$ 
22:   return  $BD$ 

```

---



---

#### Algorithm S2 Random search optimization.

---

```

 $\theta$ : { rule shutoff probability  $\rho$ , rule priority shuffle probability  $\gamma$  }
1: function RANDOM.OPTIMIZE( $X, R, \ell, p, a, BD, \lambda, \Omega^1, \theta$ )
2:    $p^{best} \leftarrow p$ 
3:    $\Omega^{best} \leftarrow \Omega^1$ 
4:   while STOPPINGCRITERIANOTMET() do
5:      $p^{rand} \leftarrow p$ 
6:     for all  $p_i \in p^{rand}$  do
7:       with  $\gamma\%$  probability, do:
8:          $p_i \leftarrow \text{RANDOM.PRIORITYSHUFFLE}(p_i, a)$ 
9:       with  $\rho\%$  probability, do:
10:         $p_i \leftarrow -1$ 
11:      $\Omega^{rand} \leftarrow \text{EVALUATE}(X, R, \ell, p^{rand}, a, \mathcal{B}, BD, \lambda)$ 
12:     if  $\Omega_{loss}^{rand} < \Omega_{loss}^{best}$  then
13:        $\Omega^{best} \leftarrow \Omega^{rand}$ 
14:        $p^{best} \leftarrow p^{rand}$ 
15:   return  $(p^{best}, \Omega^{best})$ 

```

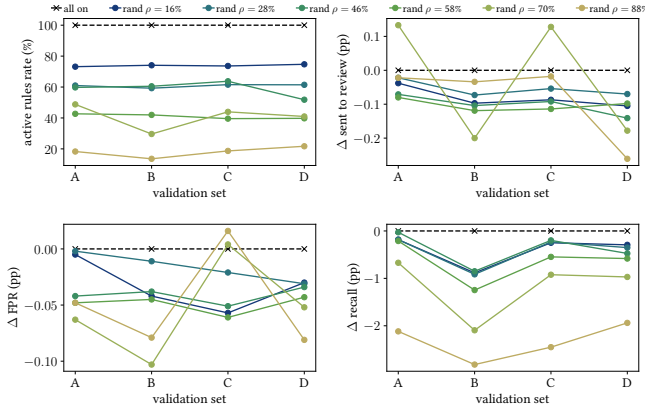
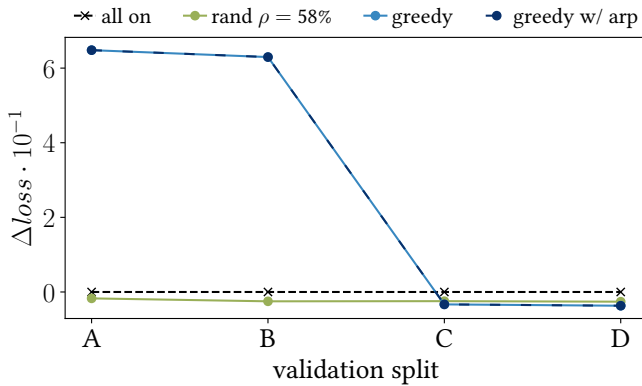
---

**Algorithm S3** Greedy expansion optimization.

```

 $\theta$ : {backtracking  $bt \in \{true, false\}$ }
1: function GREEDY.OPTIMIZE( $X, R, \ell, p, a, BD, \lambda, \Omega^1, \theta$ )
2:    $p^{best} \leftarrow p$ 
3:    $\Omega^{best} \leftarrow \Omega^1$ 
4:    $p^{keep} \leftarrow (-1, \dots, -1)$ 
5:    $p^{greedy} \leftarrow (-1, \dots, -1)$ 
6:    $Q \leftarrow \emptyset$ 
7:   while  $|Q| < |R|$  and STOPPINGCRITERIANOTMET() do
8:      $R_{keep} \leftarrow \text{None}$ 
9:      $\Omega^{keep} \leftarrow +\infty$ 
10:    for all  $\{R_j \in R \mid R_j \notin Q\}$  do
11:       $p_j^{greedy} \leftarrow p_j$ 
12:       $\Omega^{greedy} \leftarrow \text{EVALUATE}(X, R, \ell, p^{greedy}, a, \mathcal{B}, BD, \lambda)$ 
13:      if  $\Omega_{loss}^{greedy} < \Omega_{loss}^{keep}$  then
14:         $R_{keep} \leftarrow R_j$ 
15:         $\Omega^{keep} \leftarrow \Omega^{greedy}$ 
16:         $p^{keep} \leftarrow p^{greedy}$ 
17:         $p_j^{greedy} \leftarrow -1$ 
18:     $Q.ADD(R_{keep})$ 
19:    if  $\Omega_{loss}^{keep} < \Omega_{loss}^{best}$  then
20:       $\Omega^{best} \leftarrow \Omega_{loss}^{keep}$ 
21:       $p^{best} \leftarrow p^{keep}$ 
22:    if  $bt$  is true and ISBACKTRACKINGTIME() then
23:      run greedy contraction to remove  $l$  rules,  $l < |Q|$ 
24:  return  $(p^{best}, \Omega^{best})$ 

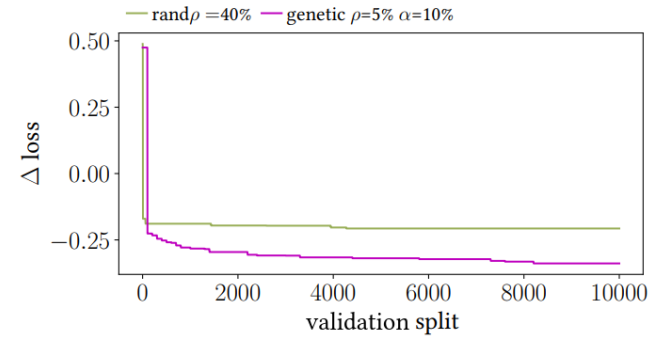
```

**Figure S1: Baseline metrics comparison in D2.****Figure S2: Greedy expansion results in D2.****Algorithm S4** Genetic programming optimization.

```

 $\theta$ : {Population size  $\psi$ , survivors fraction  $\alpha$ , mutation probability  $\rho$ }
1: function GENETIC.OPTIMIZE( $X, R, \ell, p, a, BD, \lambda, \Omega^1, \theta$ )
2:    $p^{best} \leftarrow p$ 
3:    $\Omega^{best} \leftarrow \Omega^1$ 
4:    $P \leftarrow \text{GENERATEINITIALPOPULATION}(R, p, \psi, \rho)$ 
5:   while STOPPINGCRITERIANOTMET() do
6:      $(P^{\pm}, P^{\mp}) \leftarrow \text{EVALUATEPOPULATION}(P, \alpha)$ 
7:      $P^+ \leftarrow \text{MUTATEANDCROSSOVER}(P^{\pm}, \alpha, \psi, \rho)$ 
8:      $P \leftarrow \{P^{\pm}, P^+\}$ 
9:      $(P^{\pm}, P^{\mp}) \leftarrow \text{EVALUATEPOPULATION}(P)$ 
10:     $p^{best} \leftarrow P_1^{\pm}$ 
11:     $\Omega^{best} \leftarrow \text{EVALUATE}(X, R, \ell, p^{best}, a, \mathcal{B}, BD, \lambda)$ 
12:  return  $(p^{best}, \Omega^{best})$ 
13: function GENERATEINITIALPOPULATION( $R, p, \psi, \rho$ )
14:   $P \leftarrow \emptyset$ 
15:  for  $i \in [0, \psi[$  do
16:     $p' \leftarrow p$ 
17:    for all  $p'_j \in p'$  do
18:      with  $\rho\%$  probability, do:
19:         $p'_j \leftarrow -1$ 
20:     $P[i] \leftarrow p'$ 
21:  return  $P$ 
22: function MUTATEANDCROSSOVER( $P^{\pm}, \alpha, \psi, \rho$ )
23:   $P^+ \leftarrow \emptyset$ 
24:  for  $i \in [0, (1 - \alpha) * \psi[$  do
25:     $p^{mother} \leftarrow \text{GETRANDOMVECTOR}(P^{\pm})$ 
26:     $p^{father} \leftarrow \text{GETRANDOMVECTOR}(P^{\mp})$ 
27:     $p^{child} \leftarrow p^{mother}$ 
28:    for all  $p_j^{child} \in p^{child}$  do
29:      with 50% probability, do:
30:         $p_j^{child} \leftarrow p_j^{father}$ 
31:    for all  $p_j^{child} \in p^{child}$  do
32:      with  $\rho\%$  probability, do:
33:         $p_j^{child} \leftarrow \text{RANDOMPRIORITYSHUFFLE}(p_i, a)$ 
34:     $P^+.ADD(p^{child})$ 
35:  return  $P^+$ 

```

**Figure S3: Genetic programming loss versus random search by number of evaluations in fold A of D1 (zoomed in the first 10000 rule evaluations; the methods nearly converge eventually)**