



Cost-sensitive reinforcement learning for credit risk

Jorge C-Rella ^{a,c}, David Martínez Rego ^b, Juan M. Vilar ^a

^a Department of Mathematics, Research Group MODES, CITIC, University of A Coruña, A Coruña, Spain

^b DataSpartan, London, United Kingdom

^c Department of Risks, ABANCA Financial Services, A Coruña, Spain

ARTICLE INFO

Keywords:

Cost-sensitive classification
Reinforcement learning
Bandit algorithms
Online learning
Credit risk
Decision making

ABSTRACT

Credit risk problems are dynamic because customer behavior is not stable, and they are cost-sensitive because the impact of a decision depends on the amount of the loan. Online learning algorithms, which evolve as more information becomes available, are an appropriate tool to study these dynamic problems. However, only information on approved transactions is available, which can lead to unfair biases and opportunity costs. Within reinforcement learning, bandit algorithms address this by balancing exploitation (acting according to the current model) and exploration (considering an action with limited information to improve predictions). The only remaining gap is to address the problem taking into account the different classification costs. This paper introduces cost-sensitive reinforcement learning algorithms to solve the credit risk problem from a dynamic perspective maximizing long-term benefits, proposing a cost-sensitive passive-aggressive algorithm and a cost-sensitive logistic bandit. Experiments on benchmark datasets and extensive simulation studies demonstrate the effectiveness and efficiency of the proposed algorithms.

1. Introduction

Lenders and investors face the challenge of selecting potential borrowers to whom they should extend credit in order to control risk while maximizing profit (Bahnsen et al., 2014; Lessmann et al., 2015; Lu et al., 2013). Traditional credit scoring systems assign scores based on socio-economic characteristics, with higher scores indicating a higher risk of default. Comprehensive experimental benchmarking studies have recently been conducted by Lessmann et al. (2015) and Vanderschueren et al. (2022), showing that the industry standard, logistic regression, is significantly outperformed by more advanced methods, which also have limitations. Their reliance on historical data results in a lack of information about potential borrowers who were previously deemed ineligible for credit. This can lead to customers being unfairly misjudged, opportunity costs and dangerous errors when the underlying relationships between variables change (concept drift). In addition, credit risk leads to an exploration vs. exploitation dilemma (Sutton & Barto, 2018): should a new customer be employed to gain more information and improve predictions at the potential cost of incurring losses, or should the current model's assessment dictate the decision? As a result, financial institutions face two challenges that are often overlooked: adapting to dynamic environments while optimizing both immediate risk and long-term benefit.

Online learning (OL) addresses the first challenge of dealing with dynamic environments by continuously updating models as new information becomes available. The fundamental difference between OL and classical offline methods lies in their approach to data handling. Static learning algorithms, such as logistic regression or decision trees, use a fixed set of data for model training, achieving high accuracy in offline settings. However, they cannot adapt to evolving patterns. OL processes data incrementally, learning continuously as new observations become available. As a result, models remain relevant and effective over time, even as data distributions and patterns evolve (Hoi et al., 2018; Lattimore & Szepesvári, 2020; Sutton & Barto, 2018). This continuous learning loop also leads to superior resource efficiency with minimal retraining requirements as noted in Li et al. (2010). The most prominent algorithms in the OL literature are the Online Gradient Descent algorithm (Zhao et al., 2015), the Passive-Aggressive algorithm (Crammer et al., 2006), the perceptron (Rosenblatt, 1958), and many other proposals that follow the principle of large-margin learning (Crammer et al., 2006; Dredze et al., 2008; Kivinen et al., 2004; Shalev-Shwartz, 2012; Wang et al., 2012).

OL deals with dynamic settings but does not consider exploration. Within the reinforcement learning paradigm, bandit algorithms are designed to balance exploration and exploitation. They assume that multiple actions with different rewards are available. Then, an agent

* Corresponding author at: Department of Mathematics, Research Group MODES, CITIC, University of A Coruña, A Coruña, Spain.
E-mail addresses: jorge.crella@udc.es (J. C-Rella), david@dataspartan.com (D. Martínez Rego), juan.vilar@udc.es (J.M. Vilar).

explores the different actions to decide on the optimal strategy that maximizes its reward in the long run. Due to the increasing number of problems in this context, several bandit algorithms have been developed in recent years, such as Epsilon-Greedy algorithms (Lattimore & Szepesvári, 2020; Sutton & Barto, 2018), Upper Confidence Bound (UCB) algorithms (Auer et al., 2002; Kaufmann et al., 2014) and sampling approaches (Kaufmann et al., 2012; Korda et al., 2013; Lattimore & Szepesvári, 2020; Russo et al., 2020; Thompson, 1933). The drawback of these algorithms is that they do not take into account any information other than the expected reward of each action. In practice, however, auxiliary information is usually available. For example, in credit risk, do you lend money to people regardless of their socioeconomic characteristics? Contextual bandits include the additional context of each action to make informed decisions. This approach has gained recognition in applications such as online advertising (Tang et al., 2013), recommender systems (Lai & Robbins, 1985; Li et al., 2010), personalized medicine (Rabbi et al., 2015), and credit risk (Shen et al., 2015; Visantavarakul, 2022), where tailoring decisions to specific situations leads to improved performance.

Most classification models, both online and offline, assume an uniform impact of misclassification errors. Nevertheless, this assumption can lead to suboptimal results in cost-sensitive (CS) scenarios where error costs vary across error types and instances. In credit risk, a borrower who defaults on a loan causes a greater loss than rejecting a good loan application, depending on the amount requested. By explicitly training models to adapt to different misclassification costs, they are aligned with the goal of maximizing benefits, leading to better results as shown in Bahnsen et al. (2014, 2013), C-Rella et al. (2025), Elkan (2001), Höppner et al. (2021), Vanderschueren et al. (2022) and Wang et al. (2012). Nevertheless, there is a lack of approaches for optimizing dynamic cost-sensitive problems as claimed in Achab et al. (2018) and Zhang et al. (2018).

In this paper, two novel algorithms are proposed to tackle the CS problem from a reinforcement learning perspective. An OL algorithm is proposed, extending the passive-aggressive (PA) algorithm introduced in Grammer et al. (2006). With this proposal, an updated model is obtained with all the information available to predict the optimal CS decision. In order to consider exploration, a CS version of the logistic bandit is also proposed, training a logistic model to maximize a reward function while considering exploration beyond the model's predictions. In this way, the risk of falling into bias or losing opportunities due to past decisions is minimized.

By combining a CS approach with the adaptability of OL and bandit algorithms, better results are expected with respect to current classification algorithms. This is empirically demonstrated with an extensive set of experiments on some real-world benchmark datasets and several simulations. Although exploration generally improves long-term performance, some companies may be more conservative or operate under stricter regulations that limit their ability to explore beyond the predictions of the current model. Therefore, both approaches are of practical interest. Throughout the paper the credit risk problem is considered to make the problem easier to understand, but note that all the proposals can be extended to any other CS problem by adapting the prediction costs specification.

The paper is organized as follows. Section 2 presents the CS classification problem to be solved. Section 3 introduces the online learning framework and the cost-sensitive passive-aggressive algorithm proposed in this paper. Section 4 presents state-of-the-art bandit algorithms and introduces the novel proposed CS logistic bandit. Section 5 summarizes the results of the simulation study and five real data sets. Conclusions and future extensions are included in Section 6.

2. Cost-sensitive classification

In cost-sensitive (CS) classification, the goal is to predict a binary response variable $Y \in \{0, 1\}$ (0 indicating a good operation and 1

Table 1

Instance-dependent cost matrix summarizing the benefit/cost of a true negative (TN), false positive (FP), false negative (FN) and a true positive (TP).

$\hat{y}_i = 0$	$\hat{y}_i = 1$
$y_i = 0$	$C_i^{TN} = aw_i$
$y_i = 1$	$C_i^{FP} = -cw_i - b$ $C_i^{FN} = -aw_i$ $C_i^{TP} = -b$

default for the credit risk problem) from a set of features (e.g. income, credit history) summarized in a covariate vector $\mathbf{X} = (X_1, \dots, X_d)$ taking into account prediction error costs. These costs depend on an exogenous variable W (the loan amount in credit risk), which is assumed independent of \mathbf{X} and Y according to Zadrozny and Elkan (2001). The objective is to maximize benefits.

Classical classification metrics are not valid in CS settings, as there is no direct relationship between accuracy-oriented metrics and CS metrics, as noted in C-Rella et al. (2024, 2025). Extending the proposals in Bahnsen et al. (2014, 2013), Elkan (2001) and Höppner et al. (2021), an instance-dependent *reward function* is defined from the cost matrix in Table 1 below. Considering the problem of credit risk, it is assumed that approving a good loan application implies a gain of aW , where a is the interest rate after deducting the operational costs. The impact of false negatives (FN) is given by cW , where c is the loss given default. Following the reasonableness condition introduced in Elkan (2001), it is assumed that $c > a$, i.e. that the impact of FN is greater than that of false positives (FP). Finally, there is a fixed cost b of labeling an observation as a case regardless of its true class, due to administrative costs. Negative values represent losses, positive values represent profits and it is assumed that $a, b, c > 0$.

Despite the cost matrix in Table 1 is very intuitive, we need a metric to evaluate model performance. By aggregating all the costs in Table 1 into a function, the reward function is defined as:

$$r(\hat{y}_i, w_i, y_i) = (1 - y_i)(1 - \hat{y}_i)C_i^{TN} + (1 - y_i)\hat{y}_iC_i^{FP} + y_i(1 - \hat{y}_i)C_i^{FN} + y_i\hat{y}_iC_i^{TP}. \quad (1)$$

The reward function in (1) is an individual metric. However, the goal is to maximize benefits over the entire population, for which model performance is evaluated considering the *mean reward*:

$$\mathcal{R} = \frac{1}{n} \sum_{i=1}^n r(\hat{y}_i, w_i, y_i), \quad (2)$$

which serves as an estimate of the expected benefit per operation and facilitates comparison between different models as it is a standardized metric. This is the function that is intended to be optimized throughout this work.

3. Online learning

Online learning (OL) algorithms incrementally learn a predictive model by sequentially observing data, making a prediction, receiving feedback, and adapting the model with the newly obtained information. This paper focuses on the Passive-Aggressive (PA) algorithm, a maximum margin classifier as introduced in Grammer et al. (2006), Shalev-Shwartz (2012) and Wang et al. (2012). This is motivated by the fact that it is generally a more effective approach than the Perceptron algorithm (see Zhao et al., 2015). Regarding previous CS OL approaches, the main proposals in the literature, introduced below in this section, are designed to deal with class unbalance considering constant entries in the cost matrix in Table 1. This is an oversimplification that may lead to suboptimal results, for which an instance-dependent cost-sensitive (IDCS) version of the PA algorithm is introduced to overcome the shortcomings of previous algorithms.

Throughout this section, for convenience, the binary response variable is represented as $Y \in \{-1, 1\}$, where 1 represents cases (default) and -1 non-cases (good operations). We restrict our discussion to classification functions of the form $\hat{Y} = sign(\theta' \mathbf{X})$, where $sign(z) = -1$

if $z < 0$ and $\text{sign}(z) = 1$ if $z \geq 0$, and $\theta \in \mathbb{R}^d$ is the parameter vector to be fitted.

At round t , the learner is given an observation \mathbf{x}_t of the covariate vector \mathbf{X} and makes a prediction $\hat{y}_t = \text{sign}(\theta' \mathbf{x}_t)$. The vector $\theta \in \mathbb{R}^d$ is estimated considering all the information available up to step t , denoted as $H_t = \{(\mathbf{x}_i, y_i, w_i)\}_{i=1}^{t-1}$. After the prediction \hat{y}_t is made, the correct answer $y_t \in \{-1, 1\}$ is revealed and the learner suffers a loss $\ell(\theta' \mathbf{x}_t, y_t)$, which measures the discrepancy between his prediction and the correct label (Crammer et al., 2006; Shalev-Shwartz, 2012). The algorithm uses the newly obtained instance-label pair (\mathbf{x}_t, y_t) to update the parameter vector θ for the next rounds, and the process is repeated. The goal of the learner is to minimize the cumulative loss along its path up to the horizon T , $\sum_{t=1}^T \ell(\hat{y}_t, w_t, y_t)$.

3.1. Passive-aggressive algorithm

The parameter θ in the linear classifier $\text{sign}(\theta' \mathbf{X})$ can be interpreted as the normal vector of an hyperplane dividing the data space in two half-spaces, where each side of the hyperplane corresponds to a class $\hat{Y} \in \{-1, 1\}$. The PA algorithm iteratively adjusts the parameter vector θ by maximizing the distance between the prediction for the last instance received, $\theta' \mathbf{x}_t$, and the hyperplane $\theta' \mathbf{X} = 0$. The goal is to get the maximum distance while maintaining proximity to the previous classifier, thereby ensuring a balance between adapting to new data and maintaining stability.

The value $|\theta' \mathbf{x}_t|$ represents the degree of confidence in the prediction (Crammer et al., 2006). The loss $\ell(\theta' \mathbf{x}_t, y_t)$ is proportional to the degree to which the prediction was wrong. If a prediction is correct, $y_t \hat{y}_t$ is positive. Consequently, the goal is to find $\theta \in \mathbb{R}^d$ such that $y_t \theta' \mathbf{x}_t > 0$. However, the focus remains on obtaining a classifier capable of predicting with high confidence. By imposing a margin on the classification, the following hinge loss function is obtained:

$$\ell_H(\theta; (\mathbf{x}_t, y_t)) = \begin{cases} 0, & y_t \theta' \mathbf{x}_t \geq 1 \\ 1 - y_t \theta' \mathbf{x}_t, & \text{otherwise.} \end{cases} \quad (3)$$

Considering the hinge loss function in Eq. (3), the parameter of the PA algorithm is actualized after each step t as the solution to the constrained optimization problem:

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{2} \|\theta - \theta_t\|^2 \text{ s.t. } \ell_H(\theta; (\mathbf{x}_t, y_t)) = 0, \quad (4)$$

where $\|\cdot\|$ represents the Euclidean L^2 norm. This notation is considered throughout the work for ease of reading.

If the prediction with the current parameter θ_t is correct, the hinge loss in (3) is zero and the algorithm remains passive, i.e. it keeps θ_t unchanged. If the prediction is wrong, $\ell_H(\theta_t; (\mathbf{x}_t, y_t)) > 0$ and the algorithm corrects θ_t in order to classify with margin, but changing the previous parameter θ_t as little as possible. This balances learning from new data while retaining knowledge from previous observations. The solution to the optimization problem in Eq. (4) is obtained in Crammer et al. (2006) as:

$$\theta_{t+1} = \theta_t + \tau_t y_t \mathbf{x}_t, \text{ where } \tau_t = \frac{\ell_H(\theta_t; (\mathbf{x}_t, y_t))}{\|\mathbf{x}_t\|^2}. \quad (5)$$

When an observation is misclassified, the update strategy in Eq. (5) updates the current estimate as necessary to assign the correct class to the last observation. Although reasonable, this can be too aggressive and is not robust to outliers. Therefore, following Vapnik (1998), it is considered a soft-margin approach introducing a slack variable ξ into the optimization problem in (4):

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{2} \|\theta - \theta_t\|^2 + C \xi \text{ s.t. } \ell_H(\theta; (\mathbf{x}_t, y_t)) \leq \xi \text{ and } \xi \geq 0. \quad (6)$$

where C is a positive constant that controls the effect of the slack variable. The updating process, outlined in Algorithm PA, shares the closed form of (5), but with a different τ_t term as demonstrated in Crammer

et al. (2006):

$$\theta_{t+1} = \theta_t + \tau_t y_t \mathbf{x}_t, \text{ where } \tau_t = \min \left\{ C, \frac{\ell_H(\theta_t; (\mathbf{x}_t, y_t))}{\|\mathbf{x}_t\|^2} \right\}. \quad (7)$$

Algorithm PA (Soft-margin) Passive-aggressive online learning.

```

1: Input Parameter  $C > 0$ 
2: Initialize  $\theta_1 = (0, \dots, 0)$ 
3: for  $t = 1, \dots, T$  do
4:   Receive covariate,  $\mathbf{x}_t$ 
5:   Compute prediction,  $\hat{y}_t = \text{sign}(\theta'_t \mathbf{x}_t)$ 
6:   Receive feedback,  $y_t$ 
7:   Compute  $\ell_H(\theta_t; (\mathbf{x}_t, y_t))$  as defined in Equation (3)
8:   Compute  $\tau_t = \min \left\{ C, \frac{\ell_H(\theta_t; (\mathbf{x}_t, y_t))}{\|\mathbf{x}_t\|^2} \right\}$ 
9:   Update  $\theta_{t+1} = \theta_t + \tau_t y_t \mathbf{x}_t$ 
10: end for
11: Output Parameter vector of the PA classifier,  $\theta_T$ 
```

The introduced Algorithm PA consists of a linear predictor, which may seem restrictive. However, it can be further extended to non-linear predictors using Mercer kernels as demonstrated in Vapnik (1998). Recall that the classifier can be represented as a sum of inner products:

$$\theta'_t \mathbf{x}_t = \sum_{i=1}^{t-1} \tau_i y_i (\mathbf{x}'_i \mathbf{x}_t),$$

and consider a Mercer kernel K such that $K(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})' \Phi(\mathbf{x}')$, where Φ is a “feature map” mapping \mathbf{x} to another space (possibly with higher dimensions). An explicit representation for Φ is not necessary as long as it defines an inner product. Replacing the inner products in Algorithm PA with $K(\mathbf{x}, \mathbf{x}')$ gives a nonlinear version of the PA algorithm, represented as:

$$\hat{y}_t = \text{sign}(\theta'_{t-1} \Phi(\mathbf{x}_t)), \text{ where } \theta'_{t-1} \Phi(\mathbf{x}_t) = \sum_{i=1}^{t-1} \tau_i y_i K(\mathbf{x}_i, \mathbf{x}_t). \quad (8)$$

3.2. State-of-the-art cost-sensitive online learning algorithms

The PA algorithm solves the problem of updating the estimation of θ as new information becomes available. However, it overlooks the impact of different misclassification error costs, which can lead to suboptimal results in cost-sensitive settings as noted in Bahnsen et al. (2014), Elkan (2001), Höppner et al. (2021) and Vanderschueren et al. (2022). Consequently, CS algorithms appear in the OL literature, which are introduced in this section.

3.2.1. Passive-aggressive prediction-based algorithm

The passive-aggressive prediction-based (PAPB) algorithm proposed in Crammer et al. (2006) extends Algorithm PA considering a classification margin defined by an asymmetric loss function. The PAPB algorithm considers the following optimization problem:

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{2} \|\theta - \theta_t\|^2 \text{ s.t. } y_t \theta' \mathbf{x}_t - \hat{y}_t \theta' \mathbf{x}_t \geq \sqrt{\ell_{PB}(\hat{y}_t, y_t)}, \quad (9)$$

where ℓ_{PB} is a loss function with different values depending on the error type (FN or FP). This corresponds to the constant version of the loss function in (13) below. Considering a soft-margin approach as in (6) in the optimization problem in (9) leads to the update strategy (Crammer et al., 2006):

$$\theta_{t+1} = \theta_t + \tau_t (y_t - \hat{y}_t) \mathbf{x}_t, \quad \tau_t = \min \left\{ C, \frac{(\hat{y}_t - y_t) \theta'_t \mathbf{x}_t + \sqrt{\ell_{PB}(\hat{y}_t, y_t)}}{\|(\hat{y}_t - y_t) \mathbf{x}_t\|^2} \right\}. \quad (10)$$

3.2.2. Cost-sensitive online gradient-descent algorithm

The cost-sensitive online gradient descent (CSOGD) algorithm, studied in Wang et al. (2012) and Zhao et al. (2015), extends the on-

line gradient descent approach introduced in Hazan et al. (2007) by considering misclassification costs in the parameter updating process:

$$\theta_{t+1} = \begin{cases} \theta_t + \lambda \rho_t y_t \mathbf{x}_t, & \ell_{CH}(\theta_t; (\mathbf{x}_t, y_t)) > 0 \\ \theta_t, & \text{otherwise.} \end{cases} \quad (11)$$

where I is the characteristic function, equals 1 when the expression between brackets is true and 0 otherwise, $\rho_t = \rho I(y_t = 1) + I(y_t = -1)$, $\rho = \frac{C^{FN}}{C^{FP}}$ is the ratio between the impact of a FN error and a FP error, $\lambda > 0$ is a learning rate parameter and $\ell_{CH}(\theta_t; (\mathbf{x}_t, y_t))$ is the cost-sensitive version of the hinge loss function in (3):

$$\ell_{CH}(\theta_t; (\mathbf{x}_t, y_t)) = (\rho I(y_t = 1) + I(y_t = -1)) \max \{0, 1 - y_t \theta_t' \mathbf{x}_t\}$$

3.2.3. Perceptron algorithm with uneven margin

The perceptron algorithm with uneven margin (PAUM) proposed in Li et al. (2002) is the CS version of the perceptron algorithm (Rosenblatt, 1958). The predicted class for an observation \mathbf{x}_t is assigned as:

$$\hat{y}_t = \text{sign}(\theta_t' \mathbf{x}_t + b)$$

where $\theta \in \mathbb{R}^d$ and $b \in \mathbb{R}$. Considering different margins for the positive and negative classes, τ_{+1} and τ_{-1} respectively, $R = \max_t \|\mathbf{x}_t\|$ and a given learning rate η , starting with all variables equal to 0, the updating process of the PAUM algorithm is given by:

$$\theta_{t+1} = \begin{cases} \theta_t + \eta y_t \mathbf{x}_t, & y_t(\theta_t' \mathbf{x}_t + b_t) \leq \tau_{y_t} \\ \theta_t, & \text{otherwise.} \end{cases} \quad (12)$$

$$b_{t+1} = \begin{cases} b_t + \eta y_t R^2, & y_t(\theta_t' \mathbf{x}_t + b_t) \leq \tau_{y_t} \\ b_t, & \text{otherwise.} \end{cases}$$

3.3. Instance-dependent cost-sensitive passive-aggressive algorithm

Algorithm PA has shown good performance in some settings, such as online advertising, where the agent is only concerned with whether the customer clicks or not (Li et al., 2010). As for the CS approaches presented previously, they are designed to solve the class imbalance problem considering constant costs depending on the type of error (FP or FN). However, in settings where the misclassification cost depends on the characteristics of the instance, these assumptions are too restrictive and naive.

This section proposes a novel instance-dependent cost-sensitive passive-aggressive (IDCSPA) algorithm that extends the previously introduced PAPB algorithm to the IDCS setting. This approach retains the advantages of OL, but further incorporates the instance-dependent impact of misclassification errors. This is expected to result in a proficient model.

Since the PA algorithm, by definition, focuses only on losses, the reward function in Eq. (1) is adapted to consider only misclassification costs:

$$\ell_{CS} = I(y_i = -1)I(\hat{y}_i = 1)C_i^{FP} + I(y_i = 1)I(\hat{y}_i = -1)C_i^{FN}, \quad (13)$$

where C_i^{FP} and C_i^{FN} are the FP and FN costs respectively, as defined in Table 1. The loss function ℓ_{CS} in Eq. (13) is introduced into the PAPB optimization problem (9) to obtain the proposed IDCSPA algorithm. Thus, at each step t , the parameter θ_t is fitted by solving the optimization problem:

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{2} \|\theta - \theta_t\|^2 \quad \text{s.t. } y_t \theta' \mathbf{x}_t - \hat{y}_t \theta' \mathbf{x}_t \geq \sqrt{\ell_{CS}(\hat{y}_t, w_t, y_t)}, \quad (14)$$

where $\hat{y}_t = \text{sign}(\theta_t' \mathbf{x}_t)$ as for Algorithm PA. The optimization problem in (14) retains the idea of varying the current estimate as little as possible, while satisfying a certain margin, which in this case depends on the loss function ℓ_{CS} in (13). The optimization problem in Eq. (14) has a close form solution generalizing the solution developed in Crammer et al.

(2006):

$$\theta_{t+1} = \theta_t + \tau_t (y_t - \hat{y}_t) \mathbf{x}_t, \quad \text{where } \tau_t = \frac{(\hat{y}_t - y_t) \theta_t' \mathbf{x}_t + \sqrt{\ell_{CS}(\hat{y}_t, w_t, y_t)}}{\|(\hat{y}_t - y_t) \mathbf{x}_t\|^2}. \quad (15)$$

As for the PA algorithm in Equation (5), the update strategy in (15) may be too aggressive, especially since we include the loss function in the estimation. To consider a soft-margin approach, extending the proposal in Eq. (7), the IDCSPA updating rule is defined by:

$$\theta_{t+1} = \theta_t + \tau_t (y_t - \hat{y}_t) \mathbf{x}_t, \quad \tau_t = \min \left\{ C, \frac{(\hat{y}_t - y_t) \theta_t' \mathbf{x}_t + \sqrt{\ell_{CS}(\hat{y}_t, w_t, y_t)}}{\|(\hat{y}_t - y_t) \mathbf{x}_t\|^2} \right\}, \quad (16)$$

where $C > 0$ is the parameter that controls the aggressiveness of the algorithm. The resulting classifier is shown in Algorithm IDCSPA, which as for the PA algorithm, can be extended to non-linear predictors by replacing the inner products with a Mercer kernel, similar to what was presented in (8).

Algorithm IDCSPA Instance-dependent cost-sensitive (soft-margin) passive-aggressive online learning (IDCSPA)

```

1: Input Parameter  $C > 0$ 
2: Initialize  $\theta_1 = (0, \dots, 0)$ 
3: for  $t = 1, \dots, T$  do
4:   Receive covariate  $\mathbf{x}_t$ 
5:   Compute prediction  $\hat{y}_t = \text{sign}(\theta_t' \mathbf{x}_t)$ 
6:   Receive feedback,  $y_t$ 
7:   Compute  $\tau_t = \min \left\{ C, \frac{(\hat{y}_t - y_t) \theta_t' \mathbf{x}_t + \sqrt{\ell_{CS}(\hat{y}_t, w_t, y_t)}}{\|(\hat{y}_t - y_t) \mathbf{x}_t\|^2} \right\}$ , with  $\ell_{CS}$  as in (13)
8:   Update  $\theta_{t+1} = \theta_t + \tau_t (y_t - \hat{y}_t) \mathbf{x}_t$ 
9: end for
10: Output Vector classifier  $\theta_T$ 
```

4. Bandit algorithms

The OL framework aims to optimize exploitation with all the available information, but does not consider exploration beyond the predictions of the model. However, many real-world problems as credit risk (Achab et al., 2018; Visantavarakul, 2022) or online advertising (Li et al., 2010) require choosing among multiple options and learning the best one through successive trial and error. Bandit algorithms are a class of reinforcement learning algorithms designed to balance the process of seeking new information (exploration) with maximizing rewards based on existing knowledge (exploitation) in sequential decision problems. Since exploration and exploitation cannot be done simultaneously, in Sutton and Barto (2018) it is concluded that there is a trade-off between these two goals. Thus, at each step, the agent must decide whether to exploit the known option with the highest expected reward or to explore a less known option, potentially discovering better strategies for future gains.

Bandit algorithms receive their name from a very illustrative example represented in Fig. 1, where there are multiple bandit machines with different (unknown) reward distributions (Achab et al., 2018). A gambler (agent) stares at this series of one-armed bandits (actions), each promising unknown payouts. The agent is faced with the dilemma of pulling the familiar lever that is currently spitting out coins (exploitation) or trying a new one for potentially bigger wins (exploration).

The agent's goal is to learn the reward distribution of each arm and find a strategy (or policy) in order to maximize the cumulative reward over time. A naive approach would be to try all arms once and then choose the one with the highest reward. Nevertheless, this can be misleading, as the initial result could be pure luck. On the other

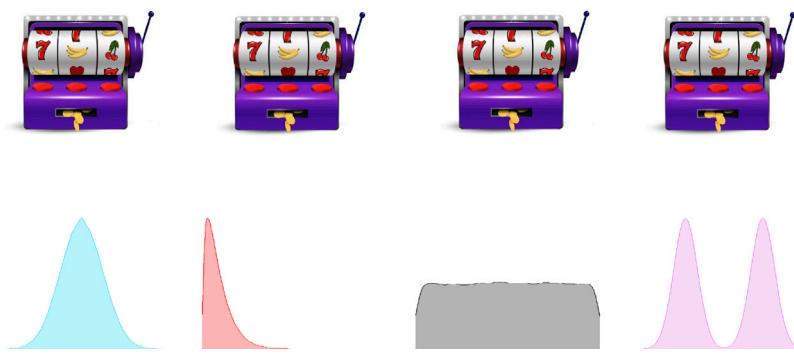


Fig. 1. Bandit machines with different reward distributions (bottom graphs).

hand, endless exploration will lead to missing the benefits of the best option. Bandit algorithms address this exploitation-exploration trade-off by strategically choosing between exploring less known options and exploiting the seemingly best one. Considering that instead of multiple slot machines there is a single one with multiple arms to pull, the problem is usually known as the multi-arm bandit (MAB) problem.

In this section, state-of-the-art MAB algorithms are presented along their drawbacks, motivating the novel CS logistic bandit algorithm proposed in this paper. There are three main bandit algorithms in the literature. The Epsilon-Greedy algorithm (see Sutton & Barto, 2018), a greedy algorithm allowing for exploration, is an intuitive way of forcing exploration. At each step, it exploits the most lucrative arm with probability $1 - \epsilon$, and explores a random new arm with probability ϵ . However, this algorithm does not explore the environment efficiently, since it selects new actions with equal probability, regardless of whether it is actually a good action to explore. Two approaches are proposed to overcome this drawback. The first is the Upper Confidence Bound (UCB) algorithm (Auer et al., 2002; Kaufmann et al., 2014; Lattimore & Szepesvári, 2020), which prioritizes regions with high potential rewards while accounting for uncertainty through confidence bounds. Based on optimism in the face of uncertainty, the action with the highest upper confidence bound for the expected reward is selected. In this way, UCB exploits the best arm while exploring those with promising rewards, minimizing the risk of overlooking hidden riches. The other approach is Thompson Sampling (TS), introduced in Thompson (1933) and further studied in Kaufmann et al. (2012), Korda et al. (2013) and Russo et al. (2020). It embraces the probabilistic nature of the problem by treating models as probabilistic terms rather than deterministic truths. TS samples model parameters from their posterior distribution, updated with the observations collected up to time t . This uncertainty-aware approach allows TS to adapt to dynamic scenarios while exploring the arms for which there is less information.

4.1. Contextual bandits

The previously introduced MABs blindly explore actions with unknown rewards. Nevertheless, real-world scenarios often provide additional information about each option. For example, the socioeconomic characteristics of the customer in credit risk or the search history in advertisement recommendation. This additional knowledge is called *context*, and it can change the optimal choice as the best action may differ in different contexts (Langford & Zhang, 2007; Panedy et al., 2007; Wang et al., 2005; Zhou, 2016).

Contextual multi-armed bandits (CMAB) take this context into account, tailoring predictions to each unique situation. Thus, improving performance and allowing the generalization of the information across actions (Russo et al., 2020). As a result, CMABs have found applications in many different areas, such as recommendation engines (Lai & Robbins, 1985; Li et al., 2010; Tang et al., 2013), (personalized)

healthcare (Rabbi et al., 2015), and portfolio selection (Shen et al., 2015), among others.

A CMAB, B , is defined by a set of available arms $\mathcal{A} \subseteq \{1, \dots, K\}$. The set of arms \mathcal{A} is assumed to be fixed at each step $t = 1, \dots, T$. Each arm $k \in \mathcal{A}$ is described by a reward function r_k mapping d -dimensional context vectors x_t to some reward $r_{t,k} = r_k(x_t)$ (Auer et al., 2002; Langford & Zhang, 2007). A policy π observes the current state of the world represented by the context feature x_t . Using some arm-selection strategy based on the previous information, $H_t = \{(x_i, r_{i,a_i})\}_{i=1}^{t-1}$, the policy chooses one of the available arms, defined as taking action $a_t \in \mathcal{A}$. Then, the policy gets the reward of the chosen arm, r_{t,a_t} . With the new observation (x_t, r_{t,a_t}) the policy updates its strategy. The goal of the policy π is to maximize its cumulative reward $\sum_{i=1}^T r_{i,a_i}$. This is the unscaled version of the expected reward in (2).

One approach when having different contexts might be to apply a non-contextual MAB algorithm to each context. However, the relationship between the contexts is lost, so learning one context does not help learning the others as noted in Zhou (2016). In addition, this approach assumes that the contexts can be enumerated, which is not true when they are continuous. The most extended approach to solving CMAB problems is to consider stochastic bandit algorithms, which assume that the reward of each arm follows an unknown probability distribution depending on the context. The main state-of-the-art proposals are presented in the following subsections.

4.1.1. Contextual Epoch-Greedy

The Contextual Epoch-Greedy (CEG) algorithm studied in Langford and Zhang (2007) considers two well differentiated steps of exploration and exploitation. During exploration, it gathers information by randomly selecting each of the available arms to obtain information about their rewards along the contextual variables. The goal is to create unbiased samples by randomly pulling all arms to improve the accuracy of learning. This step leads to large immediate regret, but can potentially be reduced for the exploitation step. With the information obtained, a model $s : X \mapsto \mathcal{A}$ is estimated in a class S of models to select the optimal action a_t given a context x_t .

Once a satisfactory regret is obtained with the estimated model, the algorithm switches to the exploitation phase up to the horizon T . For a Bernoulli bandit, the expected regret in the exploration step can be as large as $O(1)$ when all chosen actions have a reward of 0 and the optimal arm has a reward of 1. In the exploitation step, the regret is expected to be much smaller. Considering n exploration steps, where the average regret is e_n , the total regret is bounded by $n + (T - n)e_n$. Consequently, the optimal strategy is to switch from exploration to exploitation at the step n that minimizes the previous expression, as demonstrated in Langford and Zhang (2007).

4.1.2. LinUCB

The Linear Upper Confidence Bound (LinUCB) algorithm, as proposed in Chu et al. (2011) and Li et al. (2010), extends the UCB algorithm introduced in Auer et al. (2002) to accommodate to context information. It assumes that the expected reward of an arm k is linear with respect to the feature vector \mathbf{x}_t through an unknown coefficient vector θ_k :

$$E[r_{t,k} | \mathbf{x}_t] = \theta_k' \mathbf{x}_t \quad (17)$$

In Li et al. (2010) it is shown that, under the linear assumption in (17), a confidence interval for the expected reward can be efficiently computed. Considering ridge regression, the parameter vectors θ_k are estimated as:

$$\hat{\theta}_k = (\mathbf{D}'_k \mathbf{D}_k + I_d)^{-1} \mathbf{D}'_k \mathbf{r}_k$$

where $\mathbf{D}_k \in \mathcal{M}^{(t-1) \times d}$ is the design matrix whose rows correspond to the observed instances $\{\mathbf{x}_i\}_{i=1}^{t-1}$ for the arm k , \mathbf{r}_k the analog for the rewards, and I_d is the $d \times d$ identity matrix. Then, as shown in Li et al. (2010):

$$\mathbb{P}\left(|\hat{\theta}_k \mathbf{x}_t - E[r_{t,k} | \mathbf{x}_t]| \leq \alpha \sqrt{\mathbf{x}_t' (\mathbf{D}'_k \mathbf{D}_k + I_d)^{-1} \mathbf{x}_t}\right) \geq 1 - \delta \quad (18)$$

for any $\delta > 0$, where $\alpha = 1 + \sqrt{\ln(2/\delta)/2}$. The above inequality gives a bound on the expected reward of an arm k . Extending the UCB logic of optimism in the face of uncertainty, the policy in LinUCB selects the arm with the largest expected reward bound (18) as shown in Algorithm LinUCB.

Algorithm LinUCB LinUCB algorithm

```

1: Input Confidence parameter  $\delta \in [0, 1]$ 
2: Initialize  $\mathbf{A}_a \leftarrow I_d$ 
3: Initialize  $\mathbf{b}_a \leftarrow \mathbf{0}_d$ 
4: for  $t = 1, \dots, T$  do
5:   for  $a \in A_t$  do
6:     Compute  $\hat{\theta}_a \leftarrow \mathbf{A}_a^{-1} \mathbf{b}_a$ 
7:     Compute  $r_{t,a} \leftarrow \hat{\theta}_a' \mathbf{x}_t + (1 + \sqrt{\ln(2/\delta)/2}) \sqrt{\mathbf{x}_t' (\mathbf{x}_t \mathbf{x}_t' + I_d)^{-1} \mathbf{x}_t}$ 
8:   end for
9:   Compute  $a_t = \operatorname{argmax}_{a \in A_t} r_{t,a}$ 
10:  Receive Reward  $r_{a_t}$ 
11:  Update  $\mathbf{A}_a \leftarrow \mathbf{A}_a + \mathbf{x}_t \mathbf{x}_t'$ 
12:  Update  $\mathbf{b}_a \leftarrow \mathbf{b}_a + r_{a_t} \mathbf{x}_t$ 
13: end for
14: Output Reward parameters  $\hat{\theta}_k$ 

```

4.1.3. Logistic bandits

In the Bernoulli bandit problem, the agent chooses an action a_t and receives a binary feedback $y_t \in \{0, 1\}$ Bernoulli distributed. The logistic bandit (LB) model extends Bernoulli bandits by incorporating generalization across actions through the context \mathbf{x}_t . As in logistic regression, the conditional probability of the response variable is modeled with a logistic function through a linear combination of the context vector. For each arm k , the probability of receiving a reward is modeled as $P(Y = 1 | \mathbf{X} = \mathbf{x}) = s(\mathbf{x}; \theta_k) = (1 + \exp(-\theta_k' \mathbf{x}))^{-1}$, where $\theta_k \in \mathbb{R}^d$ is the parameter vector to be estimated.

At each step t , the parameter θ_k associated with each arm k is estimated as the maximizer of the likelihood function given the available information:

$$\hat{\theta}_{k,t} = \operatorname{argmax}_{\theta} \mathcal{L}_{k,t}^L(\theta) \quad (19)$$

where

$$\mathcal{L}_{k,t}^L(\theta) = \frac{1}{t} \sum_{i=1}^t I(a_i = k) \{y_i \ln[s(\mathbf{x}_i; \theta)] + (1 - y_i) \ln[1 - s(\mathbf{x}_i; \theta)]\} \quad (20)$$

is the likelihood function for an arm k and information up to the step t . The main difference with the classical logistic model is the exploration carried out by the logistic bandit and the continuous updating as more information is available. The exploration is imposed by considering Thompson Sampling (TS) as introduced in Thompson (1933). There are other proposals for allowing the agent to explore (Dumitrascu et al., 2018; Faury et al., 2020; Zhang et al., 2016). However, in Sutton and Barto (2018) TS is found to be more effective, for which further approaches are not considered.

In TS, at each step t , the parameter vector $\theta_{k,t}$ is treated as a random variable sampled from the posterior distribution of $\hat{\theta}_{k,t}$, estimated from the available information. The sampled vectors are used to estimate the reward probability on each arm k as $s(\mathbf{x}; \theta_{k,t})$, and the action that maximizes the expected reward is selected. Feedback on the selected action is received, and the new information is used to update the posterior distribution of $\hat{\theta}_{k,t}$ using Bayes' rule (Dumitrascu et al., 2018; Russo et al., 2020). For unexplored arms, the tails of the distribution are heavier, favoring exploration. As more information is obtained, the uncertainty is reduced, favoring the exploitation of the most profitable arms.

Fig. 2 shows an example of different steps of TS. The dotted line represents the estimator at each step and the curves represent the density over which it is resampled. In the early steps, there is a lot of uncertainty, favoring exploration. As more information becomes available, the density of the estimator becomes more concentrated, favoring exploitation.

The drawback of the logistic bandit algorithm is that it results in a computationally intractable posterior, which makes the implementation of TS challenging as noted in Dumitrascu et al. (2018). To address this, the Laplace approximation is used, which is based on a Taylor series expansion of the distribution function up to the second order. The Laplace approximation is expected to work well in the logistic framework, with a smooth density peaked around its mode (Dumitrascu et al., 2018; Gamerman & Lopes, 2006). There are other approaches to performing the sampling process, such as Metropolis Hasting or Langevin Monte Carlo Markov Chain. However, in Visantavarakul (2022) no significant improvement is observed with these approaches, for which no further algorithms are considered.

For the logistic bandit, given the estimator $\hat{\theta}_{k,t}$ of θ_k as in Eq. (19), considering the Laplace sampling approximation, the parameter vector used in the step $t+1$ is resampled from the distribution:

$$\theta_{k,t+1} \sim N_d(\hat{\theta}_{k,t}, H_{k,t}) \quad (21)$$

where $H_{k,t}$ is the Hessian of the likelihood function in (20) evaluated at $\hat{\theta}_{k,t}$. This process is analogous to the one represented in Fig. 2. Considering the Laplace sampling approximation in Eq. (21), the logistic bandit algorithm is illustrated in Algorithm LB.

Algorithm LB Logistic bandit algorithm

```

1: Initialize  $\theta_{k,0} \leftarrow \mathbf{0}_d$ 
2: Initialize  $H_0 \leftarrow I_d$ 
3: for  $t = 1, \dots, T$  do
4:   Compute  $\theta_{k,t} \sim N_d(\hat{\theta}_{k,t-1}, H_{k,t-1})$ ,  $k \in A_t$ 
5:   Compute  $a_t = \operatorname{argmax}_{a \in A} s(\theta_{a,t}; \mathbf{x})$ 
6:   Receive Result  $y_t$ 
7:   Update  $\hat{\theta}_{k,t} = \operatorname{argmax}_{\theta} \mathcal{L}_{k,t}^L(\theta)$ , where  $\mathcal{L}_{k,t}^L$  as defined in (20)
8:   Update  $H_{k,t} = H \mathcal{L}_{k,t}^L(\hat{\theta}_{k,t})$ 
9: end for
10: Output Parameters  $\hat{\theta}_{k,T}$ 

```

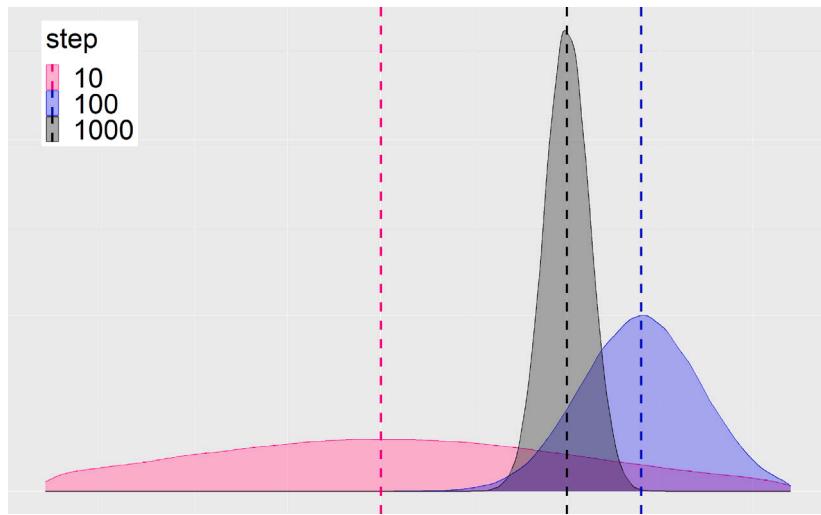


Fig. 2. Thompson sampling resampling densities for different steps.

4.2. Cost-sensitive logistic bandit

In this section, motivated by the success of the MAB policies recalled above, a cost-sensitive logistic bandit (CSLB) is proposed, extending the logistic bandit introduced in Algorithm LB, to solve the credit risk problem (a CS binary classification problem as introduced in Section 2).

The set of actions in credit risk is binary, $\hat{Y} \in \{0, 1\}$, for which it can be thought of as a Bernoulli one-armed bandit where the agent has to decide whether to pull the arm or not. If the operation is approved, the customer's behavior is observed; otherwise, no information is obtained. The goal is to maximize the total reward, which depends on the characteristics of each instance as defined by the reward function in Eq. (1).

Considering the logistic classifier, assume that $E[Y | X = x] = s(x; \theta) = (1 + \exp(-\theta'x))^{-1}$, where $\theta \in \mathbb{R}^d$. The goal is to estimate the parameter θ maximizing the expected reward $E[r(\hat{Y}, W, Y)]$. Assuming that $\hat{Y} = \mathbb{I}(s(X) > u)$, with $u \sim U[0, 1]$, then $E[\hat{Y} | X = x] = P(\hat{Y} = 1 | x) = s(x; \theta)$ as introduced in C-Rella et al. (2025). Thus, the average expected reward (AER) (Bahnsen et al., 2014; Höppner et al., 2021) of a classifier s over the set sample $\mathcal{Z}_t = \{\mathbf{z}_i = (\mathbf{x}_i, w_i, y_i)\}_{i=1}^t$ can be defined as:

$$\mathcal{L}_t^{CS}(\theta) = \frac{1}{t} E \left[\sum_{i=1}^t r(\hat{y}_i, w_i, y_i) \mid \mathcal{Z}_t \right] = \frac{1}{t} \sum_{i=1}^t r(s(\mathbf{x}_i; \theta), w_i, y_i) \quad (22)$$

where r is the reward function in (1). The AER in (22) corresponds to the empirical version of $E[r(\hat{Y}, W, Y)]$ given the sample $\mathcal{Z}_t = \{\mathbf{z}_i = (\mathbf{x}_i, w_i, y_i)\}_{i=1}^t$. Thus, the CS estimator of θ is defined as:

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}_t^{CS}(\theta) \quad (23)$$

The objective function in (22) is usually not convex nor has a unique minimum due to its dependence on W , Y and θ through \hat{Y} . To allow solving the optimization problem, a regularizer is added to the objective function to enforce convexity as suggested in Bahnsen et al. (2013), Höppner et al. (2021) and Stripling et al. (2018). This makes it easier for numerical optimization techniques to avoid suboptimal solutions and improves generalization performance. Therefore, it is considered the regularized AER:

$$\mathcal{L}_t^{\lambda}(\theta) = \frac{1}{t} \sum_{i=1}^t r(s(\mathbf{x}_i; \theta), w_i, y_i) + \lambda \|\theta\|_1 \quad (24)$$

where $\lambda \in \mathbb{R}^+$ is the parameter that controls the regularization. We consider a lasso regularization as proposed in C-Rella et al. (2025)

and Höppner et al. (2021), without loss of generality. The corresponding estimator is defined as:

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}_t^{\lambda}(\theta) \quad (25)$$

As the set of actions is assumed binary, the problem can be modeled with a single parameter θ such that $\hat{y}_t = I(s(\theta; \mathbf{x}_t) > h_t)$, where h_t is the decision threshold for the instance (\mathbf{x}_t, w_t) . For a reward function as in (1), considering the expected reward of an action,

$$R(\hat{y} | \mathbf{x}, w) = r(\hat{y}, w, 0)(1 - s(\mathbf{x})) + r(\hat{y}, w, 1)s(\mathbf{x})$$

the arm is pulled if the expected reward is greater than not pulling the arm, i.e. if $R(0 | \mathbf{x}, w) \geq R(1 | \mathbf{x}, w)$. This logic leads to the theoretically optimal decision threshold (Bahnsen et al., 2013; Elkan, 2001):

$$h_t = \frac{C_i^{FP} - C_i^{TN}}{C_i^{FP} - C_i^{TN} - C_i^{TP} + C_i^{FN}} \quad (26)$$

Introducing the regularized AER in (24) as the objective function in the logistic bandit in Algorithm LB and considering the decision rule in (26), the CS logistic bandit (CSLB) is represented in Algorithm CSLB. At each step, the CSLB selects the optimal action based on the estimator in (25) and the decision threshold in (26) considering all the information available. Similarly, by considering the TS sampling strategy, it efficiently explores all contexts to enhance long-term rewards. By combining exploration with a CS approach, the CSLB algorithm is expected to outperform previous models.

Algorithm CSLB Cost-sensitive logistic bandit algorithm

- 1: **Input** Reward function $r(\hat{y}_i, w_i, y_i)$ as in Equation (1)
 - 2: **Initialize** $\theta_0 \leftarrow \mathbf{0}_d$
 - 3: **Initialize** $H_0 \leftarrow I_d$
 - 4: **for** $t = 1, \dots, T$ **do**
 - 5: **Compute** $\theta_t \sim N_d(\hat{\theta}_{t-1}, H_{t-1})$
 - 6: **Receive** Context (\mathbf{x}_t, w_t)
 - 7: **Compute** $h_t = \frac{C_t^{FP} - C_t^{TN}}{C_t^{FP} - C_t^{TN} - C_t^{TP} + C_t^{FN}}$
 - 8: **Compute** $a_t = I(s(\theta_t; \mathbf{x}_t) \geq h_t)$
 - 9: **Receive** Result y_t
 - 10: **Update** $\hat{\theta}_t = \operatorname{argmax}_{\theta} \mathcal{L}_t^{\lambda}(\theta)$, where \mathcal{L}_t^{λ} as defined in (24)
 - 11: **Update** $H_t = H \mathcal{L}_t^{\lambda}(\hat{\theta}_t)$
 - 12: **end for**
 - 13: **Output** Reward parameter $\hat{\theta}_T$
-

5. Experimental results

The performance of the introduced OL and bandit algorithms are evaluated considering a wide range of simulated settings and the analysis of five real credit risk data sets. The objective is to understand how the different algorithms behave as the proportion of cases, the dependency relationship and the cost specification vary.

Online learning and bandit algorithms represent two different methodologies. Therefore, the algorithms are evaluated separately. After introducing the datasets, the results obtained for the different algorithms are presented, together with some comments on the performance of the classifiers. The optimization process for all the algorithms is implemented using the open source statistical software R ([R Core Team, 2021](#)).

Regarding OL algorithms, the passive-aggressive (PA) algorithm detailed in Algorithm [PA](#), the CSOGD algorithm in Eq. (11), the PAUM algorithm as depicted in Eq. (12), the PAPB algorithm described in (10) and the proposed Algorithm [IDCSPA](#) are considered. The selection of the aggressiveness parameter C in the PA algorithms is done by cross-validation (CV) over the grid $C \in (0, 0.001, 0.01, 0.1, 1, 10, 100)$. In addition, we evaluate a static approach (Static) considering a linear model, $\text{sign}(\theta' \mathbf{X})$, previously fitted to a static training sample minimizing the hinge loss in (3), in line with the classical approach in credit risk.

Regarding bandit algorithms, we consider the logistic bandit described in Algorithm [LB](#), the LinUCB algorithm summarized in Algorithm [LinUCB](#) and the CEG algorithm described in Section 4.1.1. We also consider a classical (Static) approach in the simulation study. This consists of training a logistic model on a training sample and then considering the fitted model on the new samples. All these algorithms are compared with the proposal in Algorithm [CSLB](#), with the regularization parameter λ estimated by cross-validation over the grid $\lambda \in (0, 0.001, 0.01, 0.1, 1, 2, 5, 10, 20)$.

By considering static approaches, the improvement obtained with reinforcement learning algorithms is evaluated with respect to a static approach. Similarly, by considering cost-insensitive models such as the PA algorithm or the logistic bandit, a comparison is also made of the improvement obtained by considering a CS approach. Finally, this study evaluates the behavior of the proposed algorithms and compares them with the state-of-the-art.

For the OL algorithms, for a set $\{(\hat{y}_i, y_i, w_i)\}_{i=1}^n$ of the predicted label, the case indicator and the exogenous variable, the savings ([Bahnsen et al., 2014; Höppner et al., 2021](#)) are considered to evaluate model performance:

$$\text{Savings} = 1 - \frac{\sum_{i=1}^n \ell(\hat{y}_i, w_i, y_i)}{\sum_{i=1}^n y_i w_i} \quad (27)$$

where ℓ is the loss function defined in Eq. (13). This metric provides a standardized measure of the losses prevented compared to the base case scenario, allowing for a reasonable comparison across different models and problems. The maximum value this metric can take is 1, in which case all observations are well classified. For the bandit algorithms, we evaluate model performance considering the AER introduced in Eq. (22).

Since a scaled version of the same objective function is obtained with the same a/c and b/c ratios in the cost matrix in [Table 1](#), the values $b = 10$ (the analysis cost) and $c = 0.75$ (the loss given default) are fixed and $a = \{0.05, 0.1, 0.2\}$ (the benefit per operation). This provides valuable insights into model performance under different false positive and false negative rates.

Model performance is evaluated using both cost-sensitive and classification metrics, taking into account the savings in (27) or the AER in (22) as appropriate, the positive predicted proportion (PP), the precision ($PPV = TP/PP$), the recall = $TP/(TP + FN)$ and the F score.

Table 2

Data generation models for the simulation study, where μ is the link function, $V = \theta_0' \mathbf{X}$ a linear combination of the covariates and ϵ an error term.

Model	Index	$\mu(V)$	ϵ
Logit	$V = \theta_0' \mathbf{X}$	$\frac{e^V}{1+e^V}$	$N(0, 1)$
Squared	$V = \theta_0' \mathbf{S}$	$\frac{e^V}{1+e^V}$	$N(0, 1)$
Qubic	$V = \theta_0' \mathbf{Q}/5$	$V(V - 2)(V + 2)$	$N(0, 1)$
Ker & Sam	$V = \theta_0' \mathbf{X}$	V	$.5N(3, 1) + .5N(-3, 1)$
Scobit	$V = \theta_0' \mathbf{X}$	$1 - \frac{1}{(1+e^V)^10}$	$N(0, 1)$
Klein & Spady	$V = \theta_0' \mathbf{X}$	V	$N(0, \sqrt{0.1}(1 + V^2))$

5.1. Simulation study

This section explores the behavior of the algorithms presented in a simulation study motivated by the credit risk problem introduced in Section 2. By controlling the data generation process we can observe and compare the strengths and weaknesses of each algorithm, providing valuable guidance for selecting the most appropriate method for specific real-world applications.

5.1.1. Simulation models

The data is simulated based on six different models, which are summarized in [Table 2](#). For all simulations, the underlying model is given by $P(Y = 1) = \mu(v + \epsilon)$, where μ is a link function and $V = \theta_0' \mathbf{X}$ is a linear combination of the predictor covariates $\mathbf{X} = (1, X_1, \dots, X_6) \in \mathbb{R}^7$. For the Logit model, $X_i \sim N(0, 1)$, $\mu(v) = \frac{e^v}{1+e^v}$ and $\epsilon \sim N(0, 1)$. For the Squared model, $V = \theta_0' \mathbf{S}$ where $\mathbf{S} = (1, S_1, \dots, S_6)$ and $S_j = X_j^2$. The Qubic model is constructed similarly, but defines $V = \theta_0' \mathbf{Q}$, with $\mathbf{Q} = (1, Q_1, \dots, Q_6)$ and $Q_j = X_j^3/5$. For the Ker & Sam model, introduced in [Ker and Sam \(2018\)](#), $X_j \sim N(0, 1)$, $j = 1, 3, 5$, $X_j \sim \chi_3^2$, $j = 2, 4, 6$ and ϵ is drawn from a normal mixture such that $\epsilon \sim N(3, 1)$ or $\epsilon \sim N(-3, 1)$ with probability 0.5. This corresponds to an homoscedastic model but with a bimodal error term. The Scobit model is defined as the skewed version of the Logit model, considering the link function $\mu(v) = 1 - \frac{1}{(1+e^v)^6}$. For the Klein & Spady model, introduced in [Klein and Spady \(1993\)](#), a heteroskedastic error term $\epsilon | V \sim N(0, \sigma(V))$ is considered, with $\sigma^2(V) = 0.1(1 + V^2)^2$. For all simulations, $\theta_0 = (\theta_{00}, \theta_{01}, \dots, \theta_{06}) = (\theta_{00}, 1, 2, -1, -2, 3, -3)'$ and two scenarios are considered. In the normal scenario, θ_{00} is chosen so that $\bar{y} \approx 0.1$. In the risky scenario, θ_{00} is chosen so that $\bar{y} \approx 0.25$. Thus, model performance is also explored at different levels of case proportions.

The exogenous variable is assumed to be independent from the covariates following [Zadrozny and Elkan \(2001\)](#). It is simulated as $W = 100W_1 W_2$, with $W_1 \sim \chi_4^2$ and $W_2 \sim N(7, 0.8)$, where W_1 and W_2 are independent. The idea is to replicate the asymmetry commonly encountered in financial problems such as the one motivating this work.

Six models (introduced in [Table 2](#)), three cost settings ($a = \{0.05, 0.1, 0.2\}$) and two scenarios (normal and risky) are combined to generate 36 simulation settings. This provides a comprehensive test of the performance of the introduced classifiers under different dependency relationships and cost settings. For each scenario, samples are simulated with $n = 10,000$ and the algorithms are trained as introduced in Sections 3 and 4. This is repeated 20 times for each simulation scenario and the results are summarized by averaging over these 20 runs. Due to the number of simulations performed, the results are summarized in [Appendix](#). A summary of the results obtained with each of the algorithms is included in the form of graphs in [Figs. 3–6](#). In this way we have a quick and visual representation of the performance of the models and the details are left for the [Appendix](#).

5.1.2. Online learning algorithms

This section evaluates the performance of the OL algorithms presented in Section 3 over the simulated datasets. [Figs. 3 and 4](#) show the box plot of the savings in (27) obtained for the 6 simulation models and

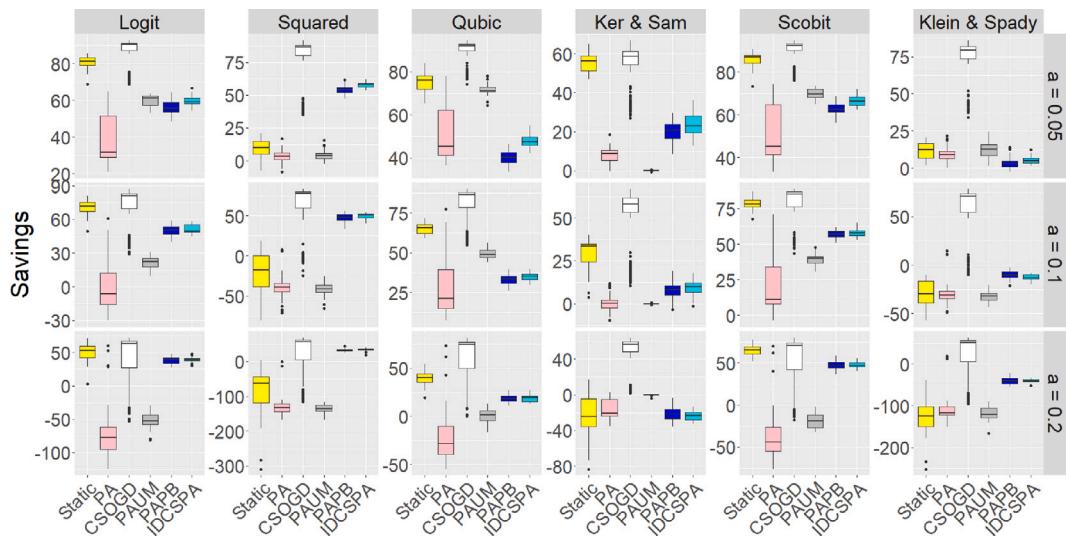


Fig. 3. Box plots of the savings in (27) for the OL algorithms in the normal scenario.

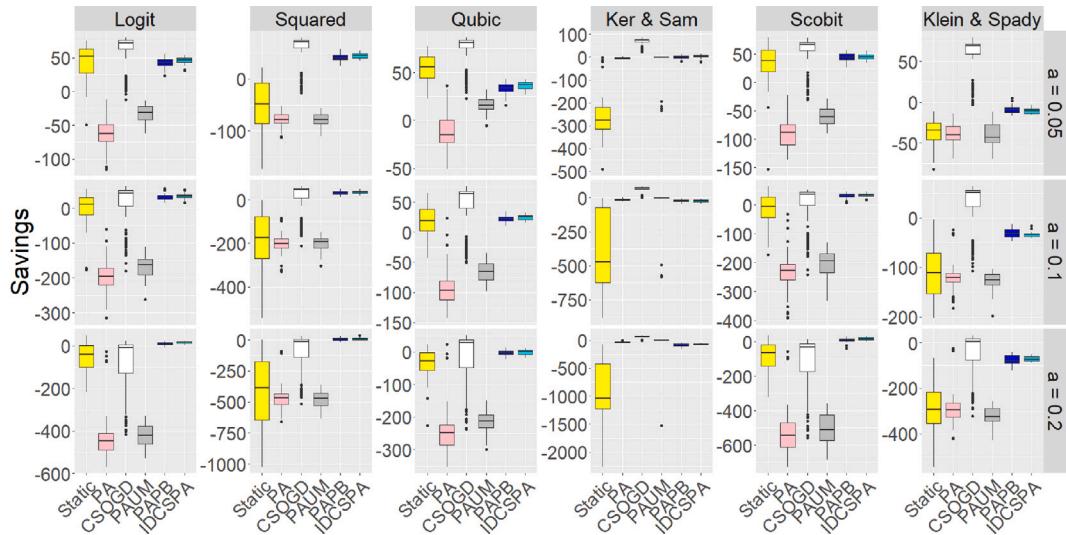


Fig. 4. Box plots of the savings in (27) for the OL algorithms in the risky scenario.

the 3 cost specifications for the normal and risky scenario, respectively. The effect of parameter α on model performance can be observed. This is due to the greater impact of false positives. Since OL algorithms only take losses into account, their performance decreases as the value of α increases. Something similar can be observed depending on the case proportion scenario. As expected, the more cases there are in the sample, the lower the savings is.

The static alternative performs well in many scenarios. However, the fact that the model is not updated with the new observations received implies a poor performance when the underlying model deviates from linearity. Similarly, the PA algorithm has the worst performance among the OL classifiers because it does not take costs into account.

The CSOGD model is the one that performs best in some scenarios, but still has a large variance in some settings. This is particularly the case for simulation models that are furthest from a linear model and where the impact of errors is greater. This variability across scenarios means that while it may be a good alternative, it must be treated with caution. Something similar happens with the PAUM model, although to a lesser extent.

As for the CS passive-aggressive algorithms, namely PAUM, PAPB, and IDCSPA, their behavior is generally quite similar, with a general

improvement with the instance-dependent algorithm. The IDCSPA algorithm does not perform better than the state-of-the-art in the normal scenario. This is because the model tends to mark many observations as cases relative to the true proportion of cases as it can be seen in . As the proportion of cases increases, the model has more information to update its predictions and achieves better performance. Therefore, the proposed model should be considered with caution in highly imbalanced scenarios. Nevertheless, it is the best in most of simulated settings and gives stable results with the least variability regardless of the problem setting as shown in Figs. 3 and 4. Lastly, it is worth noting that the proposed IDCSPA algorithm consistently performs better than the other passive-aggressive algorithms. Taking this into account, Algorithm IDCSPA is considered to be the best performing model in the simulation study. Therefore, it is an interesting model when dealing with dynamic problems, especially in scenarios where interpretability is important, such as the credit risk problem, as it provides an interpretable model with a proficient performance.

5.1.3. Bandit algorithms

This section evaluates the behavior of the introduced bandit algorithms over the simulated datasets. Figs. 5 and 6 show the mean

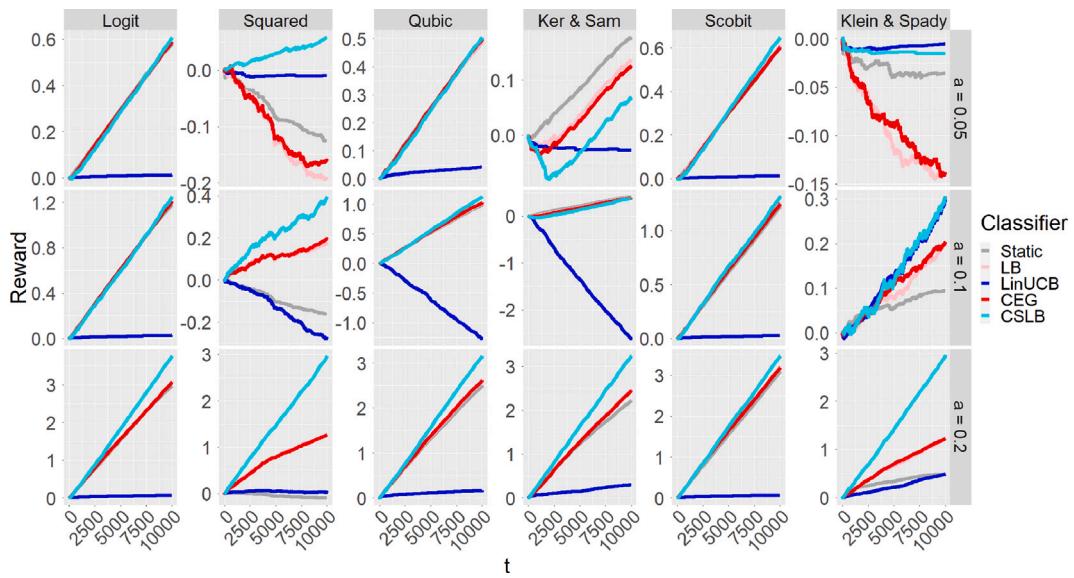


Fig. 5. Cumulative reward (in millions) for the bandit algorithms in the normal scenario.

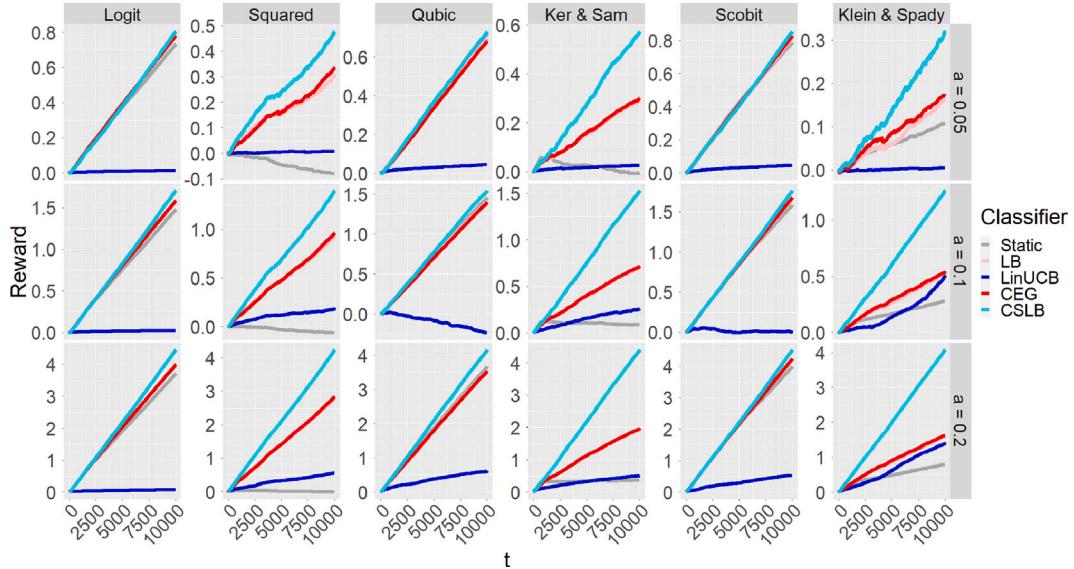


Fig. 6. Cumulative reward (in millions) for the bandit algorithms in the risky scenario.

cumulative reward obtained over the 20 simulations for the 6 simulation models and the 3 cost specifications for the normal and risky scenario, respectively. The X -axis represents the step and the Y -axis represents the cumulative reward.

For the bandit algorithms, as for the OL algorithms, we can see how the proportion of cases and the cost specification significantly affect the performance of the models. In this case, since we include both profits and losses in the objective function, the average reward increases as a does, which is the profit made on an approved good operation.

The LinUCB algorithm has a poor performance due to the fact that this algorithm only considers the context to model the expected reward. The exogenous variable is assumed to be independent of the covariates. Therefore, the model makes a naive prediction in all scenarios. Nevertheless, this optimistic approach appears effective in certain scenarios.

It is worth highlighting the good performance of the static model in some scenarios, such as Logit, Cubic or Scobit, despite the fact that it does not update predictions or take exploration into account. In the case of the Logit scenario, this good behavior is justified because the

model is well specified. In the other two, despite the fact that the model is not well specified, the link function is similar to the assumed by the logistic model, which also makes the model behave well. As the misspecification of the logistic model is further broken down, its behavior deteriorates significantly, as can be seen in the Squared or Klein & Spady scenarios. In these simulations, the exploration performed by the LB and CEG models allows them to adapt much better to the problem. This is reflected in the higher expected rewards despite not considering a CS approach for model fitting.

With respect to the proposed CSLB algorithm, it can be seen that it practically achieves the best results in all simulations when compared to the state-of-the-art algorithms, both cost-insensitive and cost-sensitive. The biggest differences are observed as the value of a increases, due to a higher potential profit per good operation. It is noteworthy that the CSLB model, unlike the other models, obtains better results in the most risky scenario. This is because the model explores a greater number of cases in the first steps. Thus, in the long run, it has a more representative sample and consequently more accurate and better results are obtained.

The only simulation scenario where the proposed model performs worse than state-of-the-art models is the Ker & Sam simulation with $a = 0.05$ in the normal scenario. This is due to a significant loss in the early iterations, and although it is corrected, it is not compensated in the run considered. Perhaps, given the results obtained in the other simulation models, the CSLB will outperform the other approaches over a longer time frame.

Taking into account that the CSLB obtains the best results regardless of the proportion of cases, relationship dependencies, and cost specification, it is considered that the proposed CSLB algorithm achieves satisfactory performance and is an interesting alternative for cost-sensitive dynamic problems.

5.2. Real datasets study

This section examines the behavior of the reinforcement learning algorithms presented on various real credit risk data sets. The objective is to evaluate the performance of the different algorithms in real scenarios, after having studied their performance in simulated environments. After introducing the data sets, the results obtained for the different algorithms are presented, together with some comments on their performance.

The different datasets are selected to investigate the impact of class imbalance and the number of covariates on model performance. As none of the state-of-the-art datasets provide guidance on the cost of misclassification errors, several cost specifications are considered. Likewise, the effect of different ratios of false positive and false negative costs on model performance is studied. The values for the parameters of the cost matrix in Table 1 are taken in line with those considered in the previous section, i.e. $a \in \{0.05, 0.10, 0.20\}$, $b = 10$ and $c = 0.75$.

Five real credit risk data sets are considered, where some covariates are available to model a response variable (fraud or default). The objective is to minimize losses/maximize benefits, which depend on the amount of the transaction and the cost specification. For each data set, a 5-fold cross-validation is performed, stratified by case proportion and amount as suggested in Höppner et al. (2021). The mean results over the test sets are summarized in Tables 3–5 for the OL algorithms, and in Fig. 7 and in Tables A.12–A.14 in the Appendix for the bandit algorithms.

5.2.1. Datasets

The first data set (Credit Card¹ data set) is a fraud data set containing 28 variables resulting from a PCA along with a “Time” variable. The data set is undersampled so that they remain 15,000 transactions containing 465 (3.19%) frauds. The second one (GCD² data set) is a real credit risk data set of size $n = 1000$ with a case proportion of 30% and 24 variables available to model the default indicator. The third one (PAKDD³ data set, as introduced in Linhart et al., 2009) consists of 11,111 prospective clients for which 4 numerical variables are available to model their probability of default. There are 3,047 (27.42%) defaults in the sample. The amount of the loan is simulated as $W \sim 100\chi_3$ in order to replicate the asymmetry commonly encountered in credit risk problems. The fourth one (ASF data set) is loaned by a collaborating bank (Abanca Servicios Financieros). It is composed of 23,377 loan requests collected between January 2022 and December 2022 with a default percentage of 3.06%. To preserve confidentiality, the number of registers in the original dataset is truncated, and so is the default proportion. It consists on 13 numerical variables which include socioeconomic information of the client and the point of sale

Table 3

Savings (Sav), accuracy (Acc), recall (Rec), positive predicted value (PPV), positive predicted proportion (PP) and F-score (FS) for the OL algorithms for the real datasets and cost setting $a = 0.05$.

Data set	Classifier	Sav	Acc	Rec	PPV	PP	FS
Credit card	Static	-52.36	43.23	11.18	0.66	54.30	1.24
	PA	-20.05	52.65	52.90	3.55	47.54	6.65
	CSOGD	63.71	73.93	98.44	11.36	29.16	20.28
	PAUM	-10.85	50.14	68.39	4.27	51.03	8.04
	PAPB	62.44	96.72	79.78	49.14	5.17	60.82
	IDCSPA	63.68	96.91	78.28	50.98	4.89	61.75
GCD	Static	27.03	53.65	25.00	23.96	31.33	24.47
	PA	31.82	65.57	35.67	41.47	25.83	38.35
	CSOGD	74.96	91.42	79.58	90.69	26.35	84.78
	PAUM	0.49	69.67	0.67	28.57	0.70	1.30
	PAPB	34.17	63.16	45.00	39.94	33.83	42.32
	IDCSPA	34.78	62.86	45.33	39.65	34.33	42.30
PAKDD	Static	34.96	54.85	42.63	28.21	41.13	33.95
	PA	22.76	61.01	26.28	27.43	26.08	26.84
	CSOGD	76.10	86.98	76.77	75.72	27.60	76.24
	PAUM	-0.01	72.77	0.00	0.00	0.01	
	PAPB	0.04	72.77	0.03	33.33	0.02	0.05
	IDCSPA	19.95	63.78	22.08	28.60	21.02	24.92
ASF	Static	-32.71	68.59	32.03	3.23	30.31	5.87
	PA	-17.03	61.64	58.47	5.05	38.93	9.31
	CSOGD	58.72	82.74	87.96	13.83	19.59	23.9
	PAUM	-3.85	63.75	77.78	6.31	37.96	11.67
	PAPB	-4.44	90.56	14.81	6.27	7.27	8.82
	IDCSPA	-4.40	90.50	14.81	6.23	7.33	8.77
Drift	Static	-22.70	83.35	5.33	1.44	13.41	2.27
	PA	-3.78	92.18	3.81	3.10	4.47	3.42
	CSOGD	72.87	98.03	71.95	73.36	3.56	72.65
	PAUM	-1.42	95.27	0.55	1.75	1.14	0.84
	PAPB	27.44	92.99	35.77	21.71	5.98	27.02
	IDCSPA	39.52	93.00	50.92	26.16	7.07	34.56

where the credit is requested and characteristics of the operation. The fifth one (Drift data set) consists of 60,000 synthetically generated observations from the ASF data set, where concept and covariate drifts are introduced to test the performance of the different models in the presence of drifts. To do this, we follow (Nowok et al., 2016) with the syn function of the R package synthpop. First, a sample of size 15,000 is synthetically generated from the ASF data. To this data set we add 15,000 observations generated with a different distribution for the covariates, but preserving the dependence between X and Y. Another 15,000 instances are added introducing concept drift, generating the response variable with a different model than the one considered for the first 30,000 observations. Finally, 15,000 synthetically generated observations from the original ASF data set are added to form the final sample.

5.2.2. Online learning algorithms

This section evaluates the performance of the online learning algorithms presented in Section 3 over the real datasets previously introduced. Tables 3–5 summarize the savings obtained with the 6 algorithms compared for the 5 datasets and the 3 cost specifications. It can be observed how the performance of the models varies significantly depending on the value of the parameter a and the case proportion as in the simulation study. As expected, when there are more bad requests in the sample, the savings obtained is lower.

In line with the findings of the simulation study, the CSOGD algorithm consistently demonstrates the best overall performance. The proposed IDCSPA model is the next best performing algorithm across all scenarios. Although the IDCSPA algorithm does not outperform the CSOGD, it achieves the best results within the passive-aggressive models. This suggests that incorporating the instance-dependent loss function into the training process improves the performance of OL algorithms. These results corroborate in practical scenarios the results of the simulation study and indicate the potential for extending other

¹ <https://kaggle.com/mlg-ulb/creditcardfraud>.

² <https://github.com/JLZml/Credit-Scoring-Data-Sets/tree/master/1.20UCI%20Repository/German>.

³ <https://github.com/JLZml/Credit-Scoring-Data-Sets/tree/master/2.20PAKDD%202009%20Data%20Mining%20Competition>.

Table 4

Savings (Sav), accuracy (Acc), recall (Rec), positive predicted value (PPV), positive predicted proportion (PP) and F-score (FS) for the OL algorithms for the real datasets and cost setting $a = 0.1$.

Data set	Classifier	Sav	Acc	Rec	PPV	PP	FS
Credit card	Static	-151.97	42.72	12.69	0.74	54.90	1.39
	PA	-103.49	51.51	50.32	3.31	48.51	6.20
	CSONGD	30.12	72.32	98.82	10.77	30.80	19.35
	PAUM	-69.31	48.34	72.26	4.34	53.08	8.19
	PAPB	54.48	97.22	79.78	54.32	4.68	64.63
GCD	IDCSPA	59.12	97.10	79.14	53.03	4.76	63.50
	Static	5.74	59.26	15.67	23.38	20.12	18.76
	PA	21.18	64.96	35.33	40.46	26.23	37.72
	CSONGD	70.22	90.22	75.00	90.82	24.80	82.15
	PAUM	-0.13	69.87	0.00	0.00	0.10	
PAKDD	PAPB	26.96	62.36	43.00	38.62	33.43	40.69
	IDCSPA	27.35	62.76	43.67	39.22	33.43	41.32
	Static	12.69	60.18	23.01	24.93	25.12	23.93
	PA	17.52	61.42	25.76	27.63	25.38	26.66
	CSONGD	73.77	86.98	77.24	75.51	27.84	76.36
ASF	PAUM	-0.01	72.74	0.03	14.29	0.05	0.05
	PAPB	-0.00	72.77	0.00	0.00	0.01	
	IDCSPA	14.16	63.76	21.45	28.21	20.70	24.37
	Static	-64.39	78.35	12.31	1.95	19.34	3.36
	PA	-91.36	60.75	61.98	5.34	40.07	9.83
Drift	CSONGD	30.67	84.11	86	13.67	17.94	23.59
	PAUM	-75.38	66	66	5.39	34.91	9.97
	PAPB	-12.8	91.19	18	7.35	6.99	10.43
	IDCSPA	-6.13	91.99	21	9.42	6.36	13
	Static	-49.56	83.35	5.33	1.44	13.41	2.27
ASF	PA	-10.99	92.18	3.81	3.10	4.47	3.42
	CSONGD	71.37	98.03	71.95	73.36	3.56	72.65
	PAUM	-3.21	95.27	0.55	1.75	1.14	0.84
	PAPB	17.72	92.54	36.96	20.59	6.52	26.45
	IDCSPA	29.05	92.69	52.11	25.35	7.46	34.10

Table 5

Savings (Sav), accuracy (Acc), recall (Rec), positive predicted value (PPV), positive predicted proportion (PP) and F-score (FS) for the OL algorithms for the real datasets and cost setting $a = 0.2$.

Data set	Classifier	Sav	Acc	Rec	PPV	PP	FS
Credit card	Static	-178.65	56.18	9.68	0.75	41.25	1.39
	PA	-239.81	51.99	50.32	3.34	48.03	6.26
	CSONGD	-53.32	71.78	98.06	10.55	31.28	18.97
	PAUM	-187.69	47.95	80.43	4.75	53.99	8.97
	PAPB	35.53	97.05	78.92	52.43	4.80	63.00
German	IDCSPA	45.16	97.59	80.65	58.96	4.36	68.12
	Static	14.14	59.76	44.00	36.07	36.64	39.64
	PA	13.62	65.17	34.67	40.62	25.63	37.41
	CSONGD	70.71	90.47	77.83	89.04	26.25	83.06
	PAUM	-0.16	69.77	0.00	0.00	0.20	
PAKDD	PAPB	15.26	63.96	43.33	40.62	32.03	41.94
	IDCSPA	14.65	61.46	48.67	38.73	37.74	43.13
	Static	7.96	56.07	33.49	26.09	34.94	29.33
	PA	7.75	61.79	25.38	27.84	24.81	26.56
	CSONGD	69.83	86.75	76.28	75.36	27.55	75.82
ASF	PAUM	0.24	72.73	0.35	39.39	0.25	0.70
	PAPB	-0.00	72.77	0.00	0.00	0.01	
	IDCSPA	6.89	63.03	21.48	27.27	21.44	24.03
	Static	-235.03	65.91	29.65	2.76	32.85	5.05
	PA	-265.14	60.87	59.84	5.23	39.82	9.62
Drift	CSONGD	-42.39	82.89	84.3	14.93	19.48	25.37
	PAUM	-243.08	64.06	68.6	6.36	37.22	11.64
	PAPB	-49.54	89.79	17.36	7.53	7.96	10.5
	IDCSPA	-42.01	90.7	12.4	6.38	6.7	8.43
	Static	-103.27	83.35	5.33	1.44	13.41	2.27
ASF	PA	-25.41	92.18	3.81	3.10	4.47	3.42
	CSONGD	68.37	98.03	71.95	73.36	3.56	72.65
	PAUM	-6.80	95.27	0.55	1.75	1.14	0.84
	PAPB	-2.16	92.54	36.96	20.59	6.52	26.45
	IDCSPA	7.93	92.82	51.56	25.67	7.29	34.27

online learning algorithms to cost-sensitive settings to achieve further improvements.

5.2.3. Bandit algorithms

This section evaluates the behavior of the bandit algorithms introduced in Section 4. Results are summarized in Fig. 7 and detailed in the Appendix in Tables A.12-A.14. Fig. 7 shows the mean cumulative reward obtained with the 5 considered algorithms for the 5 datasets and the 3 different cost specifications as in Fig. 5. As for the OL algorithms, we can see how the proportion of cases, the number of covariates and the cost specification significantly affect model performance.

As in the simulation study, the LinUCB algorithm is the worst performing approach. Regarding the performance of the Static approach, excluding the LinUCB algorithm, its performance is the worse. Thus, it is concluded that updating the models with all the available information and performing exploration significantly improves model performance, especially in real problems as those considered in this study.

With respect to the proposed CSLB algorithm, it can be seen that it achieves the best results in all the datasets and cost settings when compared to both cost-insensitive and cost-sensitive approaches. As for the simulation study, the biggest differences are observed as the value of a increases, where the CSLB algorithm is able to increase the profits even more thanks to a higher profit per good operation granted.

Regarding the Drift data set, it can be clearly seen a decline in model performance for all the algorithms as the drifts appear, specially in the step where it starts the concept drift. Nevertheless, the proposed CSLB is able to correct its estimates thanks to the exploration considered and obtain the best aggregated results in all the cost settings.

The proposed Algorithm CSLB obtains the best results in all the problems, regardless of the proportion of cases, sample size, relationship dependencies, and cost specification. Taking this into account, it is considered that the proposed model achieves satisfactory performance and is an interesting alternative for cost-sensitive dynamic problems as the credit risk problem.

6. Conclusions and future work

Motivated by the credit risk problem, this paper presents two novel algorithms designed to tackle cost-sensitive dynamic classification. In this paper, the PA algorithm (Crammer et al., 2006), is extended to the CS setting in Algorithm IDCSPA. In addition, the logistic bandit approach is extended with Algorithm CSLB for efficient CS learning considering exploration. The performance of the proposed models is empirically tested considering a wide range of simulations and the study of five real datasets. Promising experimental results demonstrate the effectiveness and efficiency of the proposed models compared to state-of-the-art algorithms. Furthermore, this study provides practitioners with valuable insights into the algorithmic performance in terms of case proportions, dependency relationships, and cost specifications.

Regarding OL algorithms, CSONGD and PAUM perform well in some scenarios but poorly in others, especially when misclassification error costs are larger. While PAPB and IDCSPA show similar behavior, the proposed IDCSPA algorithm emerges as the best model overall, offering consistent performance across different scenarios.

As for the bandit algorithms, the static approach gives a good performance in some simulations. Nevertheless, it suffers a performance decline in the real scenarios. Bandit algorithms such as ML and CEG correct model misspecification through exploration, leading to better adaptation and higher rewards. The proposed CS logistic bandit consistently outperforms state-of-the-art algorithms particularly excelling as the profit per good operation increases. This shows that the exploration performed by CSLB combined with a CS approach can significantly improve the potential benefits.

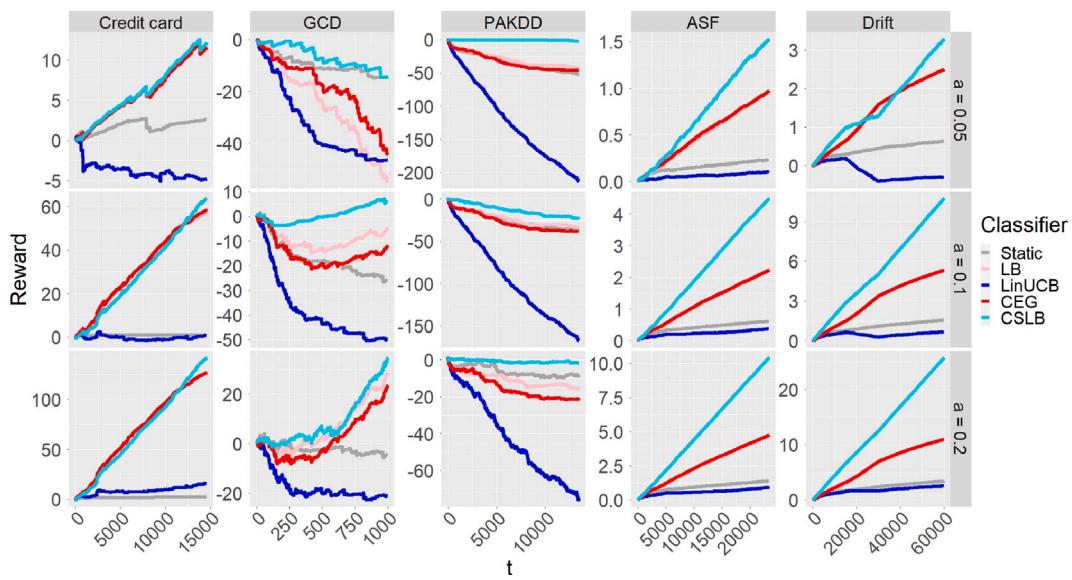


Fig. 7. Real datasets results (cumulative reward, in thousands for the credit card, GCD and PAKDD datasets, and in millions for the ASF and drift datasets) for the bandit algorithms.

The CS algorithms proposed in this paper prove to be effective and stable in most scenarios, making them good approaches for cost-sensitive dynamic problems. By proposing one approach that considers exploration and another that does not, the full range of alternatives is covered from the perspective of reinforcing learning in cost-sensitive classification problems. Thus, the dynamic cost-sensitive learning problem is solved for any setting.

Lastly, the good performance of the proposed algorithms opens the door to extend more reinforcement learning algorithms to the CS setting. In addition, one can consider extending algorithms that are not necessarily from the reinforcement learning literature. For example, consider the thresholding stage with a dynamic approach, updating the decision as more information becomes available. For this, it would be interesting to consider an empirical algorithm such as the one proposed in C-Rella et al. (2024) and C-Rella and Vilar (2024). The possible extensions are innumerable combining a cost-sensitive and dynamic approach.

CRediT authorship contribution statement

Jorge C-Rella: Conceptualization, Methodology, Software, Investigation, Writing. **David Martínez Rego:** Conceptualization, Validation, Investigation. **Juan M. Vilar:** Conceptualization, Validation.

Table A.6

Savings (Sav), accuracy (Acc), recall (Rec), positive predicted value (PPV), positive predicted proportion (PP) and F-score (FS) for the OL algorithms for the cost setting $a = 0.05$ in the simulated datasets.

Model	Classifier	Normal						Risky					
		Sav	Acc	Rec	PPV	PP	FS	Sav	Acc	Rec	PPV	PP	FS
Logit	Static	82.81	90.80	94.55	54.03	18.45	68.46	68.42	93.02	94.33	34.51	10.29	50.41
	PA	39.26	55.65	77.83	16.42	50.26	27.11	-62.73	50.83	68.03	5.11	50.54	9.50
	CSOGD	87.74	87.67	99.03	49.18	22.71	65.02	61.12	86.71	97.19	25.19	16.87	39.44
	PAUM	61.37	60.25	97.98	20.78	49.91	34.28	-30.05	52.81	96.60	7.22	50.73	13.43
	PAPB	57.86	89.71	62.47	51.25	12.95	56.22	43.25	93.99	54.38	32.92	6.35	40.71
	IDCSPA	58.66	88.02	66.03	44.67	15.31	53.28	44.47	93.90	55.30	32.08	6.50	40.34
Squared	Static	3.67	74.67	30.81	16.85	21.74	19.80	-6.55	88.18	30.99	11.53	10.49	16.57
	PA	2.50	49.46	50.07	9.87	50.56	16.49	-82.80	49.99	50.64	3.72	50.06	6.92
	CSOGD	79.10	86.68	92.34	45.98	21.77	60.84	57.62	86.72	94.58	24.55	16.56	38.52
	PAUM	3.39	50.24	49.97	10.01	49.75	16.67	-87.50	49.66	47.42	3.48	50.15	6.49
	PAPB	54.88	89.43	60.74	47.72	12.72	53.36	42.72	94.05	54.17	32.13	6.28	40.11
	IDCSPA	59.19	88.19	66.64	43.94	15.13	52.96	48.35	94.03	61.25	34.03	6.83	43.74

(continued on next page)

- Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6), 368–408. <http://dx.doi.org/10.1037/h0042519>.
- Russo, D., Roy, B. V., Kazerouni, A., Osband, I., & Wen, Z. (2020). A tutorial on thompson sampling. [arXiv:1707.02038](https://arxiv.org/abs/1707.02038).
- Shalev-Shwartz, S. (2012). Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4, 107–194. <http://dx.doi.org/10.1561/2200000018>, URL: <https://api.semanticscholar.org/CorpusID:51730029>.
- Shen, W., Wang, J., Jiang, Y.-G., & Zha, H. (2015). Portfolio choices with orthogonal bandit learning. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence*. URL: <https://www.ijcai.org/Proceedings/15/Papers/142.pdf>.
- Stripling, E., Broucke, S. v., Antonio, K., Baesens, B., & Snoeck, M. (2018). Profit maximizing logistic model for customer churn prediction using genetic algorithms. *Swarm and Evolutionary Computation*, 40, 116–130. <http://dx.doi.org/10.1016/j.swevo.2017.10.010>.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. Cambridge, MA, USA: A Bradford Book, ISBN : 978-0-262-19398-6.
- Tang, L., Rosales, R., Singh, A., & Agarwal, D. (2013). Automatic ad format selection via contextual bandits. In *Proceedings of the 22nd ACM international conference on information and knowledge management (CIKM '13)* (pp. 1587–1594). <http://dx.doi.org/10.1145/2505561.2505602>.
- Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4), 285–294. <http://dx.doi.org/10.2307/2332286>.
- Vanderschueren, T., Verdonck, T., Baesens, B., & Verbeke, W. (2022). Predict-then-optimize or predict-and-optimize? An empirical evaluation of cost-sensitive learning strategies. *Information Sciences*, 594, 400–415. <http://dx.doi.org/10.1016/j.ins.2022.02.021>.
- Vapnik, V. N. (1998). *Statistical learning theory*. Wiley, ISBN: 978-0-471-03003-4.
- Visantavarakul, K. (2022). *An application of reinforcement learning to credit scoring based on the logistic Bandit framework*. Chulalongkorn University Theses and Dissertations (Chula ETD 6050), URL: <https://digital.car.chula.ac.th/chulaetd/6050>.
- Wang, C.-C., Kulkarni, S., & Poor, H. (2005). Bandit problems with side observations. *IEEE Transactions on Automatic Control*, 50(3), 338–355. <http://dx.doi.org/10.1109/TAC.2005.844079>.
- Wang, J., Zhao, P., & Hoi, S. C. (2012). Cost-sensitive online classification. In *2012 IEEE 12th International conference on data mining* (pp. 1140–1145). <http://dx.doi.org/10.1109/ICDM.2012.116>.
- Zadrozny, B., & Elkan, C. (2001). Learning and making decisions when costs and probabilities are both unknown. vol. 2001, In *Proceedings of the seventh ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 204–213). <http://dx.doi.org/10.1145/502512.502540>.
- Zhang, Y., Wu, M., Hu, X., & Zhu, Y. (2018). An online transfer learning algorithm with adaptive cost. In *Proceedings of the 2018 international conference on signal processing and machine learning*. URL: <https://api.semanticscholar.org/CorpusID:69168504>.
- Zhang, L., Yang, T., Jin, R., Xiao, Y., & Zhou, Z. (2016). Online stochastic linear optimization under one-bit feedback. vol. 48, In *Proceedings of the 33rd international conference on machine learning* (pp. 392–401). New York, New York, USA: PMLR, URL: <https://proceedings.mlr.press/v48/zhangb16.html>.
- Zhao, P., Zhuang, F., Wu, M., Li, X., & Hoi, S. C. H. (2015). Cost-sensitive online classification with adaptive regularization and its applications. In *IEEE international conference on data mining ICDM 2015* (pp. 649–658). URL: https://ink.library.smu.edu.sg/sis_research/2923.
- Zhou, L. (2016). A survey on contextual multi-armed bandits. [arXiv:1508.03326](https://arxiv.org/abs/1508.03326).