

Offline Deep Reinforcement Learning for Dynamic Pricing of Consumer Credit

Raad Khraishi^{1, 2, *} and Ramin Okhrati¹

¹*Institute of Finance and Technology, UCL, London, United Kingdom*

²*Data Science and Innovation, NatWest Group, London, United Kingdom*

**Corresponding author: Raad Khraishi, raad.khraishi.13@ucl.ac.uk*

Abstract

We introduce a method for pricing consumer credit using recent advances in offline deep reinforcement learning. This approach relies on a static dataset and requires no assumptions on the functional form of demand. Using both real and synthetic data on consumer credit applications, we demonstrate that our approach using the conservative Q-Learning algorithm is capable of learning an effective personalized pricing policy without any online interaction or price experimentation.

Keywords: Reinforcement Learning, Finance, Pricing, Revenue Management, Consumer Credit

1 Introduction

Consumer debt in the United States alone is worth over \$15 trillion.¹ Despite the importance of this market, setting interest rates for debt products is done with varying levels of sophistication. Two common techniques used by lenders today are risk-based and profit-based pricing (Phillips, 2020). Risk-based pricing involves adding a fixed margin on top of the expected cost including default for a specific loan or pricing segment. Profit-based pricing extends this by also incorporating the estimated responsiveness of customers or customer segments to price to find the profit-maximizing interest rate.

This present work builds on profit-based pricing by introducing a model-free reinforcement learning approach to finding optimal prices. In particular, we develop an approach for pricing installment credit products such as mortgages as well as personal and student loans which form the bulk of the consumer debt market. Setting prices for these products involves factoring in short-term rewards from the underwriting of the product as well as longer-term default risk, adverse selection, market competitiveness, and customer lifetime value. With reinforcement learning we are able to formulate this problem as a sequence of pricing decisions for each loan applicant. This formulation helps induce long-term consequences of actions by allowing pricing decisions taken by the agent to affect the future state of the environment. In addition, we use a model-free reinforcement learning algorithm in which an agent must learn to implicitly estimate individual price-responsiveness as it learns a pricing policy. This is a key departure from traditional profit-based pricing approaches that make strong assumptions on the functional form of demand. Furthermore, a reinforcement learning approach also allows us to learn a dynamic pricing policy that can adapt to changes in behavioral patterns and the economy (Rana and Oliveira, 2014).

Traditional reinforcement learning algorithms learning from scratch by pricing consumer loans in a live environment may start off by setting inaccurate prices that would incur significant financial and reputational costs for a lender. The potential for these costs makes traditional online reinforcement learning difficult to

¹https://www.newyorkfed.org/medialibrary/interactives/householdcredit/data/pdf/HHDC_2021Q3.pdf

apply in practice and motivates the need for offline reinforcement learning algorithms that are able to learn a policy from a static dataset with past pricing decisions and outcomes.

Offline learning introduces several additional challenges. In particular, the agent must assess counterfactual scenarios about what might happen following an action despite not having any examples of that behavior in the training set. The decisions in the dataset used to learn a policy may be different than the decisions that it must learn to apply. This problem of distributional shift leads traditional reinforcement learning approaches to overestimate the value of unseen outcomes and propagate poor decisions (Levine et al., 2020). Given limited changes in prices historically, distributional shift and generalization are particularly problematic for pricing consumer credit products. To mitigate these issues this paper uses the conservative Q-learning (CQL) algorithm to regularize Q-values and reduce overestimation of out-of-distribution actions (Kumar et al., 2020).

The main contribution of this paper is the introduction of an approach to pricing consumer credit that uses model-free offline deep reinforcement learning to learn a sequential pricing policy. This approach makes no assumptions on demand often used in traditional pricing approaches that may introduce misspecification errors in unknown and non-stationary environments. We also extend credit pricing to the full reinforcement learning problem in which current actions may impact the future state of the environment. Using synthetic and real data on auto loans, we demonstrate that this approach is able to learn an effective policy using only a static dataset without any live interaction or price experimentation.

The rest of this paper is organized as follows. Section 2 presents a review of related pricing and reinforcement learning literature. Section 3 formulates credit pricing as a reinforcement learning problem and describes the algorithm and evaluation approaches used. Section 4 presents our result on both historic and synthetic datasets of loan applications. Finally, Section 5 presents our conclusions and possible extensions of our work.

2 Related Work

Despite many breakthroughs in reinforcement learning across several problem domains (see Li (2018) for an overview), the literature on pricing and consumer credit remains sparse. In perhaps the paper most closely related to our current work, Trench et al. (2003) used tabular value-function iteration to choose credit line increases and interest rate decreases for credit cards at regular intervals. They applied a model-based approach that required estimation of a separate static transition matrix and heavily discretized the state and action space to make the problem tractable. In addition, they introduced several constraints on the action space such as only considering interest rate reductions or credit line increases. Despite these simplifications, the authors state that their model was implemented by Bank One and estimated an associated increase in annual profit of approximately 5% (equivalent to \$75 million per year).

We improve on their work in several ways. First, we do not discretize prices or restrict actions to only price decreases. Second, we use a model-free deep reinforcement learning approach with continuous prices and states that does not involve estimation of a separate transition matrix and does not assume knowledge of the demand curve. In addition, we also treat credit pricing as a continuing infinite horizon problem which is better suited to capturing the delayed effects of adverse-selection and price competitiveness.

Several studies have applied profit-based pricing approaches to the same historical auto loans dataset used in Section 4. Phillips et al. (2015) implement a profit-based pricing approach using a logistic regression price-response model and a simplified profit objective on the same historical auto loans dataset used in this paper to estimate the effectiveness of field price discretion. Ban and Keskin (2021) extended this work by introducing an iterated approach to learning the price-response curve using a generalized linear regression model and the capacity for price experimentation. While both papers cite substantial increases in revenue, 38% and 47% respectively, their estimated performance relies heavily on their assumed demand behavior. In Section 4, we compare our reinforcement learning approach against a similar profit-based pricing approach. Several bandit approaches have also been explored recently. Bastani et al. (2019) introduced an approach that uses Thompson sampling and a linear demand model to learn a pricing policy across multiple related products through pricing experiments. They defined a discrete set of auto loan products whereas in our approach the number of products is not fixed and new products may be defined as part of the continuous

state. Luo et al. (2021) applied a contextual linear bandit approach using a modified linear upper confidence bound algorithm. As with Ban and Keskin (2021), both these approaches rely on random and potentially costly price experimentation to learn a policy. In addition, none of the papers extend credit pricing to the full reinforcement learning problem where actions may impact the future state of the environment and use simple linear function approximation.

Outside of consumer credit, many successful applications of dynamic pricing using bandit or reinforcement learning algorithms exist. For example, Cheung et al. (2014) developed pricing policies in the presence of unknown demand with limited ability to perform price experimentation and personalization. They implemented their pricing strategy for an online deal website, Groupon, and estimated a 21% increase in daily revenue. Chen et al. (2021) explored personalized pricing as well as assortment optimization for airline seating reservations in a single-period problem. Cohen et al. (2020) developed an online contextual bandit approach for pricing online fashion products with each product defined by a set of features. Trovò et al. (2018) applied multi-armed bandit algorithms to online pricing of non-perishable goods in both stationary and non-stationary environments. Several other papers have extended dynamic pricing to the full reinforcement learning problem where an agent must consider the long-term consequences of its actions. For example, Rana and Oliveira (2014) developed a model-free tabular Q-learning approach for selling a fixed inventory by a given deadline. Using simulated data, they found improved performance in the presence of demand misspecification using their model-free approach relative to classic parametrized pricing policies. Krashennikova et al. (2019) applied model-free reinforcement learning for finding a pricing policy for insurance renewals to maximize revenue subject to a minimum renewal rate constraint using a synthetic dataset.

Pricing consumer credit, however, has two unique characteristics that distinguish it from pricing many types of consumer goods and services. First, adverse selection may arise when higher risk borrowers are less sensitive to price than borrowers with lower default risk (Phillips and Raffard, 2011). This behavior leads increases in prices to increase the riskiness of a lender’s portfolio. Second, instalment loans are contractual products with terms that may vary from several months to several decades in which a borrower makes regular payments towards the original principal and interest. As such, revenue is not realized directly upon purchase as future customer behavior, such as failing to repay or early repayment, over the term of the product may impact profitability. Therefore, at the time the price of a loan is set, the loan’s profitability is a random variable.

Offline reinforcement learning techniques have been applied to many different problems where learning online is costly or inefficient. Abe et al. (2010) implemented a constrained offline reinforcement learning approach based on Q-learning for optimizing debt collections using linear function approximation. They also implemented their solution to manage collections for the New York State Department of Taxation and Finance and estimated expected savings of \$100 million over a three year period. Theocharous et al. (2015) used a classic offline reinforcement learning algorithm, Fitted Q-Iteration (Ernst et al., 2005), for personalized ad recommendations. By using reinforcement learning methods, they found that their policy was able to take into account the long-term effect of ad choice on customer lifetime value.

3 Methodology

In this section we formulate credit pricing as a Markov decision process, introduce our offline reinforcement learning approach, and describe our approach to policy evaluation.

3.1 Markov Decision Process

Following Levine et al. (2020), we define a Markov decision process to be a tuple $(\mathcal{S}, \mathcal{A}, T, d_0, r, \gamma)$, where \mathcal{S} represents the set of states $s \in \mathcal{S}$, \mathcal{A} represents the set of actions $a \in \mathcal{A}$, T is the unknown transition dynamics $T(s_{t+1}|s_t, a_t)$, d_0 is the initial state distribution, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ denotes the reward function, and $\gamma \in (0, 1)$ is the discount factor. We also define credit pricing to be a continuing problem where each time step t relates to a new application and the total number of time steps, H , equals ∞ .

States. Our state space, \mathcal{S} , includes information that is typically available in a loan application such as credit score, the type of loan, the term of the loan, the amount of the loan, and competitor rates.

For example, with the auto loans dataset described in Section 4.1.1, we define $s_t = [Term_t, Amount_t, FICO_t, PD_t, PreviousRate_t, CompetitionRate_t, PrimeRate_t, Tier_t, LoanType_t, CarType_t, PartnerBin_t, State_t, Months_t, DayOfWeek_t, MonthOfYear_t, DaysSinceApp_t]$. See Appendix A.1 for more details.

Additional features that are typically available through a credit application process such as tenure, channel, debt obligations, and income can be used to extend the state space though care must be taken to ensure any feature does not lead to biased or unfair pricing policies.

Although we refer to state for simplicity, our environment is partially observable due to incomplete information on factors that may impact a lender’s profitability. For example, a potential borrower who has recently lost their job may be more willing to accept an uncompetitive interest rate for fear of not receiving another offer later. This behavior would impact both the price sensitivity of the borrower as well as their riskiness which may induce adverse selection and is not captured by the feature set (for a detailed discussion of price-dependent risk see Phillips and Raffard (2011)).

Actions. Upon each application the agent chooses an interest rate, a_t , quoted as an annual percentage rate (APR) which determines the price of the loan. We define our action space, \mathcal{A} , to be the subset of positive real numbers, i.e. $a_t \in \mathbb{R}^+$. This may be extended to include product fees as well.

Rewards. Many alternative financial measures of profit such as Net Income or Net Interest Income may be used to define rewards. For example, with the auto loans dataset, we use a simplified expected profit measure from Phillips et al. (2015) which accounts for interest income, capital costs, and credit risk. We use the shorthand $r_t = r(s_t, a_t)$ to denote a scalar reward value, which is defined by:

$$r(s_t, a_t) = p(Accept_t | a_t, s_t) * \left[(1 - PD_t) * (TotalPayment_t - CapitalCost_t) - PD_t * LGD_t * CapitalCost_t \right], \quad (1)$$

where $p(Accept_t | a_t, s_t)$ denotes the probability of a potential borrower accepting a loan given a price, a_t , and state feature vector, s_t . PD_t is the estimated probability of default of the applicant.² We set the loss given default, LGD , to 50% for all customers as in Phillips et al. (2015). $TotalPayment_t$ is shorthand for the total value of interest and principle payments of the loan over the course of its term as a function of its interest rate, a_t , loan amount, and term. A similar function is used for the cost of capital, $CapitalCost_t$, with the prime rate, $PrimeRate_t$, used in place of a_t .

The goal of our pricing algorithm will be to learn a policy $\pi(a_t | s_t)$ that defines a distribution over actions conditioned on states. We can use this to define a trajectory, τ , to be a (potentially infinite) sequence of states and actions $(s_0, a_0, \dots, s_H, a_H)$, admitting the following distribution:

$$p_\pi(\tau) = d_0(s_0) \prod_{t=0}^H \pi(a_t | s_t) T(s_{t+1} | s_t, a_t).$$

We may also define the return, R_t , to be sum of discounted rewards $\sum_{t=0}^T \gamma^t r(s_t, a_t)$. Our objective then becomes to find a policy that maximizes the expected return, $\mathbb{E}_{\tau \sim p_\pi(\tau)} \left[\sum_{t=0}^H \gamma^t r(s_t, a_t) \right]$, denoted by $J(\tau)$.

3.2 Reinforcement Learning

Traditional pricing techniques for consumer credit such as risk-based pricing fail to capture idiosyncrasies in customer behavior and responsiveness to prices (Phillips, 2020). While more sophisticated profit-based pricing

²The probability of default, PD , is not directly available in the dataset and is estimated from the FICO scores using data from Munkhdalai et al. (2019) which covers roughly the same time period. However, as their data is not specific to auto loans we expect our PD values to be conservative estimates of default risk.

approaches attempt to estimate this they require strong assumptions on the functional form of demand (Phillips et al., 2015; Ban and Keskin, 2021). In addition, these techniques are often myopic - they use the estimated price-response curve to find the profit-maximizing prices without taking into account the long-term consequences of actions. Reinforcement learning algorithms help address these issues by allowing us to train an agent to learn to make sequential decisions based on past experience in an unknown and non-stationary environment (Sutton and Barto, 2018). In addition, model-free reinforcement learning techniques do not assume knowledge of the reward function, $r(s_t, a_t)$, price-response function, $p(\text{Accept}_t | s_t, a_t)$, or transition probabilities, $T(s_{t+1} | s_t, a_t)$.

We use the value function, V_π , to refer to the estimated value of being in a state, s_t , and following policy $\pi(a_t | s_t)$. We also define the action-value function, Q_π , to be the estimated value of being in state, s_t , taking action a_t and then following policy π . These functions are defined as follows:

$$V_\pi(s_t) = \mathbb{E}_{\tau \sim p_\pi(\tau | s_t)} \left[\sum_{t'=t}^H \gamma^{t'-t} r(s_{t'}, a_{t'}) \right],$$

$$Q_\pi(s_t, a_t) = \mathbb{E}_{\tau \sim p_\pi(\tau | s_t, a_t)} \left[\sum_{t'=t}^H \gamma^{t'-t} r(s_{t'}, a_{t'}) \right].$$

Writing the action-value function only in terms of Q :

$$Q_\pi(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim T(s_{t+1} | s_t, a_t), a_{t+1} \sim \pi(a_{t+1} | s_{t+1})} [Q_\pi(s_{t+1}, a_{t+1})],$$

gives rise to the Q-learning algorithm which iteratively applies the Bellman optimality operator, $\mathcal{B}^* Q(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim T(s_{t+1} | s_t, a_t)} [\max_{a_{t+1}} Q(s_{t+1}, a_{t+1})]$, to find the optimal Q-function. These estimates are then used to take actions usually following an implicit policy (e.g., ϵ -greedy). The Q-values may be approximated using a deep neural network with parameters θ , in which case we use Q_θ to denote the action-value function.

Actor-critic algorithms are another class of popular reinforcement learning algorithms and closely resemble classic policy iteration from dynamic programming (Sutton and Barto, 2018). These algorithms directly estimate both a parameterized value function, Q_θ , as well as a policy parameterized by ϕ , π_ϕ , by alternating policy evaluation and policy improvement steps:

$$\hat{Q}_\theta^{k+1} \leftarrow \arg \min_{\theta} \mathbb{E}_{s_t, a_t, s_{t+1} \sim \mathcal{D}} \left[\left((r(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1} \sim \pi_\phi^k(a_{t+1} | s_{t+1})} [\hat{Q}_\theta^k(s_{t+1}, a_{t+1})]) - Q(s_t, a_t) \right)^2 \right],$$

$$\hat{\pi}_\phi^{k+1} \leftarrow \arg \max_{\phi} \mathbb{E}_{s_t \sim \mathcal{D}, a_t \sim \pi_\phi^k(a_t | s_t)} \left[\hat{Q}_\theta^{k+1}(s_t, a_t) \right],$$

where $\mathcal{D} = \{(s_t, a_t, r_t, s_{t+1})\}$ refers to a dataset with states, actions, and rewards. The arg min and arg max in the equations are often approximated using gradient descent steps.

Standard Q-learning and actor-critic methods are off-policy algorithms that typically rely on some online interaction with the environment when learning.

3.2.1 Offline Reinforcement Learning

Learning to price consumer credit through trial-and-error using standard online reinforcement learning techniques can lead to costly mistakes at the detriment of both the lender and consumers. More generally, organizations may also restrict the ability to perform any sort of live pricing experiments. These restrictions on learning from live interaction motivate the need for training a pricing policy using a dataset of historical pricing decisions and outcomes.

In offline reinforcement learning, the dataset, \mathcal{D} , is gathered by a different policy than what our agent learns - in our case, a static dataset of historical pricing decisions. To denote the behavioural policy used to collect the dataset we use $\pi_\beta(a_t | s_t)$. The goal of our reinforcement learning agent then becomes to learn

a policy, π , from a dataset, \mathcal{D} , collected from behavioral policy π_β , to maximize expected return, J , when interacting with the environment at test time.

While many online reinforcement learning techniques can be used to learn from an offline dataset they often suffer from distributional shift which results in poor performance (Levine et al., 2020). Learning a policy from a static dataset requires assessing counterfactual actions that may not have been observed historically. With the absence of corrective feedback, standard online techniques may end up over-estimating the value of unseen actions and states which in turn propagate further errors as the agent uses these inaccurate estimates at test time.

Offline (or batch) reinforcement learning algorithms mitigate these issues by attempting to learn to improve on the historical policy seen in the data while reducing deviation from historical actions. This is often done through policy constraints that limit deviation from the historical policy, accounting for uncertainty in action-value estimates, or using conservative value function estimates that penalize out-of-distribution actions (Levine et al., 2020). After initial training, many offline algorithms may be combined with online fine-tuning to boost performance.

3.2.2 Conservative Q-Learning

We use the conservative Q-learning algorithm (CQL) introduced by Kumar et al. (2020) for our pricing agent. CQL is a model-free offline reinforcement learning algorithm which has achieved state-of-the-art performance on a number of offline tasks.

CQL mitigates distributional shift by penalizing values for out-of-distribution state-action tuples during training. In the implementation of CQL that we use this is done by adding an additional regularization term, $\mathbb{E}_{s \sim \mathcal{D}}[\log \sum_a \exp Q(s, a)]$, when learning the Q-function. This helps mitigate distributional shift by pushing down Q-values for action-values sampled from the current policy as they are more likely to be out-of-distribution. To avoid underestimation or producing estimates that are too conservative a second term is added to maximize the values of state-action tuples seen in the historical data, $\mathbb{E}_{s, a \sim \mathcal{D}}[Q_{\theta_i}(s, a)]$. Combined these two terms produce a conservative Q-function that is a lower bound in expectation on the true value (Kumar et al. (2020, Theorem 3.2)).

We use the soft actor-critic (SAC) (Haarnoja et al., 2018) version of the CQL algorithm implemented in the `d3rlpy` package (Seno, 2020) and use π_{CQL} to refer to the learned policy. This version modifies the soft actor-critic loss to include the additional CQL regularization term:

$$L(\theta_i) = \alpha \mathbb{E}_{s_t \sim \mathcal{D}} \left[\log \sum_a \exp Q_{\theta_i}(s_t, a) - \mathbb{E}_{s, a \sim \mathcal{D}}[Q_{\theta_i}(s, a)] - \kappa \right] + L_{SAC}(\theta_i), \quad (2)$$

where L_{SAC} refers to the soft actor-critic loss without CQL regularization. As in the original soft actor-critic paper, we use two critics parameterized by θ_i to speed up training and reduce positive bias in the policy improvement step (Haarnoja et al., 2018). The parameter α controls the degree of conservativeness and is automatically tuned by Lagrangian dual gradient descent. κ is a user-defined threshold value that helps control the magnitude of α . The $\log \sum_a \exp Q(s, a)$ term used to penalize Q-values is estimated using samples from the current policy:

$$\log \sum_a \exp Q(s, a) \approx \log \left(\frac{1}{2N} \sum_{a_i \sim \text{Unif}(a)} \left[\frac{\exp Q(s, a_i)}{\text{Unif}(a)} \right] + \frac{1}{2N} \sum_{a_i \sim \pi_\phi(a|s)} \left[\frac{\exp Q(s, a_i)}{\pi_\phi(a_i|s)} \right] \right),$$

where N refers to the user-defined number of sampled actions, and Unif is the uniform distribution.

3.3 Performance Evaluation

We are unable to test our policy through live interaction and instead evaluate our new approach using historical data. To do so, we apply our model to an out-of-sample test set and predict performance using an estimated price-response model. We compare our results against historical pricing decisions and a popular

profit-based pricing approach. In Section 4.2, we also evaluate the performance of our policy on a synthetic dataset where the true demand function is known.

3.3.1 Offline Policy Evaluation

Without knowledge of the true price-response function, $p(\text{Accept}_t|a_t, s_t)$, or the transition probabilities, $T(s_{t+1}|s_t, a_t)$, and without the ability to test our policy through live interaction it is difficult to assess model performance on the static auto loans data set. Offline policy evaluation is an area of active research with a variety of different methods that include importance sampling, models of the transitions and rewards, and value-based estimators (Paine et al., 2020; Voloshin et al., 2019). Following Phillips et al. (2015); Ban and Keskin (2021); Luo et al. (2021), we use a model-based approach and focus on estimating the price-response function in Equation (1) to predict the probability of a customer accepting a new price and use that to calculate expected reward. Although a potentially significant advantage of our approach is in treating credit pricing as a sequential decision problem, we do not develop a model for the transition dynamics, $T(s_{t+1}|s_t, a_t)$, to evaluate this.

Our primary price-response model is a logistic regression model trained on the full dataset, however, the choice of model and feature set used to estimate the price-response function impacts estimates of performance. To account for the impact of misspecification we also estimate several alternative models using parametric and non-parametric techniques to provide a range of performance.

3.3.2 Comparison

Our first point of comparison is to the historical pricing policy, π_β , captured in the data. Unfortunately, we could find no information on how pricing decisions were made by the lender at the time. However, using the data we are able to directly calculate rewards and compare total expected reward from each of the policies over the test set, see the Appendix A.2. In addition, we also measure the mean absolute percentage deviation (MAPD) in prices from the historical policy to capture how different the pricing policies are.

Our second comparison is to profit-based price optimization, an approach commonly used in practice (Phillips, 2020; Ban and Keskin, 2021; Besbes and Zeevi, 2009; Phillips et al., 2015). At a high-level, this approach involves two steps. We first estimate a price-response model to estimate how likely a potential borrower is to take out a loan at a given price point, $p(\text{Accept}_t|a_t, s_t)$. This price-response model is then used to find the greedy price that maximizes the reward function, $\max_{a_t} r(s_t, a_t)$, defined in Equation (1). We use π_{Opt} to denote this profit-based price optimization policy. An example of this approach for an arbitrary customer is shown in Appendix D. As with our CQL model this approach requires the offline policy evaluation technique described in Section 3.3.1 to estimate performance.

4 Results

4.1 Online Auto Lending

4.1.1 Data Description

For our first experiment, we use a static dataset on auto loan applications from an online lender in the United States provided by the Center for Pricing and Revenue Management.³ This dataset contains information on approximately 200,000 approved auto loan applications including the interest rate charged, the term of the loan, the approved amount, the FICO risk score, and whether the customer accepted the offer. A description of the data, features, and pre-processing steps is available in Appendix A. We also refer the reader to Phillips et al. (2015) for more detail on this dataset.

³<https://www8.gsb.columbia.edu/cprm/research/datasets>

4.1.2 The CQL Policy

We evaluate the performance of the CQL policy, π_{CQL} , on the test set using a logistic regression model to estimate the response probabilities, $p(\text{Accept}_t | s_t, a_t)$, at the new prices and calculate expected reward.⁴ Averaged over three seeds, our results indicate that the CQL agent is able to learn an effective policy that improves on historical pricing, π_β , by approximately 21% in expected profit while maintaining a less than 15% mean absolute percentage deviation in prices from the existing policy. Figure 1b shows that despite the strong overlap between π_{CQL} and π_β , the new policy pushes the average price downwards from 6.8% to 5.9%. This is consistent with the results of Phillips et al. (2015) who find evidence of historical over-pricing in this dataset due to a lack of subvented deals that were prevalent at many dealerships at the time which reduced the competitiveness of the lender’s rates.

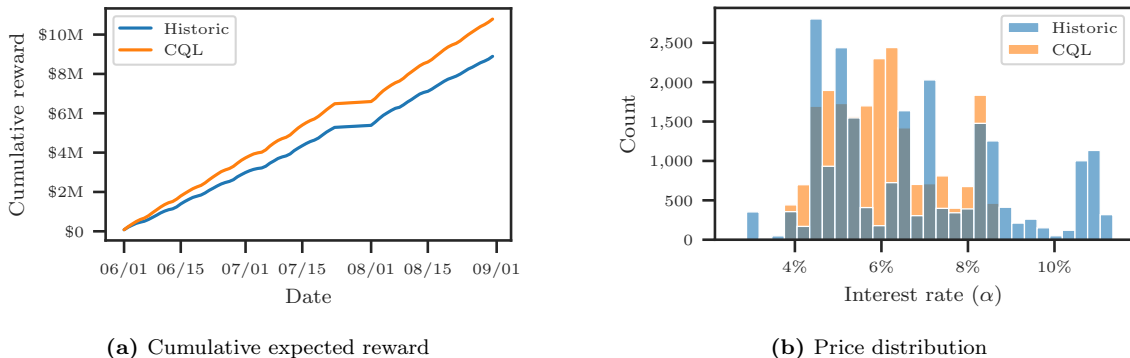


Figure 1: Performance of the CQL policy on the test set for a single seed. Cumulative expected reward is estimated using the baseline logistic regression price-response model.

We also compare the CQL model against the profit-based optimization policy, π_{Opt} , described in Section 3.3.⁵ This approach uses the same logistic regression price-response model used to evaluate performance but fitted on the training set to identify the price that maximizes the expected reward in Equation (1) for each application in the test set. With this policy we find an overall 34% increase in expected profit which is associated with an average 24% absolute percentage difference in price relative to the historical policy. These results are consistent with the results of Phillips et al. (2015) and Ban and Keskin (2021) who estimate increases of 38% and 47% on the same auto loans dataset, respectively. The differences in performance may be explained by the different test sets and the simplified reward functions they used that did not factor in credit risk or capital costs in the case of Ban and Keskin.

Figure 2 demonstrates the sensitivity of π_{Opt} performance to the assumed functional form of demand. Using alternative non-parametric and parametric price-response models with comparable classification performance, we find high variance in π_{Opt} performance with estimates ranging from a -7% decrease to a 34% increase in total expected reward relative to π_β with an average estimate of 12.6%. There is also strong evidence of overfitting with the highest performance estimated by the logistic regression used to optimize the prices. This sensitivity to the assumed functional form of demand is understated in previous literature and particularly problematic within this dataset due to the limited explainability of the feature set. For example, the baseline logistic regression model only achieves a pseudo R^2 of 0.29 and AUC of 0.83. In contrast, while π_{CQL} still exhibits sensitivity to the price-response model used for evaluating performance, the estimates lie in a narrower range from 5% to 24% with an average of 13%. These results are consistent with the work of Rana and Oliveira (2014) who test the effect of model misspecification on model-free reinforcement learning and

⁴See the Appendix for details on hyperparameters, train/val/test splits, and runtime as well as the estimated parameters of the logistic regression price-response model.

⁵We restrict prices to be between 2.5% to 12.5% to speed up optimization, however, these constraints are not binding.

parametric learning algorithms using synthetic data and find similar sensitivity of parametric methods that assume knowledge of demand.

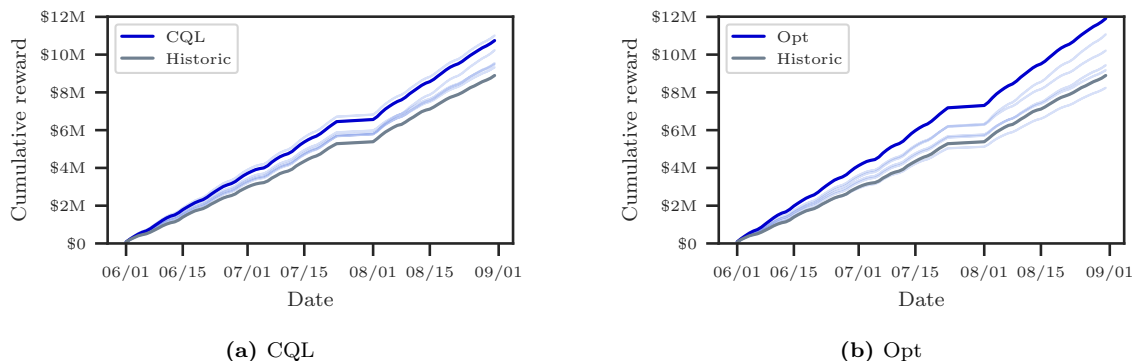


Figure 2: Cumulative expected rewards over the test set for both the CQL and Opt policies evaluated using different price-response models with a single seed. The dark blue line represents the estimated performance using the logistic regression price-response model and the lighter blue lines represent results using alternative price-response models. These include parametric and non-parametric approaches such as gradient boosting, regularized logistic regression, and neural network models.

The fact that π_{CQL} achieves performance comparable to π_{Opt} is impressive given the relatively limited change in prices from the historical policy. From Figure 3 we also observe that π_{CQL} is able to implicitly estimate the sensitivity of applicants to prices. Without knowledge of responses or rewards in the test set, it produces a policy that reduces prices for applicants who would have rejected the historical price by an average of 28% while for applicants that would have accepted it, the algorithm maintains prices on average 95% of what would have been quoted under π_{β} .

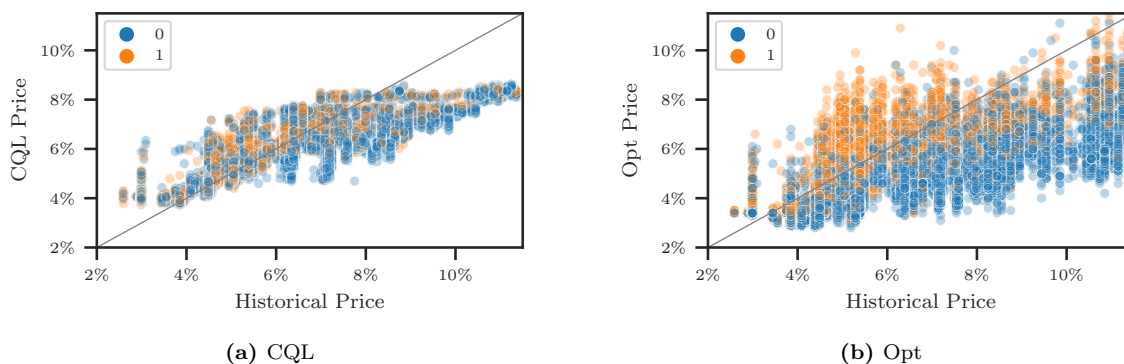


Figure 3: Comparison of π_{Opt} and π_{CQL} against π_{β} on the test set for a single seed. The blue points represent applications that did not accept the original price under π_{β} and the orange points represent loans that were funded.

In addition to the conservative improvements over the historic pricing policy, the benefits of the CQL algorithm also lie in its simplicity. Although we use estimates of the price-response curve and reward function for evaluating its policy, all that is required for training the CQL algorithm is a static dataset with features, offered prices, and rewards. In contrast, a profit-based pricing approach often requires estimation of a separate price-response model, access to the reward function, and an optimization framework to price each application.

4.1.3 The value of conservatism

We assess the effect of distributional shift on performance by fixing different values of the trade-off parameter, α , used to control the CQL regularization term in Equation (2). In the extreme case, where α is close to zero, we recover an offline version of the SAC algorithm. In Figure 4 we observe that removing the penalty on out-of-distribution actions results in substantial variation in prices compared to the historical policy as a result of over-estimating the value of unseen state-action pairs. This dramatically reduces performance and the algorithm is unable to learn an effective policy. As we increase the value α , we find that performance improves dramatically and the agent is able to learn stably.

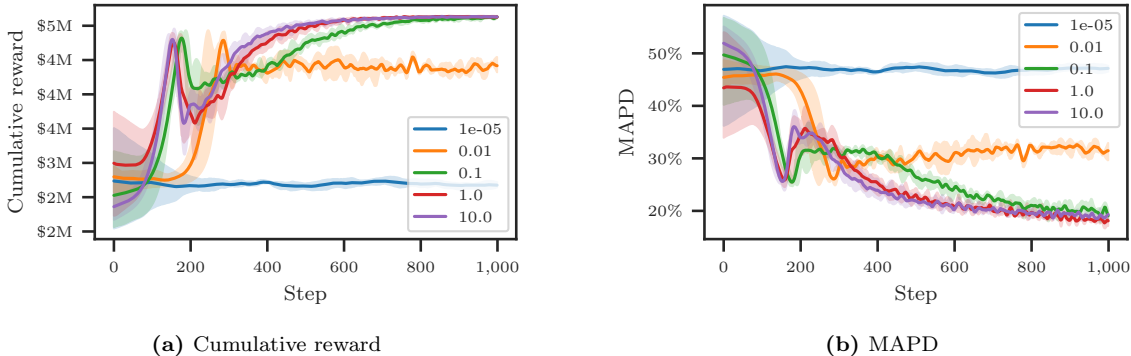


Figure 4: CQL performance on the first 10,000 test set applications as a function of training step with different fixed values of the trade-off parameter, α , averaged over three seeds. Cumulative reward is evaluated using the baseline price-response model described in Section 3.3.1.

4.2 Synthetic Price Responses

For our second set of experiments, we create synthetic datasets from the historical auto loans data using various price-response models. That is, we train a price-response model to estimate $p(\text{Accept}_t | s_t, a_t)$ on the historical dataset and use the predicted probabilities to sample from the Bernoulli distribution and replace the original accept decisions, Accept_t . We isolate California which is the largest market and has approximately 30,000 applications but otherwise follow the same approach as in the previous sections. We also compare our approach against an additional profit-based pricing model with a demand formula similar to the one introduced by Ban and Keskin (2021) that includes feature-dependent price effects (FDPE) by adding an interaction term for each of the features with price.

4.2.1 Comparison against the optimal policy

With access to the true price-response model, we are able to measure the performance of the different pricing policies against the optimal policy obtained by applying the profit-based optimization approach with the true model. In Table 1 we see the performance of the policies across different synthetic datasets. When there is no demand misspecification, i.e., the profit-based model’s price-response curve matches the dataset’s, the profit-based pricing approaches are able to achieve near-perfect performance. In real-world pricing applications, this is a strong assumption that is violated with the presence of non-stationarity, incomplete information, or complex demand behavior (Cheung et al., 2021; Luo et al., 2021). For example, when we allow the coefficients of the price-response model to vary by customer segment or use a more expressive neural network model to represent demand behavior, the profit-based models are no longer able to achieve the same level of relative performance. In these potentially more realistic scenarios, the CQL policy is able to achieve the same level of performance with approximately 50% less price variation from the historical policy. While

this stability in prices does help mitigate the risks associated with demand misspecification, it limits the effectiveness of the CQL algorithm when large variations in price are required to achieve optimal performance.

Table 1: Performance of the pricing policies on synthetic data created from five different price-response models averaged over three seeds. π_β , π_{CQL} , π_{Opt} , and $\pi_{Opt-FDPE}$ refer to the historical policy, CQL policy, profit-based optimization policy, and profit-based optimization policy with feature dependent price effects, respectively. The *logistic* dataset uses the same price-response model from π_{Opt} to generate the new accept decisions. *logistic-fdpe* includes interactions with price for each of the features. *segmented* allows the coefficients of the price-response model to vary for customer segments assigned by an unsupervised Gaussian mixture model. *time-varying* allows the coefficients of the logistic regression in the test set to vary from the training set. *neural net* uses a deep neural network as the true price-response model. MAPD is the mean absolute percentage deviation in price relative to the historical pricing policy, π_β .

Dataset	MAPD			% of Optimal Return			
	π_{CQL}	π_{Opt}	$\pi_{Opt-FDPE}$	π_β	π_{CQL}	π_{Opt}	$\pi_{Opt-FDPE}$
logistic	16.2%	24.3%	23.9%	74.5%	90.7%	99.8%	99.3%
logistic (FDPE)	16.0%	24.4%	24.8%	72.6%	87.2%	95.2%	98.9%
segmented	16.0%	25.5%	24.8%	57.4%	86.1%	83.4%	84.0%
time-varying	16.0%	25.6%	25.5%	88.1%	92.1%	91.2%	91.1%
neural net	16.1%	23.0%	26.3%	76.9%	84.0%	83.1%	82.9%

4.3 Ethical considerations

Several key issues should be considered before deploying a pricing policy based on our proposed method. In particular, our pricing agent is trained on a dataset of historical pricing decisions and may learn to recreate any biases present within that data. This may need to be addressed with additional constraints or appropriate data cleansing. Care must also be taken when selecting features such that they do not unfairly target protected groups either directly or by proxying for protected characteristics.

5 Conclusion and Future Work

In this paper, we propose a model-free offline reinforcement learning approach to pricing consumer credit. This approach makes no assumptions on the functional form of demand and introduces a formulation in which actions may impact the future state of the environment and rewards. We also demonstrate using synthetic and real data that this approach is able to improve on the existing pricing policy while showing robustness to misspecification of the underlying demand behaviour.

One logical extension of our work is to consumer goods as well as other types of consumer debt products. Future work may also seek to extend the action space to include the underwriting decision as well as incorporate capital and maximum portfolio risk requirements. In addition, as the rewards are not realized until the end of the loan term another challenge may be to develop an approach allowing for future customer behavior to affect a previously learned policy.

Acknowledgements

Raad Khraishi’s research is currently funded by NatWest Group. We thank Greig Cowan, Graham Smith, and Zachery Anderson for their valuable feedback and support. We would also like to thank Devesh Batra for his feedback on earlier drafts.

References

- Naoki Abe, Melissa Kowalczyk, Mark Domick, Timothy Gardinier, Prem Melville, Cezar Pendus, Chandan K. Reddy, David L. Jensen, Vince P. Thomas, James J. Bennett, Gary F. Anderson, and Brent R. Cooley. Optimizing debt collections using constrained reinforcement learning. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '10*, page 75–84. ACM Press, 2010. doi: 10.1145/1835804.1835817. URL <https://doi.org/10.1145/1835804.1835817>.
- Gah-Yi Ban and N. Bora Keskin. Personalized dynamic pricing with machine learning: High-dimensional features and heterogeneous elasticity. *Manage. Sci.*, 67(9):5549–5568, September 2021. ISSN 0025-1909, 1526-5501. doi: 10.1287/mnsc.2020.3680. URL <https://doi.org/10.1287/mnsc.2020.3680>.
- Hamsa Bastani, David Simchi-Levi, and Ruihao Zhu. Meta dynamic pricing: Transfer learning across experiments. *Available at SSRN 3334629*, 2019. URL <http://dx.doi.org/10.2139/ssrn.3334629>.
- Omar Besbes and Assaf Zeevi. Dynamic pricing without knowing the demand function: Risk bounds and near-optimal algorithms. *Oper. Res.*, 57(6):1407–1420, December 2009. ISSN 0030-364X, 1526-5463. doi: 10.1287/opre.1080.0640. URL <https://doi.org/10.1287/opre.1080.0640>.
- Xi Chen, Zachary Owen, Clark Pixton, and David Simchi-Levi. A statistical learning approach to personalization in revenue management. *Manage. Sci.*, January 2021. ISSN 0025-1909, 1526-5501. doi: 10.1287/mnsc.2020.3772. URL <https://doi.org/10.1287/mnsc.2020.3772>.
- Wang Chi Cheung, David Simchi-Levi, and He Wang. Dynamic pricing and demand learning with limited price experimentation. *SSRN Journal*, 65(6):1722–1731, 2014. ISSN 1556-5068. doi: 10.2139/ssrn.2457296. URL <https://doi.org/10.2139/ssrn.2457296>.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Hedging the drift: Learning to optimize under nonstationarity. *Management Science*, 2021. URL <https://doi.org/10.1287/mnsc.2021.4024>.
- Maxime C. Cohen, Ilan Lobel, and Renato Paes Leme. Feature-based dynamic pricing. *Manage. Sci.*, 66(11):4921–4943, November 2020. ISSN 0025-1909, 1526-5501. doi: 10.1287/mnsc.2019.3485. URL <https://doi.org/10.1287/mnsc.2019.3485>.
- Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning. *J. Mach. Learn. Res.*, 6:503–556, 2005. URL <https://www.jmlr.org/papers/volume6/ernst05a/ernst05a.pdf>.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1861–1870. PMLR, 10–15 Jul 2018.
- Elena Krasheninnikova, Javier García, Roberto Maestre, and Fernando Fernández. Reinforcement learning for pricing strategy optimization in the insurance industry. *Eng. Appl. Artif. Intel.*, 80:8–19, April 2019. ISSN 0952-1976. doi: 10.1016/j.engappai.2019.01.010. URL <https://doi.org/10.1016/j.engappai.2019.01.010>.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *CoRR*, abs/2006.04779, 2020. URL <https://arxiv.org/abs/2006.04779>.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *CoRR*, abs/2005.01643, 2020. URL <https://arxiv.org/abs/2005.01643>.
- Yuxi Li. Deep reinforcement learning. *CoRR*, abs/1810.06339, 2018. URL <http://arxiv.org/abs/1810.06339>.

- Yiyun Luo, Will Wei Sun, et al. Distribution-free contextual dynamic pricing. *arXiv preprint arXiv:2109.07340*, 2021. URL <https://arxiv.org/abs/2109.07340>.
- Lkhagvadorj Munkhdalai, Tsendsuren Munkhdalai, Oyun-Erdene Namsrai, Jong Lee, and Keun Ryu. An empirical comparison of machine-learning methods on bank client credit assessments. *Sustainability*, 11(3):699, January 2019. ISSN 2071-1050. doi: 10.3390/su11030699. URL <https://doi.org/10.3390/su11030699>.
- Tom Le Paine, Cosmin Paduraru, Andrea Michi, Caglar Gulcehre, Konrad Zolna, Alexander Novikov, Ziyu Wang, and Nando de Freitas. Hyperparameter selection for offline reinforcement learning. *arXiv preprint arXiv:2007.09055*, 2020. URL <https://arxiv.org/abs/2007.09055>.
- Robert Phillips and Robin Raffard. Price-driven adverse selection in consumer lending. 2011. doi: 10.13140/2.1.2901.6649. URL <http://rgdoi.net//10.13140/2.1.2901.6649>. Publisher: Unpublished.
- Robert Phillips, A Serdar Şimşek, and Garrett Van Ryzin. The effectiveness of field price discretion: Empirical evidence from auto lending. *Management Science*, 61(8):1741–1759, 2015. URL <https://doi.org/10.1287/mnsc.2014.2084>.
- Robert L Phillips. *Pricing credit products*. Stanford University Press, 2020.
- Rupal Rana and Fernando S Oliveira. Real-time dynamic pricing in a non-stationary environment using model-free reinforcement learning. *Omega*, 47:116–126, 2014. URL <https://doi.org/10.1016/j.omega.2013.10.004>.
- Takuma Seno. d3rlpy: An offline deep reinforcement library. <https://github.com/takuseno/d3rlpy>, 2020.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018. URL <http://www.incompleteideas.net/book/the-book.html>.
- Georgios Theodorou, Philip S. Thomas, and Mohammad Ghavamzadeh. Ad recommendation systems for life-time value optimization. In *Proceedings of the 24th International Conference on World Wide Web*. ACM, May 2015. doi: 10.1145/2740908.2741998. URL <https://doi.org/10.1145/2740908.2741998>.
- Margaret S. Trench, Shane P. Pederson, Edward T. Lau, Lizhi Ma, Hui Wang, and Suresh K. Nair. Managing credit lines and prices for bank one credit cards. *Interfaces*, 33(5):4–21, October 2003. ISSN 0092-2102, 1526-551X. doi: 10.1287/inte.33.5.4.19245. URL <https://doi.org/10.1287/inte.33.5.4.19245>.
- Francesco Trovò, Stefano Paladino, Marcello Restelli, and Nicola Gatti. Improving multi-armed bandit algorithms in online pricing settings. *Int. J. Approx. Reason.*, 98:196–235, July 2018. ISSN 0888-613X. doi: 10.1016/j.ijar.2018.04.006. URL <https://doi.org/10.1016/j.ijar.2018.04.006>.
- Cameron Voloshin, Hoang M Le, Nan Jiang, and Yisong Yue. Empirical study of off-policy policy evaluation for reinforcement learning. *arXiv preprint arXiv:1911.06854*, 2019. URL <https://arxiv.org/abs/1911.06854>.

A Auto Lending Dataset

We use the period from July 2002 to August 2004 which is the latest month in which all offers expire before the end of the dataset. In addition, we exclude <1,000 records with anomalous loan values less than \$5,000, missing values for state, or duplicate rows across ApproveDate, ApplyDate, Tier, FICO, CarType, State, LoanType, Rate, Amount, and Term.

A.1 Data dictionary

Table 2: Auto lending data dictionary

Field	Definition
Accept	1 for Funded, 0 for Non-funded
Rate	Customer rate
FICO	Fico score
Tier	Segmentation based on fico scores
State	Customer state
Type	Finance/Refinance
ApplyDate	Application date
ApproveDate	Approval date
DaysSinceApp	Days between application date and approval date
Term	Approved term
Amount	Loan amount approved
PreviousRate	Previous interest rate for a refinanced car
CarType	New, Used or Refinanced
CompetitionRate	Competitor's rate
PrimeRate	Prime rate
Months	Month indicator
TermClass	Segmentation based on terms
PartnerBin	Segmentation based on partners
DayOfWeek	Day of week
MonthOfYear	Month of year

A.2 Train/val/test split

Table 3: Train/val/test split. To assess out-of-sample performance we split our dataset into training, validation, and testing sets. For hyperparameter selection we use data from May 2004 which is added to the training set before assessing final performance on the test set.

Dataset	Start	End	Observations	% Accept	Avg. Reward
Train	2002-07-01	2004-04-30	161,314	21.0%	\$324
Val	2004-05-01	2004-05-31	6,482	26.1%	\$447
Test	2004-06-01	2004-08-31	21,471	25.7%	\$417

B Reproducibility

B.1 CQL Hyperparameters

Table 4: CQL Hyperparameters. We mainly use the default parameters from `d3r1py` with the exception of the hidden layers, discount rate, dropout, weight decay, and min-max scaling of the actions, features, and rewards. The hidden layers and discount rate were tuned on the validation set described in Appendix A.2. For the synthetic results in Section 4.2, we use the same hyperparameters but with only two hidden layers.

<code>n_epochs</code>	20
<code>batch_size</code>	256
<code>hidden_units</code>	[64, 64, 64, 64]
<code>n_steps</code>	1
<code>weight_decay</code>	0.0001
<code>gamma</code>	0.999
<code>alpha_threshold</code>	10
<code>conservative_weight</code>	5
<code>dropout_rate</code>	0.2
<code>n_critics</code>	2
<code>use_batch_norm</code>	False
<code>action_scaler</code>	'min_max'
<code>scaler</code>	'min_max'
<code>reward_scaler</code>	'min_max'
<code>n_action_samples</code>	10
<code>actor_learning_rate</code>	0.0001
<code>critic_learning_rate</code>	0.0003
<code>temp_learning_rate</code>	0.0001
<code>alpha_learning_rate</code>	0.0001
<code>initial_alpha</code>	1
<code>q_func_factory</code>	'mean'

B.2 Runtime

Training the algorithm on the full dataset for 20 epochs takes approximately one hour running on CPU on a 2020 MacBook Pro with a 2 GHz Quad-Core Intel Core i5 processor while a single action prediction takes approximately 640 microseconds.

B.3 Seeds

We use the following seeds in order when running our experiments: 333, 42, and 3.

C Baseline price-response model

Dep. Variable:	Accept	No. Observations:	189267			
Model:	Logit	Df Residuals:	189252			
Method:	MLE	Df Model:	14			
Date:	Sun, 09 Jan 2022	Pseudo R-squ.:	0.2946			
Time:	14:09:56	Log-Likelihood:	-69912.			
converged:	True	LL-Null:	-99105.			
Covariance Type:	nonrobust	LLR p-value:	0.000			
	coef	std err	z	P > z 	[0.025	0.975]
Intercept	44.0661	1.632	26.996	0.000	40.867	47.265
C(CarType)[T.R]	0.2901	0.046	6.242	0.000	0.199	0.381
C(CarType)[T.U]	2.4328	0.021	116.143	0.000	2.392	2.474
C(PartnerBin)[T.2]	-1.1915	0.024	-49.851	0.000	-1.238	-1.145
C(PartnerBin)[T.3]	-0.3275	0.014	-22.728	0.000	-0.356	-0.299
C(Tier)[T.2]	-0.1404	0.025	-5.710	0.000	-0.189	-0.092
C(Tier)[T.3]	-0.0351	0.034	-1.036	0.300	-0.101	0.031
C(Tier)[T.7]	0.3245	0.053	6.145	0.000	0.221	0.428
rate	-0.6599	0.010	-62.882	0.000	-0.680	-0.639
PrimeRate	0.8124	0.047	17.402	0.000	0.721	0.904
CompetitionRate	0.1619	0.024	6.744	0.000	0.115	0.209
PreviousRate	0.0021	4.78e-05	44.980	0.000	0.002	0.002
np.log(Amount)	-1.7937	0.019	-94.566	0.000	-1.831	-1.756
np.log(FICO)	-4.4786	0.243	-18.429	0.000	-4.955	-4.002
Term	0.0518	0.001	52.238	0.000	0.050	0.054

D Profit-based pricing example

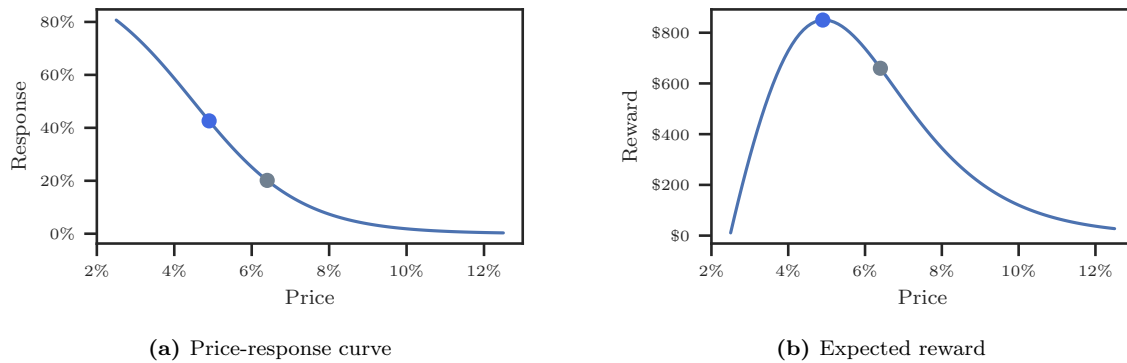


Figure 5: Profit-based price optimization for a single arbitrary customer. The plot on the left shows the probability of a customer accepting a loan at different prices estimated using logistic regression. The plot on the right uses these probabilities to estimate the expected reward at each price point. The grey dot represents the historical price and the blue dot represents the reward maximizing price.