# PAC-Bayes Meets Online Contextual Optimization

**Zhuojun Xie**
Laboratoire Génie Industriel,
CentraleSupélec, France
xie.zhuojun@centralesupelec.fr

**Adam Abdin**
Laboratoire Génie Industriel,
CentraleSupélec, France
adam.abdin@centralesupelec,fr

**Yiping Fang**[★]
Laboratoire Génie Industriel,
CentraleSupélec, France
yiping.fang@centralesupelec.fr

## Abstract

The *predict-then-optimize* paradigm bridges online learning and contextual optimization in dynamic environments. Previous works have investigated the sequential updating of predictors using feedback from downstream decisions to minimize regret in the full-information settings. However, existing approaches are predominantly frequentist, rely heavily on gradient-based strategies, and employ deterministic predictors that could yield high variance in practice despite their asymptotic guarantees. This work introduces, to the best of our knowledge, the first *Bayesian online contextual optimization* framework. Grounded in PAC-Bayes theory and general Bayesian updating principles, our framework achieves $\mathcal{O}(\sqrt{T})$ regret for bounded and mixable losses via a Gibbs posterior, eliminates the dependence on gradients through sequential Monte Carlo samplers, and thereby accommodates nondifferentiable problems. Theoretical developments and numerical experiments substantiate our claims.

## 1 INTRODUCTION

### 1.1 Problem Statement

We study a class of online contextual decision-making problems that have full-information feedback. In contrast to the widely studied online contextual bandits problem, which only assumes the observability of the reward or loss of the chosen action, the full-information setting reveals the outcome of all uncertainties relevant to the decision. Specifically, at each stage $t$, the decision-maker must determine a decision $z_t$ for stage-wise uncertainty $\xi_t$ of which the probability distribution $\mathbb{P}_{\xi_t}$ is unknown but correlated to observed contextual information $x_t$. The uncertainty $\xi_t$ will be realized after the decision is made. To prescribe a decision $z_t$ for all $t \in [T]$, the decision-maker relies on both the current context $x_t$ and historical information $\mathcal{H}_{t-1} = \{(x_i, z_i, \xi_i)\}_{i=1}^{t-1}$.

In this work, we address the online contextual decision-making problem with full-information by combining *predict-then-optimize* (PtO) and online learning. In particular, we extend decision-oriented learning from the frequentist to the Bayesian setting. We design a Bayesian update mechanism for the posterior over the model parameters that is aligned with the PAC-Bayes theory and the general belief updating framework, both of which are detailed in Section 2. This mechanism yields smooth online updates and a Gibbs posterior that minimizes a type of PAC-Bayes bound.

### 1.2 Related Work

**Online contextual decision-making**. Online contextual decision-making has been widely studied in Operations Research (OR), with applications ranging from inventory management (Cheung et al., 2022), medical services (Keyvanshokooh et al., 2025), pricing (Chen et al., 2025), and advertising (Ye et al., 2023). Unlike its offline counterpart, online contextual decision-making must address the dynamic uncertainty encapsulated in the sequential data stream (Hoi et al., 2021). This feature has outlined the importance of online learning for its methodological adaptability and efficiency in dynamic environments. Among the intersections of online learning and contextual decision-making, the contextual (multi-armed) bandits problem—where only the reward of the chosen action is observed—has developed the most extensive body of literature (Li et al., 2010). We refer readers to the work by Bietti et al. (2021) and Chen et al. (2021) and references therein for more details on contextual bandits. In this work, we focus on the online contextual optimization with full information. We study this setting especially for its ability to use rich optimization models for decisions

and ML models for prediction. While a considerable body of bandits literature extends beyond linear reward functions (Li et al., 2010; Krause and Ong, 2011; Li et al., 2017; Chowdhury and Gopalan, 2017; Zhou et al., 2020), these works typically emphasize the exploration–exploitation trade-off within relatively simple decision-making problems. These typically focus on discrete actions, knapsacks with stage-wise constraints, or long-term resource constraints (Slivkins et al., 2024; Pacchiano et al., 2025; Castiglioni et al., 2022). By contrast, the full-information setting enables decision-makers to leverage powerful predictive and prescriptive models to better exploit problem structure when uncertainty is well-estimated. More broadly, it generalizes the bandit framework and encompasses online variants of offline contextual optimization (Sadana et al., 2025), in which the data stream is dynamic and non-stationary.

**Decision-oriented offline PtO**. It is widely acknowledged that learning a predictive model using conventional statistical criteria, such as mean squared error (MSE) or maximum likelihood estimation (MLE), may cause a misalignment between model training and inference (Kong et al., 2022). To address this, the seminal frameworks of *smart predict-and-optimize* (Elmachtoub and Grigas, 2022) and *decision-focused learning* (Mandi et al., 2024, DFL), suggested a structured learning framework that optimizes predictive models directly for decision performance. Subsequent work introduced approaches to mitigate the computational challenges in deriving the gradient of decisions w.r.t. problem parameters, enabling integration with standard gradient-based training methods for parametric models. We refer the reader to the survey by Mandi et al. (2024) for a comprehensive overview. Within this research stream, optimization problems involving uncertainty in the constraints or nonlinear, nonconvex constraints/objectives have received limited attention, largely due to the technical difficulty of deriving gradients. Moreover, applications were limited to Frequentist methods, where a single predictive model is trained from large datasets. In contrast, our framework relaxes the strict reliance on gradients, thereby accommodating a broader range of decision-making structures as long as forward computation remains feasible.

**Online PtO**. Much of the online contextual decision-making literature extends from the offline setting, where a parametric model is iteratively updated in a PtO, or estimate-then-optimize manner. In the bandits setting, for example, MSE is often used to update a deterministic reward model, typically with $\ell_1$- (Chu et al., 2011) or $\ell_2$-regularization (Zhou et al., 2020). In a Bayesian approach with Thompson sampling, the

posterior of the reward function parameters is often updated by MLE with a specified likelihood (Agrawal and Goyal, 2013; Kassraie and Krause, 2022; Clavier et al., 2024). However, both MSE and MLE serve only as indirect criteria for decision-making and can lead to suboptimal decisions under considerable model misspecification (Vaswani et al., 2023; Elmachtoub et al., 2025).

In the full-information setting, model updates are typically performed using online gradient descent (OGD) or mirror descent. For example, Lobos et al. (2021) proposed a dual mirror descent framework for joint online learning and optimization, but did not explicitly consider the effect of model learning on decision quality. Two works are particularly related to our study. Liu and Grigas (2022) studied a knapsack problem with unknown reward and consumption, where covariates were available. Rather than updating the predictive model with MSE, they employed the SPO+ loss (Elmachtoub and Grigas, 2022) within a dual mirror descent framework, thereby directly accounting for the impact of prediction on decision-making. More recently, Capitaine et al. (2025) extended contextual linear programs to the online setting with uncertainty in the cost vector, incorporating a differentiable regularization term to enable gradient-based updates within an OGD framework. Our approach infers parameters via a Gibbs posterior, enabling coherent prior–posterior updates and avoiding strict reliance on gradients, which supports non-linear and nondifferentiable models.

Our contributions are summarized as follows: (1) We introduce a Bayesian online contextual optimization (BOCO) framework that unifies general Bayesian updating with PAC-Bayes in the PtO setting. Using a Gibbs posterior and aggregation, we obtain an $\mathcal{O}(\sqrt{T})$ regret guarantee under bounded, mixable losses. (2) We develop a practical PAC-Bayes SMC scheme with any-time gradient-free weight updates and Liu–West rejuvenation, yielding a Metropolis–Hastings ratio that depends only on incremental loss. (3) On a non-differentiable knapsack with uncertain weights, the aggregated approach delivers higher, more stable reward and feasibility than Gibbs stochastic prediction and deterministic PtO/DFL baselines, with especially strong gains in early, data-scarce stages.

## 2 PRELIMINARIES

In this section, we introduce two critical ingredients of our framework: the general Bayesian updating and the PAC-Bayes learning. We review these topics in the context of data-driven contextual optimization. For clarity, the discussion is presented primarily in the

offline setting, with extensions to online setting presented in Section 3.

## 2.1 General Bayesian Updating

We first review the basic Bayesian updating rule. Given a model parameter space $\Theta \ni \theta$, a covariate space $\mathcal{X} \ni x$, an uncertainty space $\Xi \ni \xi$, a predictive model class $\mathcal{M} := \{m(\cdot; \theta) : \mathcal{X} \to \Xi, \theta \in \Theta\}$, and a batch of observations $\mathcal{D} = \{(x_i, \xi_i)\}$, the Bayes rule is given by:

$$\pi(\theta) = \frac{l(\theta|\mathcal{D})\pi_0(\theta)}{\int_\Theta l(\theta'|\mathcal{D})\pi_0(\theta')\mathrm{d}\theta'} \quad (1)$$

where $\pi_0$ and $\pi$ are the prior and posterior distribution over $\Theta$, and $l(\theta|\mathcal{D}) = \prod_{(x_i, \xi_i) \in \mathcal{D}} p_\theta(\xi_i|x_i)$ where $p_\theta(\cdot|x)$ is the probability density (or mass function) of $\xi$ under model $m(\cdot; \theta)$, denoting the likelihood of parameters $\theta$ given observations $\mathcal{D}$. The likelihood is central for Bayesian inference, as it directly reflects the modeler's *acceptance* of specific model parameters given observations. Often, an analytical likelihood is adopted, such as Gaussian likelihoods for general-purpose modeling (Kassraie and Krause, 2022) or right-censored Weibull distributions for demand modeling (Chuang and Kim, 2023). However, reliance on analytical likelihoods inevitably risks model misspecification when the true data-generating process cannot be captured by such specifications. Although there is extensive literature on likelihood-free Bayesian inference, these approaches still require an appropriate definition of *proximity* between the simulated and observed data (Thomas et al., 2022).

For data-driven optimization, it is helpful to adopt a philosophy similar to frequentist decision-oriented learning when interpreting the likelihood in Bayesian inference. That is, parameters should be judged by *prescriptiveness*—the quality of the decisions they induce. In contrast to a probability density-based likelihood, prescriptiveness has two important features. Firstly, it is tractable, though non-analytical, to compute by passing parameters through the PtO pipeline. Secondly, it provides a natural one-dimensional summary statistic reflecting the decision-making performance of the parameters.

Task-oriented likelihood constructions are well established in Bayesian statistics. For instance, Ibrahim and Chen (2000) introduced a power factor for the likelihood term, such as $l(\theta|D)^\alpha$, for robustness. Jiang and Tanner (2008) proposed a classification error-based criterion to mitigate the impact of a misspecified likelihood, and suggested a risk-based posterior updating. Other constructions were developed thereafter, among which we highlight the work by Bissiri et al. (2016)

that systematically proposed the general Bayesian updating which most directly motivates our work. General Bayesian updating addresses likelihood misspecification when a task loss is available, an insight that closely aligns with the decision-oriented learning in data-driven contextual optimization. The general updating rule can be framed in Eq (2) as a variant of Eq (1):

$$\pi(\theta) = \frac{e^{-L(\theta|\mathcal{D})}\pi_0(\theta)}{\int_\Theta e^{-L(\theta'|\mathcal{D})}\pi_0(\theta')\mathrm{d}\theta'} \propto e^{-L(\theta|\mathcal{D})}\pi_0(\theta), \quad (2)$$

where $L(\theta|\mathcal{D})$ denotes the task loss. General Bayesian updating replaces the likelihood $l(\theta|\mathcal{D})$ by a pseudo-likelihood $e^{-L(\theta|\mathcal{D})}$, resulting in a Gibbs posterior.

Computationally, a Gibbs posterior enables gradient-free Bayesian inference techniques, thereby generalizing decision-oriented learning to other nondifferentiable optimization problems. While no analytical form exists for such a pseudo-likelihood, Metropolis–Hastings Markov Chain Monte Carlo (MCMC) only requires forward computation of the acceptance probability. Additionally, approximate Bayesian computation is applicable by using the task loss $L(\theta|\mathcal{D})$ as a summary statistic and accept/reject $\theta$ based on proximity to the empirical optimum $\min_{\theta \in \Theta} L(\theta|\mathcal{D})$. On the other hand, if the gradient $\nabla_\theta L(\theta|\mathcal{D})$ is available, gradient-based techniques can be considered to improve computational efficiency, such as Hamiltonian MC and neural variational inference (Mnih and Gregor, 2014).

General Bayesian updating is a principled way to update the belief for a specific task under model misspecification. In practice, such a posterior can be interpreted as the outcome of learning from finite data given a specific loss. In data-driven settings, since the Gibbs posterior is always accompanied with a predictive model $m(\cdot; \theta)$, the generalization ability of such a combination for making good decisions/predictions is of interest in OR/ML. Section 2.2 provides a theoretical support for using the Gibbs posterior in learning tasks.

## 2.2 PAC-Bayes Learning

We refer readers to Alquier (2024) for a modern overview of PAC-Bayes. PAC, short for *Probably Approximately Correct*, is a theoretical framework that provides probabilistic guarantees on the generalization error of a learning algorithm, ensuring that with high probability the error remains close to its expected value when trained on finite data. PAC-Bayes learning particularly focuses on the generalization error of a model under *any* posterior of parameters incorporating a prior distribution. To provide a concrete example

and to motivate our approach, we present a well-known PAC-Bayes bound for bounded loss. Recent extensions to loss with more general tail behaviors can be found in Rodriguez-Galvez et al. (2024).

Alquier (2024, Theorem 2.1): Suppose the data is i.i.d. collected from $\mathbb{P}$ and loss is bounded in $[0, C]$. Then the following inequality holds with probability at least $1 - \delta$ over the draw of data for any $\lambda > 0$, any $\delta \in (0, 1)$, any data-independent prior $\pi_0 \in \mathcal{P}(\Theta)$, and any posterior $\pi \in \mathcal{P}(\Theta)$:

$$\mathbb{E}_{\theta \sim \pi}[R(\theta)] \leq \mathbb{E}_{\theta \sim \pi}[r(\theta)] + \frac{\lambda C^2}{8n} + \frac{\mathbb{D}_{KL}(\pi \| \pi_0) + \log \frac{1}{\delta}}{\lambda},$$

where $R(\theta)$ and $r(\theta)$ denote the true and empirical risk, respectively. $\mathcal{P}(\Theta)$ denotes the set of all probability measures over $\Theta$. In particular, the expectation $\mathbb{E}_{\theta \sim \pi}[R(\theta)]$ corresponds to the true risk of the Bayesian stochastic predictor, showing the generalization ability of randomly drawn predictive model $m(\cdot; \theta), \theta \sim \pi$. Since this upper bound is arbitrary for $\pi$, the practical value of PAC-Bayes lies in the optimization of $\pi$ for minimizing this bound. According to Alquier (2024, Corollary 2.3), the minimizer takes the Gibbs form:

$$\frac{e^{-\lambda r(\theta)} \pi_0}{\int_\Theta e^{-\lambda r(\theta')} \pi_0(\mathrm{d}\theta')} =$$
$$\arg\min_{\pi \in \mathcal{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim \pi}[r(\theta)] + \frac{\mathbb{D}_{KL}(\pi \| \pi_0)}{\lambda} \right\},$$

which coincides with the Gibbs posterior in the general Bayesian updating with an extra scalar $\lambda$ that controls the trade-off between the discrepancy $\mathbb{D}_{KL}(\pi \| \pi_0)$ and the pseudo-likelihood. The discrepancy $\mathbb{D}_{KL}(\pi \| \pi_0)$, i.e., the relative entropy between posterior and prior, controls the complexity of the posterior and can be replaced by alternative measures (Amit et al., 2022). Other discrepancies generally do not preserve the optimality of the Gibbs posterior. Nonetheless, Bissiri et al. (2016) showed that relative entropy is the unique choice that preserves coherent inference, a crucial property for online learning to guarantee any-time validity. Therefore, we stick to the relative entropy for our online framework due to the coherence and decision-oriented explanation of the Gibbs posterior under the general Bayesian updating.

## 3 BOCO

In this section, we introduce our Bayesian Online Contextual Optimization (BOCO) framework, including its definition, theoretical properties, and a practical algorithm.

### 3.1 Framework

Following the notations in Section 1, we recall that at each stage $t$, the decision-maker observes the covariate $x_t$ correlated to the uncertainty $\xi_t$, which further characterizes a parametric optimization problem $\mathbf{P}(\xi)$. Specifically, we assume the problem $\mathbf{P}(\xi)$ has the exact structure to reflect the real-world objective and constraints in the full-information setting. For generality, we write the problem as:

$$\mathbf{P}(\xi) = \arg\min_{z \in g(\xi)} c(z; \xi), \tag{3}$$

where $c(z; \xi)$ denotes the uncertainty-related objective function and $g(\xi)$ the uncertainty-related feasible set for decision-making. This formulation allows the uncertainty to be in the objective and/or constraints. In the experiments, we focus on a hard-constrained integer knapsack problem with uncertainty in the weights.

Consider a decision-maker who employs a parametric model $m \in \mathcal{M}$ to predict the uncertainty in each stage $t$ given covariate $x_t$, with an initial data-independent prior $\pi_0 \in \mathcal{P}(\Theta)$. For any $t \in [T]$, the decision-making and learning proceed as follows:

$$\hat{z}_t = \mathbf{P}(m(x_t; \pi_t)), \tag{4}$$
$$\pi_{t+1}(\theta) \leftarrow \frac{e^{-\lambda \ell(\theta, d_t)} \pi_t(\theta)}{\int_\Theta e^{-\lambda \ell(\theta', d_t)} \pi_t(\theta') \mathrm{d}\theta'}, \tag{5}$$

where $\lambda > 0$ controls relative importance between prior and current information, and $d_t = (x_t, \xi_t)$. Particularly, with a slight abuse of notation, we define the pushforward of $\pi_t$ by $m(x_t; \cdot)$:

$$m(x_t; \pi_t) := (m(x_t; \cdot))_\# \pi_t, \tag{6}$$

which is a probability measure over $\Xi$ induced by $\pi_t$. Furthermore, the loss function is defined as the regret:

$$\hat{z}_t(\theta) := \mathbf{P}(m(x_t; \theta)),$$
$$\ell(\theta, d_t) := c'(\hat{z}_t(\theta); \xi_t) - \min_{z \in g(\xi_t)} c(z; \xi_t), \tag{7}$$
$$\ell(\pi_t, d_t) := c'(\hat{z}_t; \xi_t) - \min_{z \in g(\xi_t)} c(z; \xi_t), \tag{8}$$

where $c'(z; \xi)$ evaluates the decision quality of a decision $z$ given uncertainty realization $\xi$, see Assumption 3.1 for example. A distinctive feature of our framework is to leverage the aggregated predictor $m(x_t; \pi_t)$ in Eq (4), rather than the stochastic predictor $m(x_t; \theta)$ with a randomly drawn $\theta \sim \pi_t$ at each stage $t$.

The aggregated predictor is common in cost-sensitive classification which would be subject to high variance from the stochastic predictor. Notably, Lacasse et al.

(2006) and following work considered a binary classification problem with asymmetric losses, and introduced the *C-Bound* for the aggregated predictor induced by the Gibbs posterior. Such aggregation strategies have been applied to the multi-armed bandits setting which shares a similar structure with multiclass classification. For instance, Sakhi et al. (2023) designed a posterior policy, which was a mixture of deterministic and parametric decision rules, to assign probabilities over finite actions.

The majority vote hardly generalizes to problems with a general decision space. By defining the pushforward measure $m(x_t; \pi_t)$ in Eq (6), our framework allows the utilization of stochastic optimization and its variants to address the uncertainty encoded in the posterior, which assigns *weights* $\pi_t(\theta)$ to many *experts* $m(\cdot; \theta)$. This procedure aligns with the Bayesian decision theory, and still balances the exploitation-exploration in online learning. In this case, exploitation corresponds to choosing the maximum a posterior parameters, while exploration suggests a stochastic predictor $m(x_t; \theta_t), \theta_t \sim \pi_t$ for PtO in each stage.

### 3.2 Guarantees

Two theoretical properties of our Gibbs posterior $\pi_t$ in Eq (5) warrant further discussion. First, we establish that this update is valid in the online learning setting, i.e., it is consistent with the general Bayesian updating rule and achieves optimality under a specific criterion. Second, we derive the regret bound for such a posterior in online PtO under mild assumptions. Our theoretical analysis is built on the work by Haddouche and Guedj (2022), to which we refer readers for more details. Note that they considered a binary classification and a linear regression task, while we extend to general parametric optimization problems.

We denote the data space $\mathscr{D} := (\mathcal{X} \times \Xi)^T$, $D = \{(x_1, \xi_1), \ldots, (x_T, \xi_T)\}$ a random sample drawn from $\mathscr{D}$ by some probability rule $\mu$, and $d_t = (x_t, \xi_t)$ a pair observed at time $t$ (i.e., a realization of $(X_t, \boldsymbol{\xi}_t)$). Furthermore, we define the filtration $\mathbb{F}$ for a finite horizon $T$, that is, $\mathbb{F} := (\mathcal{F}_t)_{t \in [T]} = (\sigma(\mathcal{H}_i)|i \leq t)$, in which $\mathcal{F}_t$ is a $\sigma$-algebra over the historical information $\mathcal{H}_t$. This filtration is mainly used to define a type of regret for online learning that differs from the widely-used static and dynamic regret. Next, we make the following assumption on the evaluation function $c'$:

**Assumption 3.1** (Bounded loss). $\exists C > 0$, such that

(a) $\forall \xi \in \Xi, \forall z \notin g(\xi), c'(z; \xi) = \max_{a \in g(\xi)} c(a; \xi)$,

(b) $\forall \xi \in \Xi, \forall z \in g(\xi), c'(z; \xi) = c(z; \xi) \leq C$.

**Theorem 3.2** (Haddouche and Guedj (2022), Corollary 3.1). *Suppose Assumption 3.1 holds. For any dis-*

*tribution $\mu$ over $D$, any $\lambda > 0$, any $\delta \in (0, 1)$, and any online posterior $\{\tilde{\pi}\}$ and prior $\{\pi\}$ sequences, the following inequality holds with at least probability $1 - \delta$ over the draw $D \sim \mu$:*

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \tilde{\pi}_{t+1}}[\mathbb{E}[\ell(\theta, \boldsymbol{\xi}_t) | \mathcal{F}_{t-1}, x_t]]$$
$$\leq \sum_{t=1}^{T} \left( \mathbb{E}_{\theta \sim \tilde{\pi}_{t+1}}[\ell(\theta, d_t)] + \frac{\mathbb{D}_{KL}(\tilde{\pi}_{t+1} \| \pi_t)}{\lambda} \right)$$
$$+ \frac{\lambda T C^2}{8} + \frac{\log(1/\delta)}{\lambda}.$$

Here $\ell(\theta, \boldsymbol{\xi}_t)$ denotes the random loss of $\theta$ given $\boldsymbol{\xi}_t$ conditioning on $\mathcal{F}_{t-1}$ and $x_t$. The online posterior sequence $\{\tilde{\pi}\}$ denotes the posterior $\tilde{\pi}_t$ is $\mathcal{F}_{t-1}$-measurable, depending only on $\tilde{\pi}_{t-1}$ and $(x_{t-1}, \xi_{t-1})$. For each stage $t$ and fixed parameters $\theta$, the target (l.h.s.) is defined as $\mathbb{E}[\ell(\theta, \boldsymbol{\xi}_t) | \mathcal{F}_{t-1}, x_t]$ rather than $\min_{\theta \in \Theta} \ell(\theta, d_t)$ such that, the decision-maker does not seek for an unrealistic oracle that achieves optimality for any stochastic realization $\xi_t \sim \boldsymbol{\xi}_t$, but a pragmatic approach to minimize the expected loss that is achievable when only information $\{\mathcal{H}_{t-1}, x_t\}$ is available.

Theorem 3.2 provides a post-hoc criterion to update the belief, considering the decision $\hat{z}_t$ is always made with prior $\pi_t$ before knowing $\xi_t$ in online PtO. Once $(x_t, \xi_t)$ is realized, and the probability distribution of $\boldsymbol{\xi}_t$ conditioning on $(x_t, \mathcal{H}_{t-1})$ is known, the decision-maker would update the prior to a posterior that optimizes the regret, and leverages this posterior for the future. To obtain the posterior for each stage $t$, we optimize the r.h.s. that relates to the posterior:

$$\tilde{\pi}_{t+1} := \arg\min_{\pi \in \mathcal{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim \pi}[\ell(\theta, d_t)] + \frac{\mathbb{D}_{KL}(\tilde{\pi} \| \pi_t)}{\lambda} \right\}$$

which justifies the updating rule in Eq (5). Next, we show the regret bound for implementing this posterior from stage $t$ as the prior for stage $t + 1$.

**Theorem 3.3** (Haddouche and Guedj (2022), Corollary 3.3). *Suppose Assumption 3.1 holds. For any distribution $\mu$ over $D$, any $\lambda > 0$, any $\delta \in (0, 1)$, and any online posterior sequences $\{\pi_t\}$, the following inequality holds with at least probability $1 - \delta$ over the draw $D \sim \mu$:*

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t}[\mathbb{E}[\ell(\theta, \boldsymbol{\xi}_t) | \mathcal{F}_{t-1}, x_t]]$$
$$\leq \sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t}[\ell(\theta, d_t)] + \mathcal{O}\left( \sqrt{\log(1/\delta) C^2 T} \right)$$

*where optimal $\lambda = \sqrt{\frac{8 \log(1/\delta)}{T C^2}}$ is adopted.*

Since Theorem 3.3 applies to any posterior sequence, it also applies to the Gibbs posterior in Eq (5) which admits the optimal updating rule by Theorem 3.2.

We emphasize that the optimality of the Gibbs posterior is w.r.t. the risk of stochastic predictor in Theorem 3.2. Additionally, the guarantee in Theorem 3.3 applies to the stochastic predictor $m(x_t; \theta)$, not the aggregated predictor $m(x_t; \pi_t)$. It is possible that, if the true risk and empirical risk in Theorem 3.2 are defined for the aggregated predictor, the Gibbs posterior may not retain the optimality. We stick to the Gibbs posterior updating in Eq (5) for two reasons. First, it aligns with the general Bayesian updating principle. Second, it admits a point-wise updating process without using distributional information in $\pi_t$, thereby reducing computational effort, allowing individual evolutions of parameters $\theta$ as the algorithms we propose in Section 3.3.

We aim to modify Theorem 3.3 to construct a generalization bound for the aggregated predictor $m(x_t; \pi_t)$ for each stage $t$. Define the target risk for $\pi_t$ as:

$$R_t(\pi_t) := \mathbb{E}[\ell(\pi_t, \boldsymbol{\xi}_t) | \mathcal{F}_{t-1}, x_t]. \tag{9}$$

The target risk in Eq (9) differs from the l.h.s. in Theorem 3.3 by evaluating the decision $\hat{z}_t$ optimized for distribution $m(x_t; \pi_t)$ under conditional distribution of $\boldsymbol{\xi}_t$, rather than the posterior expectation for $\hat{z}_t(\theta)$. We leverage the following mixability assumption on the decision-making and predictive model to provide a generalization bound for our BOCO framework with limited modification to Theorem 3.3.

**Assumption 3.4** ($\lambda$-mixable loss)**.** *The loss function $\ell$ is $\lambda$-mixable given $\mathcal{M}$, i.e., $\forall \lambda > 0, \forall d \in \mathcal{X} \times \Xi, \forall \pi \in \mathcal{P}(\Xi)$,*

$$\ell(\pi, d) \le -\frac{1}{\lambda} \log \mathbb{E}_{\theta \sim \pi}[e^{-\lambda \ell(\theta, d)}].$$

The mixability is a standard assumption for online aggregation and it has a natural justification in data-driven optimization. This property formalizes why distribution-aware (stochastic) optimization can outperform plug-in decisions that commit to a single scenario. With Assumption 3.4, we derive an extension of Theorem 3.3 for the aggregated predictor.

**Corollary 3.5.** *Suppose Assumption 3.1 and 3.4 hold. For any distribution $\mu$ over $D$, any $\lambda > 0$, any $\delta \in (0, 1)$, and any online posterior sequences $\{\pi_t\}$, the following inequality holds with at least probability $1 - \delta$ over the draw $D \sim \mu$:*

$$\sum_{t=1}^{T} R_t(\pi_t) \le \sum_{t=1}^{T} \mathbb{E}_{\theta \sim \pi_t}[\ell(\theta, d_t)] + \mathcal{O}\left(\sqrt{\log(1/\delta)C^2 T}\right)$$

*where optimal $\lambda = \sqrt{\frac{8\log(1/\delta)}{TC^2}}$ is adopted.*

Corollary 3.5 achieves the same rate for the aggregated predictor as the one for stochastic predictor in Theorem 3.3. The proof mainly depends on bounding the risk of aggregated predictor by that of the Gibbs stochastic predictor. We note that this bound may be vacuous in practice, especially in case where severe uncertainty exists and stochastic optimization provides better decision than a point-based deterministic optimization. We demonstrate this phenomenon in the experiments when the aggregated predictor remarkably outperforms the stochastic predictor.

### 3.3  Algorithms

We propose a practical sequential Monte Carlo (SMC) sampler to approximate the posterior $\pi_t$ for nondifferentiable optimization problems that only allow objective evaluation. Our algorithm, summarized in Algorithm 1, follows a classic SMC procedure with MCMC-based rejuvenation steps after importance sampling. We highlight the techniques we utilized to mitigate the computational challenges in rejuvenation using a Liu-West kernel density estimator (Liu and West, 2001). The estimator is used to approximate current posterior $\hat{\pi}_t$, and to act as an independent proposal distribution for drawing proposed parameters $\theta'$ in the MCMC process.

Consider the original definition of acceptance ratio for proposed parameters $\theta'$ given current parameters $\theta$ and posterior $\pi_t$:

$$r = \frac{\pi_{t+1}(\theta')q(\theta|\theta')}{\pi_{t+1}(\theta)q(\theta'|\theta)} = \frac{\pi_t(\theta')e^{-\lambda\ell(\theta',d_t)}q(\theta|\theta')}{\pi_t(\theta)e^{-\lambda\ell(\theta,d_t)}q(\theta'|\theta)}, \tag{10}$$

where $q$ is a proposal distribution. In the SMC, it is challenging to evaluate the density under $\pi_t$, which has no analytical form. Therefore, a practical approach is to consider the expansion of $\pi_t$, leading to:

$$r = \frac{\pi_0(\theta')e^{-\lambda \sum_{i=0}^{t} \ell(\theta',d_i)}q(\theta|\theta')}{\pi_0(\theta)e^{-\lambda \sum_{i=0}^{t} \ell(\theta,d_i)}q(\theta'|\theta)}. \tag{11}$$

Although Eq (11) allows the computation of acceptance ratio, it requires to evaluate the decision quality of parameters $\theta$ and $\theta'$ over all historical instances up to stage $t$, this can be very time-consuming when number of particles $N$ and MCMC steps $L$ are large. Therefore, we adopt the Liu-West Gaussian mixture estimator $q_t$ defined in Eq (12) to approximate the current posterior $\pi_t$, and regard $q_t$ as the independent

proposal distribution. This leads to:

$$q_t = \sum_{i=1}^{N} w_i^t \mathcal{N}(m_t^i; H_t) \tag{12}$$

$$r = \frac{\pi_t(\theta')e^{-\lambda\ell(\theta',d_t)}q_t(\theta)}{\pi_t(\theta)e^{-\lambda\ell(\theta,d_t)}q_t(\theta')} \approx \frac{q_t(\theta')e^{-\lambda\ell(\theta',d_t)}q_t(\theta)}{q_t(\theta)e^{-\lambda\ell(\theta,d_t)}q_t(\theta')}$$
$$= e^{-\lambda\ell(\theta',d_t)+\lambda\ell(\theta,d_t)}, \tag{13}$$

which cancels the posterior $\pi_t$ and only requires the incremental loss from $\theta$ to $\theta'$ as long as $\theta'$ is drawn from $q_t$. Compared to Eq (11), this approach only requires one decision quality evaluation for proposed parameters $\theta'$. Moreover, it can be seen as that, if parameters $\theta$ and $\theta'$ have the same density in the product of *prior* and proposal, the posterior $\pi_{t+1}$ should move towards the one with better decision quality.

---

**Algorithm 1:** PAC-Bayes Sequential Monte Carlo

**Input:** $\pi_0, a, \tau, \lambda, L, N, \ell, T$.

Initialize $\forall i \in [N], \theta_0^i \sim \pi_0, w_0^i = 1/N$;

**for** $t \leftarrow 0$ **to** $T$ **do**

  $\hat{\pi}_t \leftarrow \sum_{i=1}^{N} w_t^i \delta_{\theta_t^i}$;

  $\hat{z}_t \leftarrow \mathbf{P}(\hat{\pi}_t)$;

  $\ell_t \leftarrow c'(\hat{z}_t, \xi_t)$;

  $\tilde{w}_t^i \leftarrow w_t^i e^{-\lambda\ell(\theta_t^i, d_t)}, \forall i \in [N]$;

  $w_t^i \leftarrow w_t^i / \sum_{j=1}^{N} \tilde{w}_t^j, \forall i \in [N]$;

  **if** $1/\sum_{i=1}^{N}(w_t^i)^2 \leq \tau N$ **then**

    $\bar{\theta}_t \leftarrow \sum_{i=1}^{N} w_t^i \theta_t^i$;

    $\Sigma_t \leftarrow \sum_{i=1}^{N} w_t^i(\theta_t^i - \bar{\theta}_t)(\theta_t^i - \bar{\theta}_t)^\top$;

    $m_t^i \leftarrow a\theta_t^i + (1-a)\bar{\theta}, \forall i \in [N]$;

    $H_t \leftarrow (1-a^2)\Sigma_t$;

    **for** $i \leftarrow 1$ **to** $N$ **do**

      **for** $k \leftarrow 1$ **to** $L$ **do**

        $I \leftarrow \mathrm{Cat}(w_t^1, \ldots, w_t^N)$;

        $\theta' \sim \mathcal{N}(m_t^I, H_t)$;

        $r \leftarrow e^{-\lambda\ell(\theta',d_t)+\lambda\ell(\theta_t^i,d_t)}$;

        **if** $r \geq u \sim Uniform[0,1]$ **then**

          $\theta_t^i \leftarrow \theta'$;

    $w_t^i \leftarrow 1/N, \forall i \in [N]$;

  $\theta_{t+1}^i \leftarrow \theta_t^i, \forall i \in [N]$;

  $w_{t+1}^i \leftarrow w_t^i, \forall i \in [N]$;

**Output:** $\{\ell_t\}_{t=0}^{T-1}$

---

## 4 Experiments

In this section, we test our `BOCO` framework on a nondifferentiable integer knapsack problem with uncertain item weight matrix. Mathematically, a four-dimensional decision variable $z \in \mathbb{N}^4$ denotes the quantities of different items, $c \in \mathbb{R}^4$ the unit reward of

each item, $b \in \mathbb{R}^3$ the available amount of three resources, $q \in \mathbb{R}^3$ the unit salvage value of resources, and $A_t \in \mathbb{R}^{3\times 4}$ the weight matrix of items. At each stage $t$ the decision-maker observes the covariate $x_t$ and must make the decision $\hat{z}_t$ with predicted $\tilde{A}_t$. The problem can be formulated as:

$$\max_{z \in \mathbb{N}_+^4} \quad c^\top z + q^\top [b - \tilde{A}_t z]^+$$
$$\text{s.t.} \quad \tilde{A}_t z \leq b$$

where $[a]^+ = [\max\{a_i, 0\}]^\top$ for a vector $a$. The reward for decision $\hat{z}_t$ is revealed given $A_t$. In particular, if $\hat{z}_t$ is not feasible for $A_t$, the reward will be zero. Otherwise, the reward will be $r_t = c^\top \hat{z}_t + q^\top[b - A_t\hat{z}_t]^+$. In this problem, we take values $c = [12, 12, 12, 12]^\top, b = [8, 8, 8]^\top, q = [3, 3, 3]^\top$. The data-generation process for $(x_t, A_t)$ is a variant of the ARMA(2,2) proposed by Bertsimas and Kallus (2020). We reshape and rescale the demand to fit our case study.
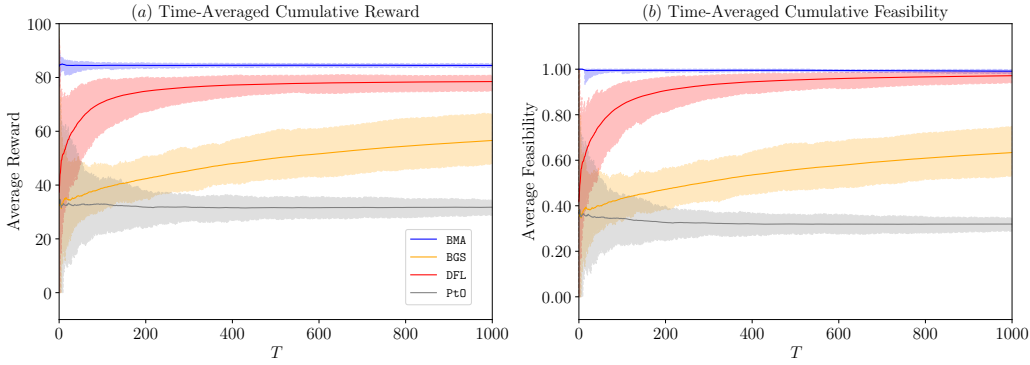
In our `BOCO` framework, we have $N$ scenarios $\{\tilde{A}_t^i\}_{i=1}^N$ generated by $N$ models $\{m(\cdot; \theta_i)\}_{i=1}^N$ for each stage $t$ given $x_t$. Therefore, we frame a chance-constrained stochastic program given the empirical posterior distribution as follows:

$$\max_{z \in \mathbb{N}_+^4} \quad \sum_{i=1}^{N} w_t^i \mathbb{I}[\tilde{A}_i z \preceq b]\{c^\top z + q^\top[b - \tilde{A}_t^i z]^+\}$$
$$\text{s.t.} \quad \sum_{i=1}^{N} w_t^i \mathbb{I}[\tilde{A}_i z \preceq b] \geq \alpha$$

where $\alpha \in (0, 1)$ is the feasibility target. We derive a deterministic equivalent mixed-integer linear program in the supplementary material to solve it. For all approaches considered in the experiments, we compared `BMA` (our `BOCO` with pushforward distribution and stochastic optimization), `BGS` (one predictive model drawn from $\pi_t$ per step), `PtO` (MSE-trained deterministic predictor), and `DFL` (decision loss-trained deterministic predictor). All experimental details are provided in the complementary material.

Figure 1 depicts the performance of different frameworks on the time-averaged cumulative reward and feasibility from 100 individual experiments with $T = 1000$. The uncertainty range is plotted using the 10-90 percentiles within the 100 experiments. We summarize the statistics of results in Table 1. In the table, $\bar{r}_{1000}$ and $\bar{\alpha}_{1000}$ denote the averaged reward and feasibility over 1000 steps, while $\bar{r}_{500}$ and $\bar{\alpha}_{500}$ denote that over the last 500 steps. First, we note the stability of `BMA` framework indicated by the uncertainty of reward and feasibility across 100 runs. We annotate the average of cumulative reward and feasibility at the last stage for each framework. Comparing `BMA`

Figure 1: Time-averaged cumulative reward and feasibility of four frameworks in 100 trials.



*Note*: The uncertainty range is plotted using 10-90 percentiles from 100 trials.

Table 1: Framework performances in average reward and feasibility

| Framework | $\bar{r}_{1000}$ | $\bar{r}_{500}$ | $\bar{\alpha}_{1000}$ | $\bar{\alpha}_{500}$ |
|---|---|---|---|---|
| BMA | $84.52 \pm 0.05$ | $84.53 \pm 0.03$ | $0.99 \pm 0.00$ | $0.99 \pm 0.00$ |
| DFL | $75.63 \pm 5.27$ | $78.16 \pm 0.24$ | $0.92 \pm 0.08$ | $0.96 \pm 0.01$ |
| BGS | $48.45 \pm 6.29$ | $53.63 \pm 1.89$ | $0.54 \pm 0.07$ | $0.60 \pm 0.02$ |
| PtO | $31.85 \pm 0.39$ | $31.67 \pm 0.05$ | $0.32 \pm 0.01$ | $0.32 \pm 0.00$ |

*Note*: The uncertainty is computed using one standard deviation from 100 trials.

with BGS highlights the benefit of combining the posterior distribution and the stochastic optimization for decision-making under uncertain environments. Although the BGS framework leverages the parameters that constitute the SMC sample, it suffers from infeasible decisions. Similarly, the DFL framework suffers from the infeasibility of decisions especially at the early stage, then improves the feasibility by making conservative predictions and leads to conservative decisions. Therefore, its cumulative reward is consistently smaller than that of BMA. Additionally, for the deterministic approaches, the initialization of model parameters has nonnegligible impact on the model's performance, while BMA mitigates such impact by leveraging distributional information. As the baseline, the PtO framework performs the worst in terms of reward and feasibility. This is due to the fact that the reward is strongly asymmetric in the prediction. Consequently, the MSE-based learning is highly misaligned with the decision-making target. Overall, these results highlight the advantages of the BOCO framework for decision-making under uncertainty. Specifically, BOCO framework demonstrates its stability and robustness especially at the early stage, when data are insufficient for the deterministic DFL approach.

## 5 Conclusions

This work proposes a Bayesian online contextual optimization (BOCO) framework that unifies general Bayesian updating with PAC-Bayes to bring principled, task-oriented learning into online predict-then-optimize. By updating beliefs via a Gibbs posterior, the method provides coherent any-time updates and achieves $\mathcal{O}(\sqrt{T})$ regret for bounded and mixable losses. Computationally, a sequential Monte Carlo sampler with Liu–West rejuvenation delivers gradient-free inference, enabling nondifferentiable and structured optimization models to be handled seamlessly within the online loop. Empirically, on a nondifferentiable knapsack with uncertain weights, the BOCO framework attains higher, more stable reward and feasibility than Gibbs stochastic predictor and deterministic PtO and DFL baselines, particularly in the data-scarce early stages. These results show, for the first time, that coupling PAC-Bayes with PtO yields a robust, theoretically grounded, and practically effective approach to online contextual decision-making under uncertainty.

# References

Agrawal, S. and Goyal, N. (2013). Thompson sampling for contextual bandits with linear payoffs. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pages 127–135.

Alquier (2024). User-friendly introduction to pac-bayes bounds. *Foundations and Trends® in Machine Learning*, 17(2):174–303.

Amit, R., Epstein, B., Moran, S., and Meir, R. (2022). Integral probability metrics PAC-Bayes bounds. In *Advances in Neural Information Processing Systems*, volume 35, pages 3123–3136.

Bertsimas, D. and Kallus, N. (2020). From predictive to prescriptive analytics. *Management Science*, 66(3):1025–1044.

Bietti, A., Agarwal, A., and Langford, J. (2021). A contextual bandit bake-off. *Journal of Machine Learning Research*, 22(133):1–49.

Bissiri, P. G., Holmes, C. C., and Walker, S. G. (2016). A general framework for updating belief distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):1103–1130.

Capitaine, A., Haddouche, M., Moulines, E., Jordan, M. I., Boursier, E., and Durmus, A. (2025). Online decision-focused learning. *arXiv preprint*, 2505.13564.

Castiglioni, M., Celli, A., Marchesi, A., Romano, G., and Gatti, N. (2022). A unifying framework for online optimization with long-term constraints. In *Advances in Neural Information Processing Systems*, volume 35, pages 33589–33602.

Chen, H., Lu, W., and Song, R. (2021). Statistical inference for online decision making: In a contextual bandit setting. *Journal of the American Statistical Association*, 116(533):240–255.

Chen, X., Simchi-Levi, D., and Wang, Y. (2025). Utility fairness in contextual dynamic pricing with demand learning. *Management Science*, Article in advance.

Cheung, W. C., Ma, W., Simchi-Levi, D., and Wang, X. (2022). Inventory balancing with online learning. *Management Science*, 68(3):1776–1807.

Chowdhury, S. R. and Gopalan, A. (2017). On kernelized multi-armed bandits. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 844–853.

Chu, W., Li, L., Reyzin, L., and Schapire, R. (2011). Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15, pages 208–214.

Chuang, Y.-T. and Kim, M. J. (2023). Bayesian inventory control: Accelerated demand learning via exploration boosts. *Operations Research*, 71(5):1515–1529.

Clavier, P., Huix, T., and Oliviero Durmus, A. (2024). VITS : Variational inference Thompson sampling for contextual bandits. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pages 9033–9075.

Elmachtoub, A. N. and Grigas, P. (2022). Smart "predict, then optimize". *Management Science*, 68(1):9–26.

Elmachtoub, A. N., Lam, H., Lan, H., and Zhang, H. (2025). Dissecting the impact of model misspecification in data-driven optimization. In *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics*, volume 258, pages 1594–1602.

Haddouche, M. and Guedj, B. (2022). Online PAC-Bayes learning. In *Advances in Neural Information Processing Systems*, volume 35, pages 25725–25738.

Hoi, S. C., Sahoo, D., Lu, J., and Zhao, P. (2021). Online learning: A comprehensive survey. *Neurocomputing*, 459:249–289.

Ibrahim, J. G. and Chen, M.-H. (2000). Power prior distributions for regression models. *Statistical Science*, pages 46–60.

Jiang, W. and Tanner, M. A. (2008). Gibbs posterior for variable selection in high-dimensional classification and data mining. *The Annals of Statistics*, 36(5).

Kassraie, P. and Krause, A. (2022). Neural contextual bandits without regret. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151, pages 240–278.

Keyvanshokooh, E., Zhalechian, M., Shi, C., Van Oyen, M. P., and Kazemian, P. (2025). Contextual learning with online convex optimization: Theory and application to medical decision-making. *Management Science*, Article in advance.

Kong, L., Cui, J., Zhuang, Y., Feng, R., Prakash, B. A., and Zhang, C. (2022). End-to-end stochastic optimization with energy-based model. In *Advances in Neural Information Processing Systems*, volume 35, pages 11341–11354.

Krause, A. and Ong, C. S. (2011). Contextual Gaussian process bandit optimization. In *Advances in Neural Information Processing Systems*, volume 24, pages 2447–2455.

Lacasse, A., Laviolette, F., Marchand, M., Germain, P., and Usunier, N. (2006). PAC-Bayes bounds for the risk of the majority vote and the variance of the

gibbs classifier. In *Advances in Neural Information Processing Systems*, volume 19, page 769–776. MIT Press.

Li, L., Chu, W., Langford, J., and Schapire, R. E. (2010). A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th International Conference on World Wide Web*, 10, page 661–670.

Li, L., Lu, Y., and Zhou, D. (2017). Provably optimal algorithms for generalized linear contextual bandits. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 2071–2080.

Liu, H. and Grigas, P. (2022). Online contextual decision-making with a smart predict-then-optimize method. *arXiv preprint*, 2206.07316.

Liu, J. and West, M. (2001). Combined parameter and state estimation in simulation-based filtering. In *Sequential Monte Carlo methods in practice*, pages 197–223. Springer.

Lobos, A., Grigas, P., and Wen, Z. (2021). Joint online learning and decision-making via dual mirror descent. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 7080–7089.

Mandi, J., Kotary, J., Berden, S., Mulamba, M., Bucarey, V., Guns, T., and Fioretto, F. (2024). Decision-focused learning: Foundations, state of the art, benchmark and future opportunities. *Journal of Artificial Intelligence Research*, 80:1623–1701.

Mnih, A. and Gregor, K. (2014). Neural variational inference and learning in belief networks. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, pages 1791–1799.

Pacchiano, A., Ghavamzadeh, M., and Bartlett, P. (2025). Contextual bandits with stage-wise constraints. *Journal of Machine Learning Research*, 26(170):1–57.

Rodriguez-Galvez, B., Thobaben, R., and Skoglund, M. (2024). More PAC-Bayes bounds: From bounded losses, to losses with general tail behaviors, to anytime validity. *Journal of Machine Learning Research*, 25(110):1–43.

Sadana, U., Chenreddy, A., Delage, E., Forel, A., Frejinger, E., and Vidal, T. (2025). A survey of contextual optimization methods for decision-making under uncertainty. *European Journal of Operational Research*, 320(2):271–289.

Sakhi, O., Alquier, P., and Chopin, N. (2023). PAC-Bayesian offline contextual bandits with guarantees. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pages 29777–29799.

Slivkins, A., Zhou, X., Sankararaman, K. A., and Foster, D. J. (2024). Contextual bandits with packing and covering constraints: A modular lagrangian approach via regression. *Journal of Machine Learning Research*, 25(394):1–37.

Thomas, O., Dutta, R., Corander, J., Kaski, S., and Gutmann, M. U. (2022). Likelihood-free inference by ratio estimation. *Bayesian Analysis*, 17(1):1–31.

Vaswani, S., Kazemi, A., Babanezhad Harikandeh, R., and Le Roux, N. (2023). Decision-aware actor-critic with function approximation and theoretical guarantees. In *Advances in Neural Information Processing Systems*, volume 36, pages 66451–66498.

Ye, Z., Zhang, D. J., Zhang, H., Zhang, R., Chen, X., and Xu, Z. (2023). Cold start to improve market thickness on online advertising platforms: Data-driven algorithms and field experiments. *Management Science*, 69(7):3838–3860.

Zhou, D., Li, L., and Gu, Q. (2020). Neural contextual bandits with UCB-based exploration. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 11492–11502.

# PAC-Bayes Meets Online Contextual Optimization: Supplementary Materials

## 1 MISSING PROOFS

### 1.1 Proofs of Theorem 3.2 and 3.3

Theorems 3.2 and 3.3 are direct results of Corollary 3.1 and Corollary 3.3 from Haddouche and Guedj (2022) with modifications on (i) the condition in the conditional distribution of $\boldsymbol{\xi}_t$, and (ii) the constant term in the bound. We start by presenting the original Corollary 3.1 in Haddouche and Guedj (2022) and Theorem 3.2 in this work using our notation.

**Corollary A** (Haddouche and Guedj (2022), Corollary 3.1). *Suppose Assumption 3.1 holds. For any distribution $\mu$ over $D$, any $\lambda > 0$, any $\delta \in (0,1)$, and any online posterior $\{\tilde{\pi}\}$ and prior $\{\pi\}$ sequences, the following inequality holds with at least probability $1 - \delta$ over the draw $D \sim \mu$:*

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \tilde{\pi}_{t+1}}[\mathbb{E}[\ell(\theta, d_t)|\mathcal{F}_{t-1}]] \leq \sum_{t=1}^{T} \left( \mathbb{E}_{\theta \sim \tilde{\pi}_{t+1}}[\ell(\theta, d_t)] + \frac{\mathbb{D}_{KL}(\tilde{\pi}_{t+1} \| \pi_t)}{\lambda} \right) + \frac{\lambda T C^2}{2} + \frac{\log(1/\delta)}{\lambda}.$$

**Theorem 3.2.** *Suppose Assumption 3.1 holds. For any distribution $\mu$ over $D$, any $\lambda > 0$, any $\delta \in (0,1)$, and any online posterior $\{\tilde{\pi}\}$ and prior $\{\pi\}$ sequences, the following inequality holds with at least probability $1 - \delta$ over the draw $D \sim \mu$:*

$$\sum_{t=1}^{T} \mathbb{E}_{\theta \sim \tilde{\pi}_{t+1}}[\mathbb{E}[\ell(\theta, \boldsymbol{\xi}_t)|\mathcal{F}_{t-1}, x_t]] \leq \sum_{t=1}^{T} \left( \mathbb{E}_{\theta \sim \tilde{\pi}_{t+1}}[\ell(\theta, d_t)] + \frac{\mathbb{D}_{KL}(\tilde{\pi}_{t+1} \| \pi_t)}{\lambda} \right) + \frac{\lambda T C^2}{8} + \frac{\log(1/\delta)}{\lambda}.$$

The first difference lies between $\mathbb{E}[\ell(\theta, d_t)|\mathcal{F}_{t-1}]$ and $\mathbb{E}[\ell(\theta, \boldsymbol{\xi}_t)|\mathcal{F}_{t-1}, x_t]$ in the l.h.s.. The second difference lies in the denominator for the term $\lambda T C^2$.

For the first difference, since $x_t$ has already been realized, it is $\mathcal{F}_{t-1}$-measurable. Therefore, by the tower probability and measurability of $x_t$, for all $\theta$ we can write:

$$\mathbb{E}_{(X_t = x_t, \boldsymbol{\xi}_t)}[\ell(\theta, (X_t = x_t, \boldsymbol{\xi}_t))|\mathcal{F}_{t-1}] = \mathbb{E}_{X_t}[\mathbb{E}_{\boldsymbol{\xi}_t}[\ell(\theta, \boldsymbol{\xi}_t)|\mathcal{F}_{t-1}, X_t = x_t]|\mathcal{F}_{t-1}] \tag{A.1}$$

Because $X_t = x_t$ is observed, $\mathbb{P}\{X_t = x_t|\mathcal{F}_{t-1}, x_t\} = 1$. Thus, the outer expectation $\mathbb{E}_{X_t}[\cdot|\mathcal{F}_{t-1}]$ in the r.h.s. of Eq (A.1) degenerates to the form in Theorem 3.2. Note that if the target risk $\mathbb{E}_{(X_t, \boldsymbol{\xi}_t)}[\ell(\theta, (X_t, \boldsymbol{\xi}_t))|\mathcal{F}_{t-1}]$ is adopted, it suggests that we also have to predict/estimate the covariate $x_t$ for the next step, which differs from the setting of contextual optimization.

The second difference is due to the modification of the bound range in Hoeffding's lemma. According to Assumption 3.1, we know the loss $\ell(\theta, d_t) \in [0, C], \forall \theta, \forall d_t$, while Haddouche and Guedj (2022) considered the range $[-C, C]$. Therefore, $\forall \theta, \forall d_t$,

$$\Delta_t(\theta) = \mathbb{E}[\ell(\theta, \boldsymbol{\xi}_t)|\mathcal{F}_{t-1}, x_t] - \ell(\theta, d_t) \in [\mu(\theta, x_t, \mathcal{F}_{t-1}) - C, \mu(\theta, x_t, \mathcal{F}_{t-1})],$$

where $\mu(\theta, x_t, \mathcal{F}_{t-1})$ is a constant (expectation) for a given $(x_t, \theta, \mathcal{F}_{t-1})$. Therefore, Hoeffding's lemma applies as follows:

$$\mathbb{E}[e^{\lambda \Delta_t(\theta)}|\mathcal{F}_{t-1}, x_t] \leq e^{\lambda^2 C^2/8}.$$

Applying this inequality in Lemma D.2 in Haddouche and Guedj (2022) and keeping other proofs unchanged will give Theorem 3.2. Similarly, Corollary 3.3 in Haddouche and Guedj (2022) can be modified in the same way.

## 1.2 Proof of Corollary 3.5

*Proof.* By Assumption 3.4, we have $\forall \lambda > 0, \forall d \in \mathcal{X} \times \Xi, \forall \pi \in \mathcal{P}(\Xi)$:

$$\ell\left(\pi, d\right) \leq -\frac{1}{\lambda} \log \mathbb{E}_{\theta \sim \pi}\left[e^{-\lambda \ell(\theta, d)}\right] \tag{A.1}$$

In other words, this assumption ensures the pseudo-likelihood $e^{-\lambda \ell(\pi, d)}$ when implementing a decision optimized for posterior $\pi$ is always better than the expected pseudo-likelihood $e^{-\lambda \ell(\theta, d)}$ of an individual drawn $\theta \sim \pi$. We leverage this assumption to bound the risk of the aggregated predictor based on the risk of the stochastic predictor.

Applying the convexity of $e^x$ in $x$ point-wise for $d$ gives:

$$\mathbb{E}_{\theta \sim \pi}\left[e^{-\lambda \ell(\theta, d)}\right] \geq e^{-\lambda \mathbb{E}_{\theta \sim \pi}[\ell(\theta, d)]}. \tag{A.2}$$

Taking Eq (A.2) into Eq (A.1) gives:

$$\ell\left(\pi, d\right) \leq -\frac{1}{\lambda} \log\left(\mathbb{E}_{\theta \sim \pi}\left[e^{-\lambda \ell(\theta, d)}\right]\right) \leq -\frac{1}{\lambda} \log\left(e^{-\lambda \mathbb{E}_{\theta \sim \pi}[\ell(\theta, d)]}\right) = \mathbb{E}_{\theta \sim \pi}\left[\ell(\theta, d)\right]. \tag{A.3}$$

Here, Eq (A.3) is very similar to the direct application of convexity of $\ell$ in $\theta$ if one regards the $\pi$ as the mean $\mathbb{E}_{\theta \sim \pi}[m(x; \theta)]$. However, because we input the pushforward distribution into decision-making rather than the posterior mean, we thus leverage the mixability of $\ell$ to circumvent the requirement of convexity.

Taking expectation $\mathbb{E}\left[\cdot | \mathcal{F}_{t-1}, x_t\right]$ on both sides of Eq (A.3) gives:

$$\mathbb{E}\left[\ell\left(\pi, \boldsymbol{\xi}_t\right) | \mathcal{F}_{t-1}, x_t\right] \leq \mathbb{E}\left[\mathbb{E}_{\theta \sim \pi}\left[\ell\left(\theta, \boldsymbol{\xi}_t\right)\right] | \mathcal{F}_{t-1}, x_t\right]. \tag{A.4}$$

Finally, by the measurability of $\pi_t$ given $\mathcal{F}_{t-1}$ and Fubini's theorem, the expectation in the r.h.s. of Eq (A.4) can be swapped as:

$$\mathbb{E}\left[\ell\left(\pi, \boldsymbol{\xi}_t\right) | \mathcal{F}_{t-1}, x_t\right] \leq \mathbb{E}_{\theta \sim \pi}\left[\mathbb{E}\left[\ell\left(\theta, \boldsymbol{\xi}_t\right) | \mathcal{F}_{t-1}, x_t\right]\right], \tag{A.5}$$

which holds almost surely and independently of the draw of data $D \sim \mu$. Applying this upper bound in Eq (A.5) to the l.h.s. of Theorem 3.3 finishes the proof.

$$\square$$

# 2 EXPERIMENTAL DETAILS

The code of the experiments will be available soon. For all experiments, we use the open-source package SCIPY (Virtanen et al., 2020) to solve the deterministic and chance-constrained MILP problems.

## 2.1 Data Generation and Hypothesis

Our data generation follows the ARMA(2,2) process proposed in Bertsimas and Kallus (2020) to construct a time-series $(x_t, \xi_t)$. Specifically, for each stage $t$:

$$x_t = u_t + \Phi_1 x_{t-1} + \Phi_2 x_{t-2} + \Theta_1 u_{t-1} + \Theta_2 u_{t-2},$$

where $u_t \sim \mathcal{N}(\mathbf{0}, \Sigma_U)$ are innovations, and all the matrices are chosen as:

$$\Phi_1 = \begin{bmatrix} 0.5 & -0.9 & 0.0 \\ 1.1 & -0.7 & 0.0 \\ 0.0 & 0.0 & 0.5 \end{bmatrix}, \quad \Phi_2 = \begin{bmatrix} 0.0 & -0.5 & 0.0 \\ -0.5 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{bmatrix},$$

$$\Theta_1 = \begin{bmatrix} 0.4 & 0.8 & 0.0 \\ -1.1 & -0.3 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}, \quad \Theta_2 = \begin{bmatrix} 0.0 & -0.8 & 0.0 \\ -1.1 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}, \quad \Sigma_U = \begin{bmatrix} 1.0 & 0.5 & 0.0 \\ 0.5 & 1.2 & 0.5 \\ 0.0 & 0.5 & 0.8 \end{bmatrix}.$$

For each stage, once $x_t$ is generated, we generate the uncertainty $\xi_t$ as follows:

$$\tilde{\xi}_t = G(x_t + \delta_t/4) + (Bx_t) \circ \epsilon_t,$$

where $\delta_t$ and $\epsilon_t$ are independently sampled from standard Gaussian distribution, and matrices $G$ and $B$ are chosen as:

$$G = 2.5 \times \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \\ 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \\ 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \\ 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}, \quad B = 7.5 \times \begin{bmatrix} 0 & -1 & -1 \\ -1 & 0 & -1 \\ -1 & -1 & 0 \\ 0 & -1 & 1 \\ -1 & 0 & 1 \\ -1 & 1 & 0 \\ 0 & 1 & -1 \\ 1 & 0 & -1 \\ 1 & -1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}.$$

Having $\tilde{\xi}_t \in \mathbb{R}^{12}$, we do $\xi_{t,i} = \frac{\max\{-100, \tilde{\xi}_{t,i}\}}{100} + 2$ for each element in $\tilde{\xi}_t$ to ensure nonnegativity of weights. Eventually, we stack every 4 elements in the vector for a row, resulting in a matrix $A_t \in \mathbb{R}^{3 \times 4}$ as the weight matrix in optimization.

We adopt an almost-linear model with a `sigmoid` activation function for the prediction task. Denote weights $W \in \mathbb{R}^{3 \times 12}$, bias $b \in \mathbb{R}^{12}$, and `sigmoid` function: $S(x) = \frac{1}{1+e^{-x}}$. The forward prediction can be written as:

$$\hat{A}_t = 2S(Wx_t + b),$$

while model parameter $\theta$ denotes the aggregation of weights and bias. We reshape the prediction to meet the dimension of the weight matrix. In total, the dimension of $\Theta$ is $3 \times 12 + 12 = 48$.

## 2.2 Framework Details

### 2.2.1 BMA

In the implementation of `BMA` framework, we specify the following parameters: the initial prior $\pi_0 = \mathcal{N}(\mathbf{0}, I_{48})$, the shrinkage factor $a = 0.9$, the temperature $\lambda = 10^{-4}$, the Effective Sample Size threshold $\tau = 0.5$, the number of MCMC steps $L = 3$, the number of particles in SMC $N = 20$, the feasibility chance $\alpha = 0.9$. We choose relatively small $L$ and $N$ since we found they can already lead to satisfactory results while saving computational resources. Having all predictions $\hat{A}_t^i$ and the according weights $w_t^i$ for each prediction, we input these data in the following MILP problem to prescribe the decision $\hat{z}_t$ for each stage $t$ in the `BMA` framework under chance constraint:

$$\max_{u,z,V,l} \quad \sum_{i=1}^{N} w_t^i (V^i + q^\top l^i)$$

$$\text{s.t.} \quad \sum_{i=1}^{N} w_t^i u^i \geq \alpha$$

$$\hat{A}_t^i z \leq b + (1 - u^i)M, \qquad\qquad i = 1, \cdots, N$$

$$0 \leq V^i \leq u^i M, \qquad\qquad i = 1, \cdots, N$$

$$V^i \geq c^\top z - (1 - u^i)M, \qquad\qquad i = 1, \cdots, N$$

$$l^i \leq b - \hat{A}_t^i z + M(1 - u^i), \qquad\qquad i = 1, \cdots, N$$

$$l^i \geq b - \hat{A}_t^i z - M(1 - u^i), \qquad\qquad i = 1, \cdots, N$$

$$0 \leq l^i \leq u^i M, \qquad\qquad i = 1, \cdots, N,$$

$$u \in \{0,1\}^{20}, z \in \mathbb{N}^4, V \in \mathbb{R}^{20}, l \in \mathbb{R}^{20 \times 3}.$$

with $M$ here denotes a sufficiently large number to achieve the big-M modeling.

### 2.2.2 BGS

The BGS framework is a variant of BMA. Instead of leveraging the posterior approximation (by samples and weights) into the chance-constrained program for decision-making, BGS draws an indicator $I_t \sim \mathrm{Cat}(w_t^1, \ldots, w_t^N)$ to approximate the sampling from the posterior $\pi_t$, and input $\hat{A}_t^{I_t}$ into the deterministic program to prescribe a decision.

### 2.2.3 PtO

The BGS framework leverages online gradient descent (OGD) to minimize the MSE loss between the uncertainty realization $A_t$ and $\hat{A}_t$ predicted by a deterministic model. The model updated rule follows the classic OGD algorithm:

$$\theta_{t+1}^{pto} = \theta_t^{pto} - \eta_t \nabla_\theta \sqrt{\sum_i \sum_j (A_{t,i,j} - \hat{A}_{t,i,j})^2}.$$

Later, the updated model parameters will be used to predict the uncertainty $A_{t+1}$ for the next stage. In the experiments, we use standard Gaussian to initialize the parameters of PtO model.

We use Adam optimizer to update the model parameters with a step decay step size of 0.99. The initial learning rate is chosen as 0.1 in our experiments, which is optimal in terms of MSE performance among $[0.01, 0.05, 0.1, 0.5, 1, 5, 10]$ for 20 trials with length $T = 500$. Figure 1 demonstrates the time-averaged cumulative MSE loss within the PtO framework using different learning rates.
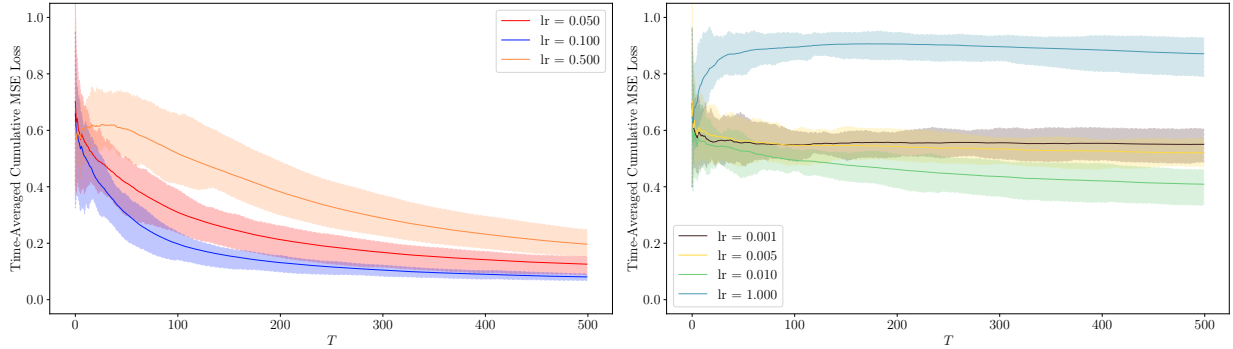


Figure 1: Time-averaged cumulative MSE loss for PtO with different learning rate. Results are computed from 20 trials with a horizon $T = 500$.

### 2.2.4 DFL

The DFL framework leverages the gradient of the decision quality of the last stage to update the deterministic predictive model parameters. Because we consider an integer knapsack problem that is nondifferentiable, we approximate the gradient for DFL using score function gradient approximation. Consider at stage $t$, the DFL model prediction is $\hat{A}_t$ and uncertainty realization is $A_t$. To approximate the gradient of $\ell$ w.r.t. $\hat{A}_t$, we draw $K$ random noise $\{\epsilon^i\}$ from $\mathcal{N}(0, I_{\dim(A_t)})$. Then, the gradient $\nabla_A \ell$ can be approximated as:

$$\nabla_A \ell \approx \frac{1}{K} \sum_{i=1}^{K} \left( c'(\mathbf{P}(\hat{A}_t + \epsilon^i), A_t) - \min_z c(z; A_t) \right) \epsilon^i.$$

In this work, we take $K = 20$ to match the computational cost of BMA in each update. The optimizer and step decay schedule are the same as PtO. 0.1 is chosen as he initial learning rate for DFL, which is optimal in terms of average reward among $[0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1]$ for 20 trials with length $T = 500$. Figure 2 demonstrates the time-averaged cumulative reward for DFL using different learning rates.
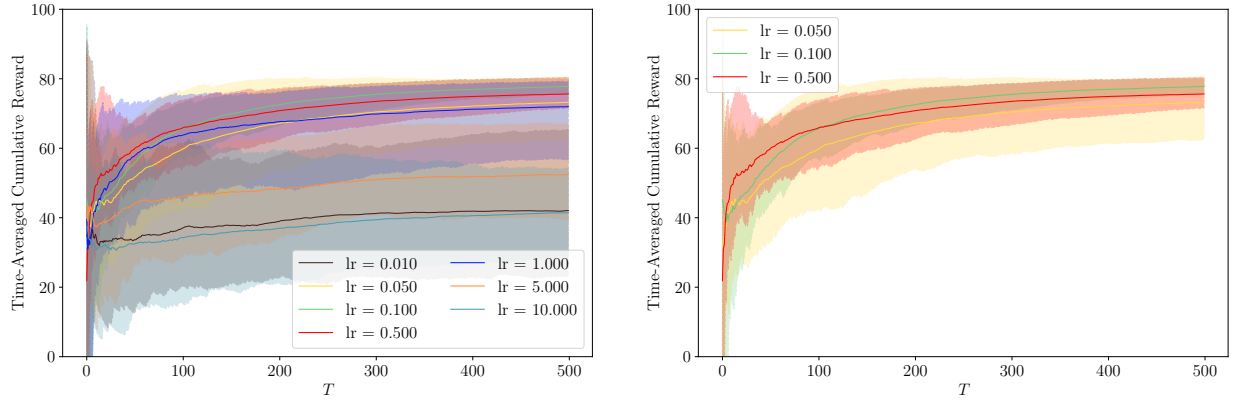
Figure 2: Time-averaged cumulative reward for `DFL` with different learning rate. Results are computed from 20 trials with a horizon $T = 500$.

## References

Bertsimas, D. and Kallus, N. (2020). From predictive to prescriptive analytics. *Management Science*, 66(3):1025–1044.

Haddouche, M. and Guedj, B. (2022). Online PAC-Bayes learning. *Advances in Neural Information Processing Systems*, 35:25725–25738.

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., et al. (2020). Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272.