

Online Optimization and Ambiguity-based Learning of Distributionally Uncertain Dynamic Systems

Dan Li¹, Dariush Fooladivanda¹ and Sonia Martínez¹

Abstract—This paper proposes a novel approach to construct data-driven online solutions to optimization problems (P) subject to a class of distributionally uncertain dynamical systems. The introduced framework allows for the simultaneous learning of distributional system uncertainty via a parameterized, control-dependent ambiguity set using a finite historical data set, and its use to make online decisions with probabilistic regret function bounds. Leveraging the merits of Machine Learning, the main technical approach relies on the theory of Distributional Robust Optimization (DRO), to hedge against uncertainty and provide less conservative results than standard Robust Optimization approaches. Starting from recent results that describe ambiguity sets via parameterized, and control-dependent empirical distributions as well as ambiguity radii, we first present a tractable reformulation of the corresponding optimization problem while maintaining the probabilistic guarantees. We then specialize these problems to the cases of 1) optimal one-stage control of distributionally uncertain nonlinear systems, and 2) resource allocation under distributional uncertainty. A novelty of this work is that it extends DRO to online optimization problems subject to a distributionally uncertain dynamical system constraint, handled via a control-dependent ambiguity set that leads to online-tractable optimization with probabilistic guarantees on regret bounds. Further, we introduce an online version of the Nesterov’s accelerated-gradient algorithm, and analyze its performance to solve this class of problems via dissipativity theory.

I. INTRODUCTION

Online optimization has attracted significant attention from various fields, including Machine Learning, Information Theory, Robotics and Smart Power Systems; see [1]–[3] and references therein. A basic online optimization setting involves the minimization of time-varying convex loss functions, resulting into Online Convex Programming (OCP). Typically, loss objectives in OCP are functions of non-stationary stochastic processes [4], [5]. Regret minimization aims to deal with non-stationarity by reducing the difference between an optimal decision made with information in hindsight, and one made as information is increasingly revealed. Thus, several online algorithms and techniques are aimed at minimizing various types of regret functions [6], [7]. More recently, and with the aim of further reducing the cost, regret-based OCP has integrated prediction models of loss functions [8]–[11]. However, exact models of evolving loss functions may not be available, while alternative data-based approximate models may require large amounts of data that are hard to obtain. This motivates the

need of developing new learning algorithms for loss functions that can employ finite data sets, while guaranteeing a precise performance of the corresponding optimization.

Literature Review. Due to recent advances in Data Science and Machine Learning, the question of learning system models as well as distributional uncertainty from data is gaining significant attention. From the early work on Systems Identification [12], Willem’s Behavioral Theory and fundamental lemma [13], [14] have been recently leveraged to learn linear, time-invariant system models in predictive control applications [14]–[18]. The aforementioned works rely on the use of Hankel system representations of the LTI system, and may be subject or not to additional uncertainty. In particular, the work [19] leverages the behavioral theory to obtain sub-linear regret bounds for the online optimization of discrete-time unknown but deterministic linear systems. Other approaches to learn LTI systems from input-output data employ concentration inequalities and finite samples, and include, for example, [20], exploiting least squares and the Ho-Kalman algorithm, [21], using subspace identification techniques for LTI systems subject to unknown Gaussian disturbances, and [22], resorting to Lasso-like methods that exploit the sparse representation of LTI systems.

On the other hand, classical online optimization relies on Sample Averaging Approximation (SAA) (with bootstrap) to derive optimal value and/or policy approximations. However, SAA usually requires large amounts of data to provide good approximations of the stochastic cost, which leads to non-robust solutions to unseen data. In contrast, recent developments on measure-of-concentration results [23] have led to a new type of Distributionally Robust Optimization (DRO) [24]–[26], which aims to bridge this gap. Particularly, the DRO framework enables finite-sample, performance-guaranteed optimization under distributional uncertainty [24], [25], and paves the way to dealing with the control and estimation of system dynamics subject to distributional uncertainty. Motivated by this, the works [27], [28] consider the time evolution of Wasserstein ambiguity sets and their updates under streaming data for estimation. However, the nominal dynamic constraints defined in these problems are assumed to be known, while in practice, these models also need to be identified. The previous work [29] proposes a method for integrating the learning of an unknown and nominal parameterized system dynamics with Wasserstein ambiguity sets. These ambiguity sets are given by a parameter and control-dependent ambiguity ball center as well as a corresponding radius. Taking this as a starting point, and motivated by the direct use of these ambiguity sets in a type of “distributionally robust control”, here we further extend this setup in connection with online

* This research was developed with funding from ONR N00014-19-1-2471, and AFOSR FA9550-19-1-0235.

¹ D. Li and S. Martínez are with the Department of Mechanical and Aerospace Engineering, University of California San Diego, La Jolla, CA 92092, USA. D. Fooladivanda is with the Department of Electrical Engineering and Computer Sciences, University of California at Berkeley, Berkeley, CA 94720, USA. dal027.me@gmail.com; dfooladi@berkeley.edu; soniamd@ucsd.edu

optimization problems. Precisely, what distinguishes this work from other approaches is the focus on learning the transition system dynamics itself via control-dependent ambiguity sets. The control method is derived from an online optimization method [6], and, therefore, it does not aim to calculate exactly an optimal control, but to find an approximate solution that leads to a low instantaneous regret function value w.r.t. standard, online and regret-based optimization problems. Finally, this manuscript connects with the topic online optimization using decision-dependent distributions [30], [31], where the uncertainty distribution changes with the decision variable. As these problems are intractable, [30], [31] solve for alternative *stable solutions*, or optimal solutions wrt to the distribution they induce. In addition to this, and while [30], [31] can handle dynamic systems, a main difference with this work is that a dynamic system structure that is being learned is not exploited, which can help reduce uncertainty more effectively.

Statement of Contributions. In this work, we propose a novel approach to solve a class of online optimization problems subject to distributionally uncertain dynamical systems. Our end goal is to produce an online controller that results in bounded instantaneous regrets with high confidence. Our proposed framework is unique in that it enables the online learning of the underlying nominal system, maintains online-problem tractability, and simultaneously provides finite-sample, probabilistic guarantee bounds on the resulting regret. This is achieved by considering a worst-case-system formulation that employs novel parameterized and control-dependent, Wasserstein ambiguity sets. Our learning method precisely consists of updating this ambiguity set. The proposed formulation is valid for a wide class of problems, including but not limited to 1) a class of optimal control problems subject to distributionally uncertain dynamical system, and 2) online resource allocation under distributional uncertainty. To do this, we first obtain tractable problem reformulations for these two cases, which results in online and non-smooth convex problem optimizations. For each of these categories, and smoothed-out versions of these problems, we propose an online control algorithm dynamics, which extends Nesterov's accelerated-gradient method. Adapting dissipativity theory, we prove optimal first-order convergence rate for these algorithms under smoothness and convexity assumptions. This result is crucial to guarantee that the online controller can provide probabilistic guarantees on their regret bounds via the control-dependent ambiguity set. We thus finish our work by quantifying these dynamic regret bounds, and by explicitly characterizing the effect of learning parameters with finite historical samples.

II. NOTATIONS

We denote by \mathbb{R}^m , $\mathbb{R}_{\geq 0}^m$, $\mathbb{Z}_{\geq 0}^m$ and $\mathbb{R}^{m \times n}$ the m -dimensional real space, nonnegative orthant, nonnegative integer-orthant space, and the space of $m \times n$ matrices, respectively. The transpose of a column vector $\mathbf{x} \in \mathbb{R}^m$ is \mathbf{x}^\top , and $\mathbf{1}_m$ is a shorthand for $(1, \dots, 1)^\top \in \mathbb{R}^m$. We index vectors with subscripts, i.e., $\mathbf{x}_k \in \mathbb{R}^m$ with $k \in \mathbb{Z}_{\geq 0}$, and given $\mathbf{x} \in \mathbb{R}^m$ we denote its i^{th} component by x_i . We denote by $\|\mathbf{x}\|$ and $\|\mathbf{x}\|_\infty$ the 2-norm and ∞ -norm, respectively. The inner product of \mathbb{R}^m is given as $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^\top \mathbf{y}$, $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$; thus, $\|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$. The gradient of a real-valued function $\ell : \mathbb{R}^m \rightarrow \mathbb{R}$ is denoted as

$\nabla \ell(\mathbf{x})$ and $\nabla_x \ell(\mathbf{x})$ is the partial derivative w.r.t. x . In what follows, $\text{dom}(\ell) := \{\mathbf{x} \in \mathbb{R}^m \mid -\infty < \ell(\mathbf{x}) < +\infty\}$. A function $\ell : \text{dom}(\ell) \rightarrow \mathbb{R}$ is M -strongly convex, if for any $\mathbf{y}, \mathbf{z} \in \text{dom}(\ell)$ there exists $\mathbf{g} \in \mathbb{R}^m$ such that $\ell(\mathbf{y}) \geq \ell(\mathbf{z}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{z}) + M \|\mathbf{y} - \mathbf{z}\|^2/2$, for some $M > 0$. The function ℓ is convex if $M \geq 0$. We call a vector \mathbf{g} a subgradient of ℓ at \mathbf{z} and denote by $\partial \ell(\mathbf{z})$ the subgradient set. If ℓ is differentiable at \mathbf{z} , then $\partial \ell(\mathbf{z}) = \{\nabla \ell(\mathbf{z})\}$. Finally, the operation $\Pi_{\mathcal{U}}(\mathcal{X}) : \mathcal{X} \rightarrow \mathcal{U}$ projects the set $\mathcal{X} \subseteq \mathbb{R}^m$ onto $\mathcal{U} \subseteq \mathbb{R}^m$ under the Euclidean norm. We write $\Pi_{\mathcal{U}}(\mathbf{x}) := \text{argmin}_{\mathbf{z}} \|\mathbf{x} - \mathbf{z}\|^2/2 + \chi_{\mathcal{U}}(\mathbf{z})$, where $\mathbf{x} \in \mathcal{X}$, and $\chi_{\mathcal{U}}(\mathbf{z}) = 0$ if $\mathbf{z} \in \mathcal{U}$, otherwise $+\infty$. Endow \mathbb{R}^n with the Borel σ -algebra \mathcal{B} , and let $\mathcal{P}(\mathbb{R}^n)$ be the set of probability measures (or distributions) over $(\mathbb{R}^n, \mathcal{B})$. The set of probability distributions with bounded first moments is $\mathcal{M} = \{\mathbb{Q} \in \mathcal{P}(\mathbb{R}^n) \mid \int_{\mathbb{R}^n} \|\mathbf{x}\| d\mathbb{Q} < +\infty\}$. We use the Wasserstein metric [32] to define a distance in \mathcal{M} , and the dual version of the 1-Wasserstein metric $d_W : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}_{\geq 0}$, is defined by $d_W(\mathbb{Q}_1, \mathbb{Q}_2) := \sup_{f \in \mathcal{L}} \int f(\mathbf{x}) d\mathbb{Q}_1 - \int f(\mathbf{x}) d\mathbb{Q}_2$, where \mathcal{L} is the space of all Lipschitz functions with Lipschitz constant 1. We denote a closed Wasserstein ball of radius ϵ (also called an ambiguity set) centered at a distribution $\mathbb{P} \in \mathcal{M}$ by $\mathbb{B}_\epsilon(\mathbb{P}) := \{\mathbb{Q} \in \mathcal{M} \mid d_W(\mathbb{P}, \mathbb{Q}) \leq \epsilon\}$. The Dirac measure at $\mathbf{x}_0 \in \mathbb{R}^n$ is a distribution in $\mathcal{P}(\mathbb{R}^n)$ denoted by $\delta_{\{\mathbf{x}_0\}}$. Given $A \in \mathcal{B}$, we have $\delta_{\{\mathbf{x}_0\}}(A) = 1$, if $\mathbf{x}_0 \in A$, otherwise 0. A random vector $\mathbf{x} \in \mathbb{R}^m$ with probability distribution \mathbb{Q} is sub-Gaussian if there are positive constants C, v such that $\mathbb{Q}(\|\mathbf{x}\| > t) \leq Ce^{-vt^2}$. Equivalently, a zero-mean random vector $\mathbf{x} \in \mathbb{R}^n$ is sub-Gaussian if for any $a \in \mathbb{R}^n$ we have $\mathbb{E}[\exp(a^\top \mathbf{x})] \leq \exp(\|a\|^2 \nu^2/2)$ for some ν .

III. PROBLEM STATEMENT, MOTIVATION, AND APPROACH BASED ON AMBIGUITY SET LEARNING

We start by introducing a class of online optimization problems, where the objective function is time-varying according to an unknown dynamical system subject to an unknown disturbance. Consider a dynamical system that evolves according to unknown stochastic dynamics

$$\mathbf{x}_{t+1} = f(t, \mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}_t, \quad \text{from a given } \mathbf{x}_0 \in \mathbb{R}^n, \quad (1)$$

where $\mathbf{u}_t \in \mathcal{U} \subset \mathbb{R}^m$ is an online decision or control action at time t , $f : \mathbb{R}_{\geq 0} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ is a measurable, but unknown state transition function, and $\mathbf{w}_t \in \mathbb{R}^n$, is an unknown, random, disturbance vector. Due to the Markov assumption, $\mathbf{x}_t \in \mathbb{R}^n$ can be described by an unknown transition probability measure $\mathbb{P}_{t|t-1} \in \mathcal{P}(\mathbb{R}^n)$, conditioned on the system state and control at time $t-1$. Denote by $\ell : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$, $(\mathbf{u}, \mathbf{x}) \mapsto \ell(\mathbf{u}, \mathbf{x})$ an *a-priori* selected, measurable loss function. Assume that \mathcal{U} is compact, and we are interested in selecting $\mathbf{u}_t \in \mathcal{U}$ that minimizes the loss

$$\min_{\mathbf{u}_t \in \mathcal{U}} \left\{ \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}_t, \mathbf{x})] := \int_{\mathbb{R}^n} \ell(\mathbf{u}_t, \mathbf{x}) \mathbb{P}_{t+1|t}(d\mathbf{x}) \right\}.$$

This objective value is inaccessible since the state distribution $\mathbb{P}_{t+1|t}$ is unknown, and its evolution is highly dependent on the system, disturbance, and as well as on the decisions taken. In this work, we aim to propose an effective online optimization and learning algorithm which tracks the minimizers of the time-varying objective function with low regret in high

probability. Thus, at each time t , we aim to find $\mathbf{u} := \mathbf{u}_t$ that minimizes the loss in the immediate future at $t + 1$

$$\begin{aligned} \min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}, \mathbf{x})], \\ \text{s. t. } \mathbf{x} \sim \mathbb{P}_{t+1|t}, \text{ evolves according to (1).} \end{aligned} \quad (\text{P})$$

This problem formulation is similar to a one-stage optimization problems with unknown system transitions [33]. The expectation operator with respect to $\mathbb{P}_{t+1|t}$ is conditional on the historical realizations $\hat{\mathbf{x}}_k, k \leq t$, the adopted decisions $\hat{\mathbf{u}}_k, k \leq t - 1$, the yet-to-be-learned unknown dynamical system f , and realizations $\hat{\mathbf{w}}_k, k \leq t - 1$. We will identify $\mathbb{P}_{t+1|t}(d\mathbf{x}) \equiv \mathbb{P}_{t+1}(d\mathbf{x}|\mathbf{u}_t, \mathbf{x}_t = \hat{\mathbf{x}}_t, \mathbf{x}_k = \hat{\mathbf{x}}_k, \mathbf{u}_k = \hat{\mathbf{u}}_k, k \leq t - 1)$ which, by the Markovian property, satisfies $\mathbb{P}_{t+1|t}(d\mathbf{x}) \equiv \mathbb{P}_{t+1}(d\mathbf{x}|\mathbf{u}_t, \mathbf{x}_t = \hat{\mathbf{x}}_t)$. At time t , let $\mathbf{u}^* := \mathbf{u}_t^*$ denote an optimizer of Problem (P) and consider the instantaneous regret

$$R_t := \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}, \mathbf{x})] - \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}^*, \mathbf{x})],$$

which is the loss incurred if the selected \mathbf{u} is different from an optimal decision. Our goal will be to develop a robust online algorithm which ensures a probabilistic bound on the regret. That is, with high probability ρ , the regret R_t is upper bounded by a sum of terms, a first one depending on the initial condition \mathbf{x}_0 ; a second one depending on the instantaneous variation of the loss of (P); and a third term related to how well the unknown system f and the uncertainty are characterized; please see Theorem V.1. While the second and third terms are inherent to the system, the effect of the second one can be reduced by considering a predicted loss of the system [11]. In this work, we aim to bound the third term and minimize it by estimating the distribution $\mathbb{P}_{t+1|t}$ via an ambiguity set of distributions. We will show that, as historical data are assimilated over time, this third term asymptotically decays to zero. This is achieved under the following assumption

Assumption III.1 (Independent and stationary sub-Gaussian distributions) The vectors $\mathbf{w}_t \in \mathbb{R}^n, t \in \mathbb{Z}_{\geq 0}$, are i.i.d. with $\mathbf{w}_t \sim \mathbb{Q}$ and zero-mean σ sub-Gaussian¹.

Remark III.1 (On sub-Gaussian distributions) Sub-Gaussian distributions include Gaussian random variables and all distributions of bounded support.

Example 1 (Vehicle path planning and tracking): A two-wheeled vehicle moves in an unknown 2D environment. Assume that an accessible path-planner provides a control signal for the vehicle to track a desired reference trajectory under ideal conditions, see Fig. 1. Fig. 2 shows two examples where, first, the vehicle implements a series of lane changes, and, second, navigates through a planned circular/loopy route. Since both the environment and dynamics are uncertain, exact tracking is rare. Our goal is to learn the real-time road conditions, and by solving the online problem (P), derive a control signal that enables path following minimizing the tracking error with high probability.

Example 2 (Online resource allocation in the stock market): An agent aims to achieve a target profit of $r_0 = 130\%$ in a highly-fluctuating trading market. Thus, it actively allocates

¹That is, for all unit vector \mathbf{v} , we have $\mathbb{E}[e^{\lambda \mathbf{v}^\top \mathbf{w}_t}] \leq e^{\lambda^2 \sigma^2 / 2}, \forall \lambda \in \mathbb{R}$. Equivalently, $\mathbb{Q}(\|\mathbf{w}_t\| > \lambda) \leq e^{-\lambda^2 / (4\sigma^2)}, \forall \lambda$.

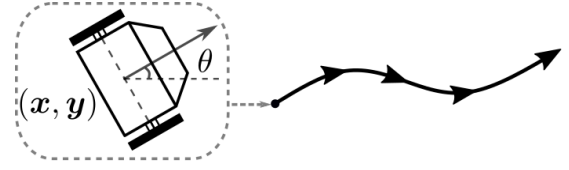


Fig. 1: A two-wheeled vehicle model with $(x, y) \in \mathbb{R}^2$ the position of the center and θ the direction.

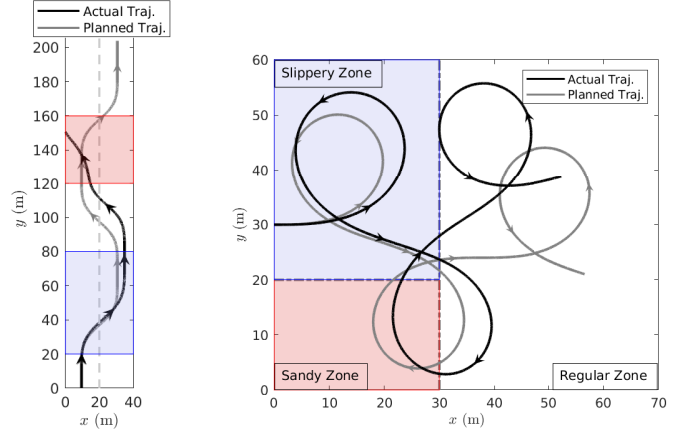


Fig. 2: The (gray) planned trajectory and (black) actual system trajectory in various road zones, with the system state $\mathbf{x} = (x, y, \theta) \in \mathbb{R}^2 \times [-\pi, \pi)$. The red region indicates sandy zone while the blue region indicates the slippery zone. Due to unknown road conditions, the actual system trajectories deviate from planned trajectories.

wealth to multiple risky assets while trying to balance resources among assets. As asset-prices are uncertain, modeling the return rate of each asset is specially challenging. To solve this, an agent can aim to learn the real-time returns responsively, estimate the distributions of immediate returns, and then allocate wealth wisely to maximize the expected profit with high probability. This problem fits in the proposed formulation, resulting in online, balanced resource allocation with low regrets.

A. Online Constructions of Ambiguity Sets

Our main approach to obtain a suitable control signal is based on learning a set of distributions or ambiguity set that characterizes system uncertainty. More precisely, we employ the dynamic ambiguity set \mathcal{P}_{t+1} proposed in [29]. The set \mathcal{P}_{t+1} contains a class of distributions, which is, in high probability, large enough to include the unknown $\mathbb{P}_{t+1|t}$ under certain conditions. Thus, we can use it to formulate a robust version of the problem at each time instant t . Such characterization enables an online-tractable reformulation of (P) later. We summarize next the construction of these ambiguity sets \mathcal{P}_{t+1} . First, we assume the following on the unknown f .

Assumption III.2 (System parametrization) Given $p \in \mathbb{Z}_{>0}$, the system f can be expressed as

$$f(t, \mathbf{x}, \mathbf{u}) = \sum_{i=1}^p \alpha_i^* f^{(i)}(t, \mathbf{x}, \mathbf{u}),$$

where $\boldsymbol{\alpha}^* := (\alpha_1^*, \dots, \alpha_p^*)^\top \in \mathbb{R}^p$ is an unknown parameter, and $f^{(i)} : \mathbb{R}_{\geq 0} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n, (t, \mathbf{x}, \mathbf{u}) \mapsto f^{(i)}(t, \mathbf{x}, \mathbf{u})$,

$i \in \{1, \dots, p\}$ is a set of p linearly independent known basis functions or *predictors* chosen *a priori*.

Now, given arbitrary (α, \mathbf{u}) , the set \mathcal{P}_{t+1} is a Wasserstein ball centered at a parametric-dependent distribution $\hat{\mathbb{P}}_{t+1|t}$ for each t ; that is,

$$\mathcal{P}_{t+1} := \mathbb{B}_{\hat{\epsilon}}(\hat{\mathbb{P}}_{t+1|t}) = \{\mathbb{Q} \mid d_W(\mathbb{Q}, \hat{\mathbb{P}}_{t+1|t}) \leq \hat{\epsilon}\}.$$

Here, $\hat{\epsilon}$ will be a time-varying function $\hat{\epsilon} \equiv \hat{\epsilon}(t, T, \beta, \alpha, \mathbf{u})$ which depends on a number of T measurements, and a confidence $\beta \in (0, 1)$. More precisely,

$$\hat{\mathbb{P}}_{t+1|t}(\alpha, \mathbf{u}) := \frac{1}{T} \sum_{k \in \mathcal{T}} \delta_{\{\sum_{i=1}^p \alpha_i \xi_k^{(i)}(\alpha, \mathbf{u})\}}; \quad (2)$$

see the footnote², where $\mathcal{T} = \{t-T, \dots, t\}$, for $t \geq T+1$. If $\alpha = \alpha^*$, then $\sum_{i=1}^p \alpha_i \xi_k^{(i)}(\alpha, \mathbf{u}) \in \mathbb{R}^n$ provides an outcome $\mathbf{x}_{t+1}^{(k)} := f(t, \mathbf{x}_t, \mathbf{u}) + \mathbf{w}_k = \sum_{i=1}^p \alpha_i^* f^{(i)}(t, \mathbf{x}_t, \mathbf{u}) + \mathbf{w}_k$, for each k . For a general $\alpha \approx \alpha^*$, the value $\sum_{i=1}^p \alpha_i \xi_k^{(i)}(\alpha, \mathbf{u})$ provides ‘‘approximated’’ outcomes $\mathbf{x}_{t+1}^{(k)}$, for each $k = 1, \dots, T$. Then, we claim the probabilistic guarantee of \mathcal{P}_{t+1} by a selection of the parameter α and $\hat{\epsilon}$ for any \mathbf{u} .

Theorem III.1 (Online probabilistic guarantee [29, Application of Theorem 1]) *Let Assumptions III.1 and III.2 hold. For a given $T \in \mathbb{Z}_{>0}$, historical data $\{\hat{\mathbf{x}}_k\}_{k \in \mathcal{T}}$ and $\{\mathbf{u}_k\}_{k \in \mathcal{T} \setminus \{t\}}$, $\mathcal{T} = \{t-T, \dots, t\}$, we select $\hat{\mathbb{P}}_{t+1|t}$ as in (2) where α is selected in [29, Theorem 2 (Learning of α^*)]³. Then, for given \mathbf{u} and a confidence-related value $\beta \in (0, 1)$, a radius $\hat{\epsilon} := \hat{\epsilon}(t, T, \beta, \alpha, \mathbf{u})$ can be chosen such that*

$$\text{Prob}(\mathbb{P}_{t+1|t} \in \mathcal{P}_{t+1}) \geq \rho(t). \quad (3)$$

Here, the left-hand-side expression is a shorthand for the probability of the event $\{(\mathbf{x}_{t+1}^{(1)}, \dots, \mathbf{x}_{t+1}^{(T)}) \in \mathbb{R}^n \times \dots \times \mathbb{R}^n \mid \mathbb{P}_{t+1|t} \in \mathbb{B}_{\hat{\epsilon}}(\hat{\mathbb{P}}_{t+1|t})\}$ and $\text{Prob} := \mathbb{P}_{t+1|t}^T$ denotes the probability measure defined on the T -fold product of $\mathbb{P}_{t+1|t}$, which evaluates the probability that the selection of samples define an ambiguity ball which contains the true distribution. In particular, the confidence value is

$$\rho(t) := (1 - \beta) \left(1 - \exp \left(- \frac{(\gamma^2 - \sqrt{2}c\gamma)T}{2\sqrt{2}(c\gamma + \sqrt{2}c^2)} \right) \right),$$

where c is a data-dependent positive constant and $\gamma > \sqrt{2}c$ is a user selected parameter. Further, the radius is

$$\hat{\epsilon} := \sqrt{\frac{2nM\sigma^2}{T} \ln\left(\frac{1}{\beta}\right)} + C_1 T^{-1/\max\{n, 2\}} + \gamma H(t, T, \mathbf{u}), \quad (4)$$

where M and C_1 are positive constants, and

$$H(t, T, \mathbf{u}) := \frac{1}{T} \sum_{i=1}^p \sum_{k \in \mathcal{T}} \|f^{(i)}(k, \hat{\mathbf{x}}_k, \mathbf{u}_k) - f^{(i)}(t, \hat{\mathbf{x}}_t, \mathbf{u})\|,$$

which bounds the variation of predicted system trajectories.

Idea of the Proof. The probabilistic guarantees (3) are a consequence of Lemma 1, Theorem 1, Theorem 2

² $\xi_k^{(i)}(\alpha, \mathbf{u}) := f^{(i)}(t, \hat{\mathbf{x}}_t, \mathbf{u}) + \hat{\mathbf{x}}_{k+1}/(\alpha^\top \mathbf{1}_p) - f^{(i)}(k, \hat{\mathbf{x}}_k, \mathbf{u}_k)$, with $\hat{\mathbf{x}}_t, \hat{\mathbf{x}}_{k+1}, \hat{\mathbf{x}}_k$ being the state measurements at time $t, k+1, k$ and \mathbf{u}_k being the past input at $k, k \in \mathcal{T}$.

³ In [29, Theorem 2], the value \mathbf{d} plays the role of \mathbf{u} in this work.

and Eqn. (7) in [29] with Assumptions III.1 and III.2. Precisely, we achieve this by upper bounding the metric $d_W(\mathbb{P}_{t+1|t}, \hat{\mathbb{P}}_{t+1|t}(\alpha, \mathbf{u}))$ using $d_W(\mathbb{P}_{t+1|t}, \hat{\mathbb{P}}_{t+1|t}(\alpha^*, \mathbf{u}))$ plus $d_W(\hat{\mathbb{P}}_{t+1|t}(\alpha^*, \mathbf{u}), \hat{\mathbb{P}}_{t+1|t}(\alpha, \mathbf{u}))$. Then, the first distance is handled via [29, Lemma 1] using standard measure of concentration results⁴, contributing to the first two terms of the radius $\hat{\epsilon}$ in (4). Next, the second distance $d_W(\hat{\mathbb{P}}_{t+1|t}(\alpha^*, \mathbf{u}), \hat{\mathbb{P}}_{t+1|t}(\alpha, \mathbf{u}))$ can be bounded in terms of the difference $\|\alpha - \alpha^*\|$ via [29, Theorem 1], contributing to the third term in $\hat{\epsilon}$. Notice that the third term depends on Assumption III.2 and the selected parameter γ which relies on the selection of α via [29, Theorem 2 (Learning of α^*)]. The confidence value $\rho(t)$ is achieved by Assumption III.1 applying to the same procedure as in [29, Theorem 2], which essentially bounds $\|\alpha - \alpha^*\|_\infty$ in probability. Precisely, by Assumption III.1, we have $\mathbb{Q}(\|\mathbf{w}_t\|_\infty > \eta) \leq e^{-\eta^2/(4\sigma^2)}$, $\forall \eta$, resulting in $\mathbb{E}[\|\mathbf{w}_t\|_\infty^l] \leq 2^{\frac{l}{2}-1} \sigma^l l^{\frac{l}{2}+1}$, $\forall l \in \mathbb{Z}_{>0}$, analogous to [34, Lemma 2]. Then, with the proof similar to [34, Theorem IV.2], we achieve

$$\text{Prob} \left(\frac{1}{T} \sum_{k \in \mathcal{T}} (\|\mathbf{w}_k\|_\infty) \geq \gamma \right) \leq \exp \left(-\gamma\lambda + \frac{T\sqrt{2}\sigma e\lambda}{2T - 2\sqrt{2}\sigma e\lambda} \right).$$

By selecting

$$\lambda = \begin{cases} \frac{T}{2\sqrt{2}\sigma e} - \frac{T}{2\gamma}, & \text{if } \gamma \geq \sqrt{2}\sigma e, \\ 0, & \text{if } \gamma < \sqrt{2}\sigma e, \end{cases}$$

we follow the proof [34, Theorem IV.2] to achieve

$$\begin{aligned} \text{Prob} \left(\frac{1}{T} \sum_{k \in \mathcal{T}} (\|\mathbf{w}_k\|_\infty) \geq \gamma \right) & \\ & \leq \begin{cases} \exp \left(- \frac{T(\gamma^2 - \sqrt{2}c\gamma)}{2\sqrt{2}\sigma e(\gamma + \sqrt{2}c)} \right), & \text{if } \gamma \geq \sqrt{2}\sigma e, \\ 1, & \text{if } \gamma < \sqrt{2}\sigma e. \end{cases} \end{aligned}$$

By bound propagation, we have

$$\text{Prob}(\|\alpha - \alpha^*\|_\infty \leq \gamma) \geq 1 - \exp \left(- \frac{(\gamma^2 - \sqrt{2}c\gamma)T}{2\sqrt{2}(c\gamma + \sqrt{2}c^2)} \right),$$

with $\gamma > \sqrt{2}c$ and c is selected as in [29, Theorem 2]. Finally, the combination of all the above considerations complete the proof. ■

Theorem III.1 provides a methodology to construct online ambiguity sets with guarantees in probability. In general, $\rho(t)$ is strictly smaller than 1 unless there is a way of making $\alpha(t) \rightarrow \alpha^*$. This is implemented in [29] via an online learning algorithm which leads to $\rho(t) \rightarrow 1 - \beta$ via Eqn. (7) in the same work. Notice how these constructions are related to the decision variable \mathbf{u} and, in the following, we leverage the probabilistic characterization $\mathcal{P}_{t+1} := \mathcal{P}_{t+1}(\alpha, \mathbf{u})$ of the distribution $\mathbb{P}_{t+1|t}$ for solutions to (P).

IV. A TRACTABLE PROBLEM REFORMULATION AND ITS SPECIALIZATION TO TWO PROBLEM CLASSES

In this section, we start by describing how to deal with the unknown $\mathbb{P}_{t+1|t}$ in Problem (P), via ambiguity sets, which

⁴ Lemma 1 in [29] makes use of a stronger Assumption III.1, which requires \mathbf{w}_k to be white. However, this can be relaxed to the current assumption by multiplying the upper bound in the lemma with a constant $M > 0$ associated with noise whitening via an appropriate linear transformation.

results in (P1). By doing this, the solution of (P1) provides guarantees on the performance of (P). Unfortunately, this results into an online intractable problem. Thus, we find a tractable reformulation (P2) which is equivalent to (P1) under certain conditions. After this, we focus the rest of our work on two problem sub-classes, which allows us to present and analyze the online algorithms for these problems in the following section. Formally, let us consider

$$\min_{\mathbf{u} \in \mathcal{U}} \sup_{\mathbb{Q} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{u}, \mathbf{x})], \quad (\text{P1})$$

where, for a fixed $\boldsymbol{\alpha} := \boldsymbol{\alpha}_t$ and $\mathbf{u} := \mathbf{u}_t \in \mathcal{U}$, it holds that $\mathbb{P}_{t+1|t} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})$ with high probability. This results in

$$\text{Prob} \left(\mathbb{E}_{\mathbb{P}_{t+1|t}}[\ell(\mathbf{u}, \mathbf{x})] \leq \sup_{\mathbb{Q} \in \mathcal{P}_{t+1}} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{u}, \mathbf{x})] \right) \geq \rho(t).$$

Observe that, the probability measure Prob and the bound $\rho(t)$ coincides with that in (3) and notice how the value $\rho(t)$ changes for various data-set sizes T in Theorem III.1.

The solution \mathbf{u} and the objective value of (P1) ensure that, when we select \mathbf{u} to be the decision for (P), the expected loss of (P) is no worse than that from (P1) with high probability. The formulation (P1) still requires expensive online computations due to its semi-infinite inner optimization problem. Thus, we propose an equivalent reformulation of (P1) for a class of loss functions as in the following assumption.

Assumption IV.1 (Lipschitz loss functions) Consider the loss function $\ell : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$, $(\mathbf{u}, \mathbf{x}) \mapsto \ell(\mathbf{u}, \mathbf{x})$. There exists a Lipschitz function $L : \mathbb{R}^m \rightarrow \mathbb{R}_{\geq 0}$ such that for each $\mathbf{u} \in \mathbb{R}^m$, it holds that $\|\ell(\mathbf{u}, \mathbf{x}) - \ell(\mathbf{u}, \mathbf{y})\| \leq L(\mathbf{u})\|\mathbf{x} - \mathbf{y}\|$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.

With this, we obtain the following upper bound:

Lemma IV.1 (An upper bound of (P1)) *Let Assumption IV.1 hold. Then, for each \mathbf{u} , $\boldsymbol{\alpha}$, β , T and t , we have*

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{u}, \mathbf{x})] \\ & \leq \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] + \hat{\epsilon}(t, T, \beta, \boldsymbol{\alpha}, \mathbf{u})L(\mathbf{u}), \end{aligned}$$

where the empirical distribution $\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})$ and scalar $\hat{\epsilon}(t, T, \beta, \boldsymbol{\alpha}, \mathbf{u})$ are described as in Section III-A.

Hereafter, see the appendix for all proofs.

Next, we claim that the upper bound in Lemma IV.1 is tight if the following assumption holds.

Assumption IV.2 (Convex and gradient-accessible functions) The loss function ℓ is convex in \mathbf{x} for each \mathbf{u} . Further, for each time t with given $\mathbf{u}(= \mathbf{u}_t) \in \mathcal{U}$ and $\boldsymbol{\alpha}(= \boldsymbol{\alpha}_t) \in \mathbb{R}^p$, there is a system prediction $\sum_{i=1}^p \alpha_i \xi_k^{(i)}(\boldsymbol{\alpha}, \mathbf{u})$ for some $k \in \mathcal{T}$ such that $\nabla_{\mathbf{x}} \ell$ exists and $L(\mathbf{u})$ is equal to $\|\nabla_{\mathbf{x}} \ell\|$ at $(\mathbf{u}, \sum_{i=1}^p \alpha_i \xi_k^{(i)}(\boldsymbol{\alpha}, \mathbf{u}))$.

The above statement enables the following theorem.

Theorem IV.1 (Equivalent reformulation of (P1)) *Let Assumptions IV.1 and IV.2 hold. Let Ξ_{t+1} denote the support of the distribution $\mathbb{P}_{t+1|t}$. Then, if $\Xi_{t+1} = \mathbb{R}^n$, (P1) is equivalent to the following problem*

$$\min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] + \hat{\epsilon}(t, T, \beta, \boldsymbol{\alpha}, \mathbf{u})L(\mathbf{u}). \quad (\text{P2})$$

Remark IV.1 (Effects of Assumptions IV.1 and IV.2) We note that Assumption IV.1 on the Lipschitz requirement of loss function is mild. In fact, many engineering problems take state values in a compact set, which then only requires the loss ℓ to be continuous. Assumption IV.2 essentially requires accessible partial gradients (in \mathbf{x}) of loss functions ℓ . For simple loss functions ℓ , e.g. linear, quadratic, etc, its partial gradient can be readily evaluated. Notice that when Assumption IV.2 fails, Problem (P2) still serves as a relaxation problem of (P1), providing a solution with a valid upper bound.

Notice that the tractability of solutions to (P2) now depend on: 1) the choice of the loss function ℓ and the associated Lipschitz function L , and 2) the decision space \mathcal{U} . To be able to further analyze (P2) and further evaluate Assumption IV.2 on gradient-accessible functions, we will impose further structure on the system as follows:

Assumption IV.3 (Locally Lipschitz, control-affine system and basis functions) *The system f is locally Lipschitz in $(t, \mathbf{x}, \mathbf{u})$ and affine in \mathbf{u} , i.e.,*

$$f(t, \mathbf{x}, \mathbf{u}) := f_1(t, \mathbf{x}) + f_2(t, \mathbf{x})\mathbf{u},$$

for some unknown $f_1 : \mathbb{R}_{\geq 0} \times \mathbb{R}^n \rightarrow \mathbb{R}^n$, $f_2 : \mathbb{R}_{\geq 0} \times \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$, $\mathbf{u} \in \mathcal{U}$ and $t \in \mathbb{Z}_{\geq 0}$. Similarly, for each $i \in \{1, \dots, p\}$, the basis function $f^{(i)}$ is selected to be

$$f^{(i)}(t, \mathbf{x}, \mathbf{u}) := f_1^{(i)}(t, \mathbf{x}) + f_2^{(i)}(t, \mathbf{x})\mathbf{u},$$

for some known locally Lipschitz functions $f_1^{(i)}$ and $f_2^{(i)}$.

Assumption IV.4 (Convex decision oracle) *The set \mathcal{U} is convex and compact. Furthermore, the projection operation of $\mathbf{u} \in \mathbb{R}^m$ onto \mathcal{U} , $\Pi_{\mathcal{U}}(\mathbf{u})$, admits $O(1)$ computation complexity.*

For simplicity of the discussion, we rewrite (P2) as

$$\min_{\mathbf{u} \in \mathcal{U}} G(t, \mathbf{u}) := G(t, \mathbf{u} | \ell, L, T, \beta, \boldsymbol{\alpha}, \hat{\mathbb{P}}_{t+1|t}, \hat{\epsilon}),$$

where G represents the objective function of (P2), depending on variables ℓ , L , β , $\boldsymbol{\alpha}$, and \mathcal{P}_{t+1} , which are kept fixed in the optimization. Then, Assumption IV.3 allows an explicit expression of G w.r.t. $\mathbf{u} := \mathbf{u}_t$ and Assumption IV.4 characterizes the convex feasible set of (P2). Note that $G(t, \mathbf{u})$ is locally Lipschitz in t .⁵

In the following, we consider two classes of general problems in the form of (P2): 1) an optimal control problem under the uncertainty; 2) an online resource allocation problem with a switch. These problems leverage the probabilistic characterization of the system and common loss functions ℓ . Then, we propose an online algorithm to achieve tractable solutions with a probabilistic regret bound in the next section.

Problem 1: (Optimal control under uncertainty) We consider a problem in form (P), where the system is unknown and is to be optimally controlled. In particular, we employ the following separable loss function

$$\ell(\mathbf{u}, \mathbf{x}) := \ell_1(\mathbf{u}) + \ell_2(\mathbf{x}), \quad \ell_1 : \mathbb{R}^m \rightarrow \mathbb{R}, \ell_2 : \mathbb{R}^n \rightarrow \mathbb{R},$$

with ℓ_1 the cost for the immediate control and ℓ_2 the optimal cost-to-go function. We assume that both ℓ_1 and ℓ_2 are convex,

⁵This can be verified by the local Lipschitz condition on $f^{(i)}$, ℓ , and finite composition of local Lipschitz functions are locally Lipschitz.

and in addition, ℓ_2 is Lipschitz continuous with a constant $\text{Lip}(\ell_2) \in \mathbb{R}_{\geq 0}$, resulting in $L(\mathbf{u}) \equiv \text{Lip}(\ell_2)$. Then, by selecting the ambiguity radius $\hat{\epsilon}$ and center $\mathbb{P}_{t+1|t}$ of \mathcal{P}_{t+1} as in Section III-A, the objective function of (P2) becomes

$$G(t, \mathbf{u}) = \ell_1(\mathbf{u}) + \frac{1}{T} \sum_{k \in \mathcal{T}} \ell_2(\mathbf{p}_{k,t}) \\ + \text{Lip}(\ell_2)\epsilon + \frac{\gamma \text{Lip}(\ell_2)}{T} \sum_{i=1}^p \sum_{k \in \mathcal{T}} \|\mathbf{H}_k^{(i)}\|,$$

where $\mathbf{p}_{k,t}, \mathbf{H}_k^{(i)} \in \mathbb{R}^n$ are affine in \mathbf{u} , for each i, k , as

$$\mathbf{p}_{k,t} := \sum_{i=1}^p \alpha_i \left(f_1^{(i)}(t, \hat{\mathbf{x}}_t) - f_1^{(i)}(k, \hat{\mathbf{x}}_k, \mathbf{u}_k) \right) \\ + \hat{\mathbf{x}}_{k+1} + \left(\sum_{i=1}^p \alpha_i f_2^{(i)}(t, \hat{\mathbf{x}}_t) \right) \mathbf{u},$$

$$\mathbf{H}_k^{(i)}(\mathbf{u}) := f_1^{(i)}(k, \hat{\mathbf{x}}_k, \mathbf{u}_k) - f_1^{(i)}(t, \hat{\mathbf{x}}_t) - f_2^{(i)}(t, \hat{\mathbf{x}}_t)\mathbf{u},$$

and parameters $\alpha \in \mathbb{R}^p$, $\epsilon \in \mathbb{R}_{\geq 0}$ and $\gamma \in \mathbb{R}_{\geq 0}$ are selected as in [29, Section IV]. Intuitively, $\mathbf{p}_{k,t}$ is the k^{th} projected outcome of the random variable \mathbf{x}_{t+1} and $\mathbf{H}_k^{(i)}$ quantifies the variation of predictor $f^{(i)}$ with respect to its previous k^{th} value. Notice that the objective function G is convex in \mathbf{u} and therefore online problems (P2) are tractable. In addition, if ℓ_2 has a constant gradient almost everywhere, then Assumption IV.2 on accessible gradients holds and (P2) is equivalent to (P1).

Problem 2: (Online resource allocation) We consider an online resource allocation problem with a switch, where a decision maker aims to make online resource allocation decisions in an uncertain environment. This problem is in form (P) and its objective is

$$\ell(\mathbf{u}, \mathbf{x}) = \max\{0, 1 - \langle \mathbf{u}, \phi(\mathbf{x}) \rangle\}, \quad \phi: \mathbb{R}^n \rightarrow \mathbb{R}^m,$$

where ϕ is an affine feature map selected in advance. The decision maker updates the decision \mathbf{u} online when $\langle \mathbf{u}, \phi(\mathbf{x}) \rangle < 1$, otherwise switches off. Notice that this type of objective functions appears in many classification problems. In particular, we assume that the system f is independent from the allocation variable, i.e., $f_2 \equiv 0$. See Section VI-B for a more explicit problem formulation involving resource allocation with an assignment switch.

Then, problem (P2) has the objective function

$$G(t, \mathbf{u}) = \frac{1}{T} \sum_{k \in \mathcal{T}} \max\{0, 1 - \langle \mathbf{u}, \phi(\mathbf{p}_{k,t}) \rangle\} + q_t L(\mathbf{u}),$$

where time-dependent parameters $\mathbf{p}_{k,t} \in \mathbb{R}^n$, $q_t \in \mathbb{R}$ are

$$\mathbf{p}_{k,t} = \hat{\mathbf{x}}_{k+1} + \sum_{i=1}^p \alpha_i \left(f_1^{(i)}(t, \hat{\mathbf{x}}_t) - f_1^{(i)}(k, \hat{\mathbf{x}}_k) \right), \quad \forall k, t, \\ q_t = \epsilon + \frac{\gamma}{T} \sum_{i=1}^p \sum_{k \in \mathcal{T}} \|f_1^{(i)}(k, \hat{\mathbf{x}}_k) - f_1^{(i)}(t, \hat{\mathbf{x}}_t)\|, \quad \forall t,$$

with $\alpha \in \mathbb{R}^p$, $\epsilon \in \mathbb{R}_{\geq 0}$ and $\gamma \in \mathbb{R}_{\geq 0}$ as in [29, Section IV]. We characterize the function $L(\mathbf{u})$ by subgradients of the loss function ℓ .

Lemma IV.2 (Quantification of L) Consider $\ell(\mathbf{u}, \mathbf{x}) := \max\{0, 1 - \langle \mathbf{u}, \phi(\mathbf{x}) \rangle\}$, where $\phi(\mathbf{x})$ is differentiable in \mathbf{x} . Then, the function $L(\mathbf{u})$ is

$$L(\mathbf{u}) = \sup_{\mathbf{g} \in \partial_{\mathbf{x}} \ell(\mathbf{u}, \mathbf{x}), \mathbf{x} \in \mathbb{R}^n} \|\mathbf{g}\|,$$

where the set $\partial_{\mathbf{x}} \ell(\mathbf{u}, \mathbf{x})$ contains all the subgradients of ℓ at \mathbf{x} , given any \mathbf{u} in advance, i.e.,

$$\partial_{\mathbf{x}} \ell(\mathbf{u}, \mathbf{x}) := h(\mathbf{x}, \mathbf{u}) \cdot \frac{\partial \phi}{\partial \mathbf{x}}(\mathbf{x})\mathbf{u},$$

where

$$h(\mathbf{x}, \mathbf{u}) = \begin{cases} -1, & \text{if } \langle \mathbf{u}, \phi(\mathbf{x}) \rangle < 1 \\ 0, & \text{if } \langle \mathbf{u}, \phi(\mathbf{x}) \rangle > 1 \\ [-1, 0], & \text{otherwise} \end{cases}$$

In particular, if $\phi(\mathbf{x}) := J\mathbf{x}$ for some matrix $J \in \mathbb{R}^{m \times n}$, then $L(\mathbf{u}) = \|J^\top \mathbf{u}\|$. If \mathbf{x} is contained in a compact set X , then $L(\mathbf{u}) = \text{Lip}(\phi)\|\mathbf{u}\|$, where $\text{Lip}(\phi) \in \mathbb{R}_{\geq 0}$ is the Lipschitz constant of ϕ on X .

Lemma IV.2 indicates that, given a properly selected feature mapping ϕ , the objective G is convex in \mathbf{u} and therefore online problems (P2) are convex and tractable. In addition, if ϕ is a linear map almost everywhere, then Assumption IV.2 on accessible gradients holds and (P2) is equivalent to (P1).

V. ONLINE ALGORITHMS

Online convex problems (P2) are non-smooth due to the normed regularization terms in G . To achieve fast, online solutions, we propose a two-step procedure. First, we follow [35], [36] to obtain a smooth version of (P2), called (P2'). Then, we extend the Nesterov's accelerated-gradient method [37]—known to achieve an optimal first-order convergence rate for smooth and offline convex problems—to solve the problem (P2'). Finally, we quantify the dynamic regret [4] of online decisions w.r.t. solutions of (P1) in probability.

Step 1: (Smooth approximation of (P2)) To simplify the discussion, let us use the generic notation $F: \mathcal{U} \rightarrow \mathbb{R}$ for a convex and potentially non-smooth function, which can represent any particular component of the objective function $G(t, \mathbf{u})$ of (P2) at time t .

Definition V.1 (Smoothable function [35]) We call a convex function $F(\mathbf{u})$ smoothable on \mathcal{U} if there exists a $a > 0$ such that, for every $\mu > 0$, there is a continuously differentiable convex function $F_\mu: \mathcal{U} \rightarrow \mathbb{R}$ satisfying

- (1) $F_\mu(\mathbf{u}) \leq F(\mathbf{u}) \leq F_\mu(\mathbf{u}) + a\mu$, for all $\mathbf{u} \in \mathcal{U}$.
- (2) There exists $b > 0$ such that F_μ has a Lipschitz gradient over \mathcal{U} with Lipschitz constant b/μ , i.e.,

$$\|\nabla F_\mu(\mathbf{u}_1) - \nabla F_\mu(\mathbf{u}_2)\| \leq \frac{b}{\mu} \|\mathbf{u}_1 - \mathbf{u}_2\|, \quad \forall \mathbf{u}_1, \mathbf{u}_2 \in \mathcal{U}.$$

To obtain a smooth approximation F_μ of F , we follow the Moreau proximal approximation technique [35], described as in the following lemma.

Lemma V.1 (Moreau-Yosida approximation) Given a convex function $F: \mathcal{U} \rightarrow \mathbb{R}$ and any $\mu > 0$, let us denote by $\partial F(\mathbf{u})$ the set of subgradients of F at \mathbf{u} , respectively. Let $D := \sup_{\mathbf{g} \in \partial F(\mathbf{u}), \mathbf{u} \in \mathcal{U}} \|\mathbf{g}\|^2 < +\infty$. Then, F is smoothable

with parameters $(a, b) := (D/2, 1)$, where the smoothed version $F_\mu : \mathcal{U} \rightarrow \mathbb{R}$ is the Moreau approximation:

$$F_\mu(\mathbf{u}) := \inf_{\mathbf{z} \in \mathcal{U}} \left\{ F(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{u}\|^2 \right\}, \quad \mathbf{u} \in \mathcal{U}.$$

In addition, if F is M -strongly convex with some $M > 0$, then F_μ is $M/(1 + \mu M)$ -strongly convex. And further, the minimization of $F(\mathbf{u})$ over $\mathbf{u} \in \mathcal{U}$ is equivalent to that of $F_\mu(\mathbf{u})$ over $\mathbf{u} \in \mathcal{U}$ in the sense that the set of minimizers of two problems are the same.

From the definition of the smoothable function, we know that: 1) a positive linear combination of smoothable functions is smoothable⁶, and 2) the composition of a smoothable function with a linear transformation is smoothable⁷. These properties enable us to smooth each component of G , i.e., ℓ_1 , ℓ_2 , h and $\|\cdot\|$, which results in a smooth approximation of (P2) via the corresponding G_μ as follows

$$\min_{\mathbf{u} \in \mathcal{U}} G_\mu(t, \mathbf{u}). \quad (\text{P2}')$$

Note that G_μ is locally Lipschitz and minimizers of (P2') are that of (P2). We provide in the following lemma explicit expressions of (P2') for the two problem classes.

Lemma V.2 (Examples of (P2'))

Problem 1: Consider the following loss function

$$\ell(\mathbf{u}, \mathbf{x}) := \frac{1}{2} \|\mathbf{u}\|^2 + F_\mu(\mathbf{x}), \quad \text{given some } \mu > 0,$$

where $F_\mu : \mathbb{R}^n \rightarrow \mathbb{R}$ is a smoothed ℓ_2 -norm function⁸, with $\text{Lip}(F_\mu) = 1$. Then, the objective function $G_\mu(t, \mathbf{u})$ is

$$\frac{1}{2} \|\mathbf{u}\|^2 + \frac{1}{T} \sum_{k \in \mathcal{T}} F_\mu(\mathbf{p}_{k,t}) + \epsilon + \frac{\gamma}{T} \sum_{i=1}^p \sum_{k \in \mathcal{T}} F_\mu(\mathbf{H}_k^{(i)}),$$

where \mathbf{p}, \mathbf{H} are affine in \mathbf{u} , defined as in Section IV. In addition, we have the smoothing parameter of $G_\mu(t, \mathbf{u})$, $(a, b) := ((1 + p\gamma)/2, \mu + s_0 + \gamma \sum_i s_i)$, where

$$s_0 = \sigma_{\max} \left(\left(\sum_{i=1}^p \alpha_i f_2^{(i)}(t, \hat{\mathbf{x}}_t) \right)^\top \left(\sum_{i=1}^p \alpha_i f_2^{(i)}(t, \hat{\mathbf{x}}_t) \right) \right),$$

⁶If F_1 is smoothable with parameter (a_1, b_1) and F_2 with parameter (a_2, b_2) , then $c_1 F_1 + c_2 F_2$ is smoothable with parameter $(c_1 a_1 + c_2 a_2, c_1 b_1 + c_2 b_2)$, for any $c_1, c_2 \geq 0$.

⁷Let $A : \mathcal{U} \rightarrow \mathcal{X}$ be a linear transformation and let $\mathbf{b} \in \mathcal{X}$. Let $\ell : \mathcal{X} \rightarrow \mathbb{R}$ be a smoothable function with parameter (a, b) . Then, the function $F : \mathcal{U} \rightarrow \mathbb{R}$, $\mathbf{u} \mapsto \ell(A\mathbf{u} + \mathbf{b})$ is smoothable with parameter $(a, b\|A\|^2)$, where $\|A\| := \max_{\|\mathbf{u}\|=1} \|A\mathbf{u}\|$. If $\mathcal{X} = \mathbb{R}$, then $\|A\|$ is the ℓ_∞ norm.

⁸**The ℓ_2 -norm function:** Consider $\mathbf{x} \in \mathbb{R}^n$, $F : \mathbf{x} \mapsto \|\mathbf{x}\|$, and $\mu > 0$. Clearly, F is differentiable almost everywhere, except at the origin. Then,

$$\begin{aligned} F_\mu(\mathbf{x}) &:= \min_{\mathbf{z} \in \mathbb{R}^n} \left\{ \|\mathbf{z}\| + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{x}\|^2 \right\}, \\ &= \min_{r \geq 0} \min_{\|\mathbf{z}\|=r} \left\{ r + \frac{1}{2\mu} (r^2 - 2\mathbf{z}^\top \mathbf{x} + \|\mathbf{x}\|^2) \right\}, \\ &= \min_{r \geq 0} \left\{ r + \frac{1}{2\mu} (r^2 - 2r\|\mathbf{x}\| + \|\mathbf{x}\|^2) \right\}, \\ &= \begin{cases} \frac{\|\mathbf{x}\|^2}{2\mu}, & \text{if } \|\mathbf{x}\| \leq \mu, \\ \|\mathbf{x}\| - \frac{\mu}{2}, & \text{otherwise,} \end{cases} \end{aligned}$$

with the smoothing parameter $(1/2, 1)$.

with σ_{\max} denoting the maximum singular value of the matrix, and

$$s_i = \sigma_{\max} \left(f_2^{(i)}(t, \hat{\mathbf{x}}_t)^\top f_2^{(i)}(t, \hat{\mathbf{x}}_t) \right), \quad i \in \{1, \dots, p\}.$$

Problem 2: Let us select the feature map ϕ to be the identity map with the dimension $m = n$, and consider

$$\ell(\mathbf{u}, \mathbf{x}) := \max\{0, 1 - \langle \mathbf{u}, \mathbf{x} \rangle\}, \quad \text{with } L(u) = \|\mathbf{u}\|,$$

resulting in

$$G_\mu(t, \mathbf{u}) = \frac{1}{T} \sum_{k \in \mathcal{T}} F_\mu^S(\langle \mathbf{u}, \mathbf{p}_{k,t} \rangle) + q_t F_\mu(\mathbf{u}),$$

where $\mu > 0$, parameters \mathbf{p}, q are as in Section IV, and functions $F_\mu^S : \mathbb{R} \rightarrow \mathbb{R}$ and $F_\mu : \mathbb{R}^n \rightarrow \mathbb{R}$ are the smoothed switch function⁹ and ℓ_2 -norm function⁸, respectively. Note that G_μ has the smoothing parameter $(a, b) := ((1 + q_t)/2, q_t + 1/T \sum_{k \in \mathcal{T}} \|\mathbf{p}_{k,t}\|_\infty^2)$. ■

Step 2: (Solution to (P2') as a dynamical system) To solve (P2') online, we propose a dynamical system extending the Nesterov's accelerated-gradient method by adapting gradients of the time-varying objective function. In particular, let $\mathbf{u}_t, t \in \mathbb{Z}_{\geq 0}$, be solutions of (P2') and let us consider the solution system with some $\mathbf{u}_0 \in \mathcal{U}$ and $\mathbf{y}_0 = \mathbf{u}_0$, as

$$\begin{aligned} \mathbf{u}_{t+1} &= \Pi_{\mathcal{U}}(\mathbf{y}_t - \varepsilon_t \nabla G_\mu(t, \mathbf{y}_t)), \\ \mathbf{y}_{t+1} &= \mathbf{u}_{t+1} + \eta_t(\mathbf{u}_{t+1} - \mathbf{u}_t), \end{aligned} \quad (5)$$

where $\varepsilon_t \leq \mu/b_t$ with positive parameters μ and $b_t := b$ being those define $G_\mu(t, \mathbf{u})$. We denote by ∇G_μ the derivative of G_μ w.r.t. its second argument and denote by $\Pi_{\mathcal{U}}(\mathbf{y})$ the projection of \mathbf{y} onto \mathcal{U} as in Assumption IV.4 on convex decision oracle. Note that, the gradient function ∇G_μ can be computed in closed form for problems of interest, see, e.g., Appendix A for those of the proposed problems. Further, we select the moment coefficient $\eta_t \in \mathbb{R}_{\geq 0}$ as in Appendix B. In the following, we leverage Appendix B on the stability analysis of the solution system (5) for a regret bound between online decisions and optimal solutions of (P1).

⁹**The Switch function:** Consider $u \in \mathbb{R}$, $F^S : u \mapsto \max\{0, 1 - u\}$, which is differentiable almost everywhere. For a given $\mu > 0$, we compute

$$\begin{aligned} F_\mu^S(u) &:= \min_{z \in \mathbb{R}} \left\{ \max\{0, 1 - z\} + \frac{1}{2\mu} \|z - u\|^2 \right\}, \\ &= \min \left\{ \min_{z \leq 1} 1 - z + \frac{1}{2\mu} \|z - u\|^2, \min_{z \geq 1} \frac{1}{2\mu} \|z - u\|^2 \right\}. \end{aligned}$$

Given that

$$\min_{z \leq 1} 1 - z + \frac{1}{2\mu} \|z - u\|^2 = \begin{cases} \frac{1}{2\mu} \|1 - u\|^2, & \text{if } u > 1 - \mu, \\ 1 - u - \frac{\mu}{2}, & \text{if } u \leq 1 - \mu, \end{cases}$$

and

$$\min_{z \geq 1} \frac{1}{2\mu} \|z - u\|^2 = \begin{cases} \frac{1}{2\mu} \|1 - u\|^2, & \text{if } u < 1, \\ 0, & \text{if } u \geq 1, \end{cases}$$

resulting in

$$F_\mu^S(u) := \begin{cases} 1 - u - \frac{\mu}{2}, & \text{if } u \leq 1 - \mu, \\ \frac{1}{2\mu} \|1 - u\|^2, & \text{if } 1 - \mu \leq u < 1, \\ 0, & \text{if } u \geq 1, \end{cases}$$

with the smoothing parameter $(1/2, 1)$.

Theorem V.1 (Probabilistic regret bound of (P1)) Given any $t \geq 2$, let us denote by \mathbf{u}_t and \mathbf{u}_t^* the decision generated by (5) and an optimal solution which solves the online Problem (P1), respectively. Consider the dynamic regret to be the difference of the cost expected to incur if we implement \mathbf{u}_t instead of \mathbf{u}_t^* , defined as

$$R_t := \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}_t, \mathbf{x})] - \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}_t^*, \mathbf{x})].$$

Then, the regret R_t is bounded in probability as follows

$$\text{Prob} \left(R_t \leq \frac{4W_t}{(t+2)^2} + TF_t + a\mu + 2L(\mathbf{u}_t^* \hat{\epsilon}) \geq \rho(t), \right)$$

where W_t depends on the system state at time $t - T$, and F_t depends on the variation of the optimal objective values in \mathcal{T} , i.e.,

$$F_t = \max_{k \in \mathcal{T}} \{ |G_{k+1}^* - G_k^*| \} + \bar{L},$$

where $G_k^* := G(k, \mathbf{u}_k^*)$ is the optimal objective value of (P2), or equivalently that of (P1). Further, \bar{L} is the variation bound of G w.r.t. time, and the rest of the parameters are the same as before. Furthermore, if all historical data are assimilated for the decision \mathbf{u}_t , then, we have

$$\liminf_{t \rightarrow \infty} \text{Prob} (R_t \leq TF_t + a\mu) \geq 1 - \beta,$$

with β a given, arbitrary confidence value.

Theorem V.1 quantifies the dynamic regret of online decisions \mathbf{u} w.r.t. solutions to (P1) in high probability. Notice that, the regret bound is dominated by terms: TF_t , $a\mu$ and $L(\mathbf{u}_t^*)\hat{\epsilon}$, which mainly depend on three factors: the data-driven parameters ε , η and μ of the solution system (5), the variation F_t over optimal objective values, and the parameters T , β , γ and $\hat{\epsilon}$ that are related to the system and environment learning. In practice, a small regret bound is determined by 1) an effective learning procedure which contributes to small $\hat{\epsilon}$; 2) a proper selection of the loss function ℓ which results in smoothing procedure with a small parameter $a\mu$; and 3) the problem structure leading to small variations F_t of the optimal objectives values. Furthermore, when we use all the historical data for the objective gradients in the solution system (5), the effect of system ambiguity learning is negligible asymptotically.

Online Procedure: Our online algorithm is summarized in the Algorithm 1.

Online Optimization and Learning Algorithm 1 Opal(\mathcal{I})

- 1: Select $\{f^{(i)}\}_i$, ℓ , β , \mathcal{U} , \mathbf{u}_0 , μ , and $t = 1$;
 - 2: **repeat**
 - 3: Update data set $\mathcal{I} := \mathcal{I}_t$;
 - 4: Compute $\boldsymbol{\alpha} := \boldsymbol{\alpha}_t$ as in [29];
 - 5: Select $\hat{\mathbb{P}}_{t+1|t}$ in (2) and $\hat{\epsilon} := \hat{\epsilon}(t, T, \beta, \boldsymbol{\alpha}, \mathbf{u})$ in (4);
 - 6: Run dynamical system (5) for $\mathbf{u} := \mathbf{u}_t$;
 - 7: Apply \mathbf{u} to (P) with the regret guarantee;
 - 8: $t \leftarrow t + 1$;
 - 9: **until** time t stops.
-

VI. IMPLEMENTATION

In this section, we apply our algorithm to the introduced motivating examples, resulting in online-tractable, effective system learning with guaranteed, regret-bounded performance in high probability.

A. Optimal control of an uncertain nonlinear system

We consider the two-wheel vehicle driving under various road conditions, and our goal is to learn one-step prediction of the system state distribution and leverage for path tracking under various unknown road zones. In particular, we represent the two-wheel vehicle as a differential-drive robot subject to uncertainty [38]:

$$\begin{aligned} x_{t+1} &= x_t + h \cos(\theta_t) d_{1,t} + h w_{1,t}, \\ y_{t+1} &= y_t + h \sin(\theta_t) d_{1,t} + h w_{2,t}, \\ \theta_{t+1} &= \theta_t - h d_{2,t} + h w_{3,t}, \\ d_{1,t} &= \frac{r}{2} (v_{l,t} + v_{r,t} + e_{1,t}), \\ d_{2,t} &= \frac{r}{2R} (v_{l,t} - v_{r,t} + e_{2,t}), \end{aligned} \quad (6)$$

where components of states $\mathbf{x}_t := (x_t, y_t, \theta_t) \in \mathbb{R}^2 \times [-\pi, \pi] \cong \mathbb{R} \times \mathbb{S}^1$ represent vehicle position and orientation on the 2-D plane. We take the discretization parameter $h = 0.01$ and assume subGaussian uncertainty $\mathbf{w}_t := (w_{1,t}, w_{2,t}, w_{3,t}) \in \mathbb{R}^3$ to be a zero-mean, mixture of Gaussian and Uniform distributions with $\sigma = 0.5$. The intermediate variable $\mathbf{d}_t := (d_{1,t}, d_{2,t})$ depends on the wheel radius $r = 0.15$ m, the distance between wheels $R = 0.4$ m, the controlled left-right wheel speed $\mathbf{u}_t := (v_{l,t}, v_{r,t})$ and an unknown parameter $\mathbf{e}_t := (e_{1,t}, e_{2,t})$, which depends on the wheel quality and road conditions. For simplicity, we assume that the planner adapts the system (6) with $\mathbf{e}_t \equiv (0, 0)$ and $\mathbf{w}_t \equiv (0, 0, 0)$, and the vehicle can move over three types of road zones, the regular zone with $\mathbf{e}^{(1)} := (0, 0)$, the slippery zone with $\mathbf{e}^{(2)} = (4, 0)$, and the sandy zone with $\mathbf{e}^{(3)} = (-1.2, -0.2)$, where locations of these zones are described in Fig. 2.

To adapt the proposed approach, we consider Problem (P) with the following loss function

$$\begin{aligned} \ell(\mathbf{u}, x, y, \theta) &= \frac{1}{20} \|\mathbf{u} - \mathbf{u}^{\text{ref}}\|^2 + \frac{1}{14\sqrt{2}} |x - x^{\text{ref}}| + \\ &\quad \frac{1}{4\sqrt{2}} |y - y^{\text{ref}}| + \frac{289}{8} (\cos(\theta) - \cos(\theta^{\text{ref}}))^2 + \\ &\quad \frac{289}{8} (\sin(\theta) - \sin(\theta^{\text{ref}}))^2, \end{aligned}$$

where $(\mathbf{u}^{\text{ref}}, x^{\text{ref}}, y^{\text{ref}}, \theta^{\text{ref}})$ are signals generated by the planner, and we select the parameter $\mu = 10^{-4}$ for components which are not smooth. In addition, we assume $\mathcal{U} = [-20, 20]^2$ and utilize $p = 3$ basis functions $\{f^{(i)}\}_i$ in form of (6), with $\mathbf{w}_t \equiv (0, 0, 0)$, and

$$\begin{aligned} i = 1, & \quad e_1 = 0, \quad e_2 = 0, \\ i = 2, & \quad e_1 = 10, \quad e_2 = 0, \\ i = 3, & \quad e_1 = 0, \quad e_2 = 10. \end{aligned}$$

Note that the ground truth parameter $\boldsymbol{\alpha}^* := (1, 0, 0)$ in the regular zone, $\boldsymbol{\alpha}^* := (0.6, 0.4, 0)$ in the slippery zone, and $\boldsymbol{\alpha}^* := (1.14, -0.12, -0.02)$ in the sandy zone. At each time t , we have access to model sets $\{f^{(i)}\}_i$ and as well as the real-time data set \mathcal{I}_t with size $T_0 = 100$, which corresponds to the moving time window of order 0.1 second. For the system learning algorithm, notions of norm and inner product are those defined on the vector space $T(\mathbb{R}^2 \times \mathbb{S}) \cong \mathbb{R}^3$. We employ our online optimization and learning algorithm for

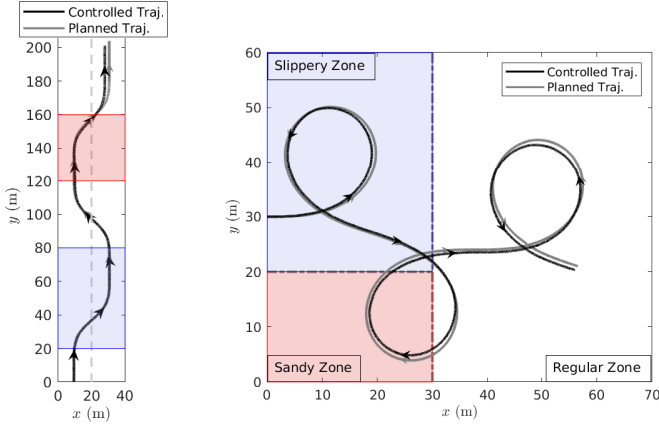


Fig. 3: An example of the (gray) planned trajectory and (black) controlled system trajectory in various road zones, with the system state $\mathbf{x} = (x, y, \theta)$. The red region indicates sandy zone while the blue region indicates the slippery zone. With the implemented control, the vehicle follows the planned path with low regrets in high probability.

the characterization of the uncertain vehicle states, learning of the unknown road-condition parameter e , and control towards planned behaviors in real time. The achieved system behaviors are demonstrated in Fig 3, contrasted with the case without the proposed approach, as in Fig. 2. In the following, we analyze each case separately and notice how the proposed approach strikes balance between the given planned control \mathbf{u}^{ref} and the actual control \mathbf{u} which reduces the weighted tracking error in road uncertainty.

Example (Lane-changing behavior adaptation) In this scenario, we assume the initial system state $\mathbf{x}_0 = (10, 0, \pi/2)$. Further, the vehicle can access path plan in Fig. 2(a) and as well as the suggested wheel speed plan as the gray signal in Fig. 4(a). To demonstrate the learning effect of the algorithm, we show in Fig. 5 components α_1 and α_2 of $\alpha = (\alpha_1, \alpha_2, \alpha_3)$, where the black lines indicate value of the ground truth α^* on the planned trajectory and the gray lines represent the learned, real-time estimate of α_1 and α_2 at the actual vehicle position. Notice that α^* is inaccessible in practice, and from this case study, the proposed approach indeed learns the system dynamics effectively. See, e.g. [29] for more analysis regarding to the effect of the learning behavior and ambiguity sets characterization on the selection of ϵ and γ .

As the proposed loss function ℓ measures the weighted tracking error, the resulting control system trajectory in Fig. 3(a) already reveals the effectiveness of the method and as well as the low regrets in probability. On the other hand, because the system is highly non-linear and uncertain, evaluating the actual optimal objective value of Problem (P) is difficult. Therefore, it's very challenging to evaluate the regret R_t in practice, even though the its probabilistic bounded is proved. Here, we provide in Fig. 4(b) the realized loss ℓ and as well as the realized objective value of Problem (P2), where the loss ℓ reveals one possible objective value of (P), and the objective value of (P2) serves as an upper-bound of that of (P) in high probability. In addition, notice that the derived (black) control signal in Fig. 4(a) has undesirable, high-oscillatory behavior. This is because the chosen loss function ℓ is only locally convex in \mathbf{x} . When the system disturbances are significant,

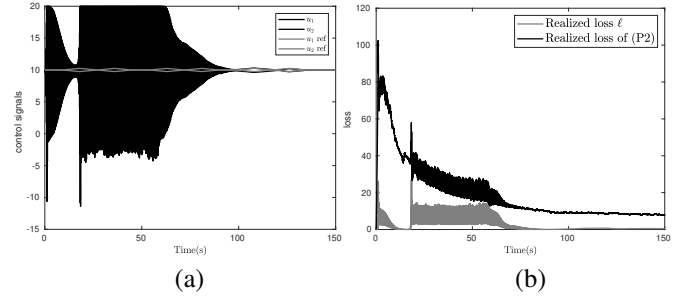


Fig. 4: (a) The (gray) control signal provided by the planner and an example of the (black) control signal derived from the proposed approach. (b) The realized loss ℓ and the achieved objective of (P2).

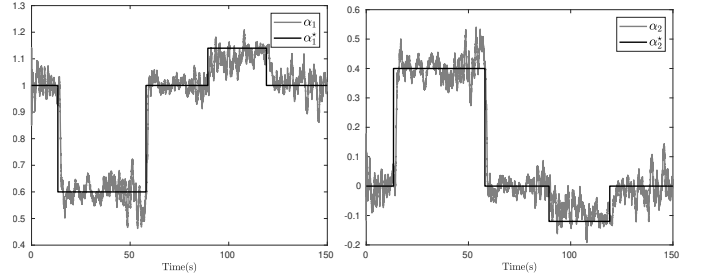


Fig. 5: The component α_1 and α_2 of the real-time parameter $\alpha := (\alpha_1, \alpha_2, \alpha_3)$ in the learning procedure.

the proposed approach then revealed certain degradation and control being oscillatory. Nevertheless, a desirable system behavior in Fig. 3(a) is achieved.

Example (Circular route tracking) In this scenario, we consider $\mathbf{x}_0 = (0, 30, 0)$. We omit the details as the analysis shares the same spirit as the last lane-changing example.

B. Online resource allocation problem

We consider an online resource allocation problem where an agent or decision maker aims to 1) achieve at least target profit under uncertainty, and 2) allocate resources as uniformly as possible. To do this, the agent distributes available resources, e.g., wealth, time, energy or human resources, to various projects or assets. In particular, for the trading-market motivating example, let us consider that the agent tries to make an online allocation $\mathbf{u} \in \mathcal{U}$ of a unit wealth to three assets. At each time t , the agent receives random return rates $\mathbf{x}_t \in \mathbb{R}_{\geq 0}^3$ of assets from some unknown and uncertain dynamics

$$\mathbf{x}_{t+1} = \mathbf{x}_t + hA(t) + h\mathbf{w}_t, \text{ with some } \mathbf{x}_0 \in \mathbb{R}^3, \quad (7)$$

where $h = 10^{-3}$ is a stepsize, the vector $A(t)$ is randomly generated, unknown and piecewise constant, and the uncertainty vector \mathbf{w}_t is assumed to be sub-Gaussian with $\sigma = 0.1$. Note that this model can serve to characterize a wide class of dynamic (linear and nonlinear) systems. In addition, we assume that the third asset is value preserved, i.e., the third component of $A(t)$ and \mathbf{w}_t are zero and $x_3 \equiv 1$. Over time, an example of the resulting unit return rates \mathbf{x} is demonstrated in Fig. 6. Then, we denote by $r_0 = 1.3$ and $\langle \mathbf{u}, \mathbf{x}_{t+1} \rangle$ the target profit and the predicted instantaneous profit, respectively. Note that the decision maker aims to obtain at least a 30% profit and allocate resources online for this purpose. In particular, the decision maker implements an

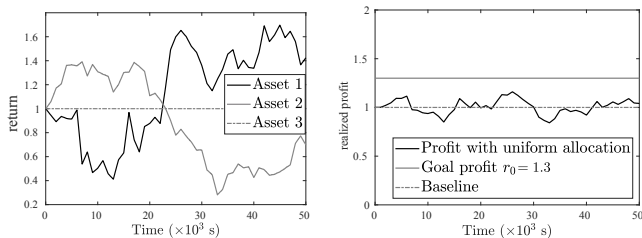


Fig. 6: An example of random returns $\mathbf{x} = (x_1, x_2, x_3)$, where returns of the first two assets $x_1, x_2 \in [0, +\infty)$ are highly fluctuating and the third is value-preserving with return $x_3 \equiv 1$. Without asset allocation, agent does not achieve the goal profit $r_0 = 1.3$ and has a chance of losing assets.

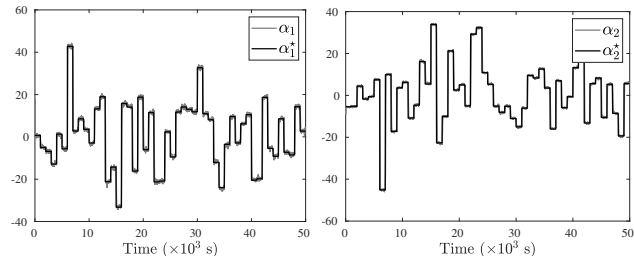


Fig. 7: The component α_1 and α_2 of the real-time parameter $\alpha := (\alpha_1, \alpha_2, \alpha_3)$ in learning, where the values α_1^* and α_2^* are the online-inaccessible ground truth. Notice the responsive behavior of the proposed learning algorithm.

allocation online if $\langle \mathbf{u}, \mathbf{x}_{t+1} \rangle \leq r_0$, otherwise does nothing. This results in (P) with the loss function

$$\ell(\mathbf{u}, \mathbf{x}) = \max\{0, 1 - \frac{1}{r_0} \langle \mathbf{u}, \mathbf{x} \rangle\},$$

and set \mathcal{U} a unit simplex. We propose $p = 3$ basis functions

$$f^{(1)} = \mathbf{x}, f^{(2)} = \mathbf{x} + 0.1h\mathbf{e}_1, f^{(3)} = \mathbf{x} + 0.1h\mathbf{e}_2,$$

where $\mathbf{e}_1 = (1, 0, 0)^\top$ and $\mathbf{e}_2 = (0, 1, 0)^\top$. At each t , we assume that only historical data are available for online resource allocations. Applying the proposed probabilistic characterization of \mathbf{x}_{t+1} as in (P1), we equivalently write it as in form (P2'), where

$$G_\mu(t, \mathbf{u}) = \frac{1}{T} \sum_{k \in \mathcal{T}} F_\mu^S(\langle \mathbf{u}, \frac{\mathbf{p}_{k,t}}{r_0} \rangle) + \frac{q_t}{r_0} F_\mu(\mathbf{u}), \mu = 0.01,$$

with functions F_μ^S and F_μ , and real-time data $\mathbf{p}_{k,t}$ and q_t determined as in Problem 2. We claim that $G_\mu(t, \mathbf{u})$ has a time-dependent Lipschitz gradient constant in \mathbf{u} given by $\text{Lip}(G_\mu) = q_t/r_0 + 1/(r_0^2 T) \sum_{k \in \mathcal{T}} \|\mathbf{p}_{k,t}\|_\infty^2$, and we use $\varepsilon := 1/\text{Lip}(G_\mu)$ in the solution system (5) to compute the online decisions.

Fig. 7 shows the real-time evolution α_1 and α_2 of the parameter $\alpha := (\alpha_1, \alpha_2, \alpha_3)$, while the behavior of α_3 can be similarly characterized. In these figures, black lines α_1^* and α_2^* are determined by the unknown signal $A(t)$ while gray lines α_1 and α_2 are those computed as in [29]. Note that α^* represents the unknown dynamics f and they are not accessible in reality. It can be seen that the proposed method effectively learns α^* .

Fig. 8 demonstrates the online resource allocation obtained by implementing (5) and the achieved real-time profit $\langle \mathbf{u}, \mathbf{x} \rangle$. The decision \mathbf{u} starts from the uniform allocation $\mathbf{u}_0 =$

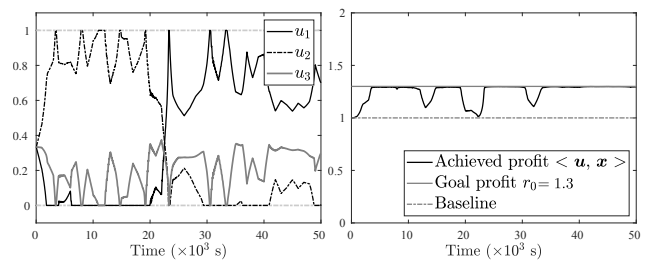


Fig. 8: Real-time resource allocation \mathbf{u} and profit $\langle \mathbf{u}, \mathbf{x} \rangle$. Notice how the decision $\mathbf{u} = (u_1, u_2, u_3)$ respects constraints and how the allocation tries to balance the assets when the goal profit r_0 is met.

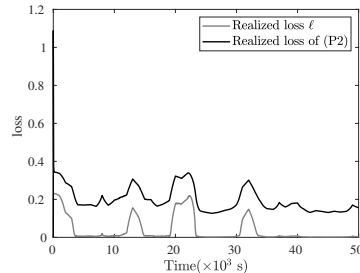


Fig. 9: The realized loss ℓ and the achieved objective of (P2).

$(1/3, 1/3, 1/3)$ and is then adjusted to approach the target profit $r_0 = 1.3$. Once the target is achieved, the agent then maintains the profit while trying to balance the allocation if possible. When the return rate \mathbf{x} is low/unbalanced, as in Fig. 6, the agent tries to improve and achieve the target profit by allocating resources more aggressively. Though did not appear in the current scenario, in case that the return rate is high and the target profit value is achieved, the agent focuses on balancing the allocation while maintaining the profit. If both the target profit and allocation balance are achieved, then the agent stops re-allocating resources and monitors the return rate \mathbf{x} until the switch turns on, e.g., when the near future profit prediction drops below r_0 again. In addition, notice how the target profit was achieved with the proposed control strategy as demonstrated in Fig. 8, which contrasts with uniform allocation case as in Fig. 6.

Fig. 9 demonstrates the evaluation of the time-varying loss ℓ as well as the realized objective value of Problem (P2). Due to the unknown time-varying distributions $\mathbb{P}_{t|t-1}$, the evaluation of the objective values of Problem (P) is intractable, and the realized loss of (P2) serves as a high-confidence upper bound of that of (P). Nevertheless, the target profit is achieved with low regret in high confidence, as revealed in Fig. 8.

VII. CONCLUSIONS

In this paper, we proposed a unified solution framework for online learning and optimization problems in form of (P). The proposed method allowed us to learn an unknown and uncertain dynamic system, while providing a characterization of the system with online-quantifiable probabilistic guarantees that certify the performance of online decisions. The approach provided tractable, online convex version of (P), via a series of equivalent reformulation techniques. We explicitly demonstrated the framework via two problem classes conforming to (P): an optimal control problem under uncertainty and an

online resource allocation problem. These two problem classes resulted in explicit, online and non-smooth convex optimization problems. We extended Nesterov's accelerated-gradient method to an online fashion and provided a solution system for online decision generation of (P). The quality of the online decisions were analytically certified via a probabilistic regret bound, which revealed its relation to the learning parameters and ambiguity sets. Two motivating examples applying the proposed framework were empirically tested, demonstrating the effectiveness of the proposed framework with the bounded regret guarantees in probability. We leave the relaxation of assumptions and the comparison of this work with other methods as the future work.

APPENDIX

A. Computation of the objective gradients

Let ℓ , G and G_μ be those in Lemma V.2 on examples of (P2'). We now derive $\nabla G_\mu := \nabla_{\mathbf{u}} G_\mu(t, \mathbf{u})$ as follows.

Problem 1 (Optimal control under uncertainty):

$$\nabla_{\mathbf{u}} G_\mu(t, \mathbf{u}) = \frac{1}{\mu} \mathbf{u} + \frac{1}{T} \sum_{k \in \mathcal{T}} \nabla_{\mathbf{u}} F_\mu(\mathbf{p}_{k,t}) + \frac{\gamma}{T} \sum_{i=1}^p \sum_{k \in \mathcal{T}} \nabla_{\mathbf{u}} F_\mu(\mathbf{H}_k^{(i)}),$$

where, for each $k \in \mathcal{T}$, the term $\nabla_{\mathbf{u}} F_\mu(\mathbf{p}_{k,t})$ is

$$\begin{cases} \frac{1}{\mu} \left(\sum_{i=1}^p \alpha_i f_2^{(i)}(t, \hat{\mathbf{x}}_t) \right)^\top \mathbf{p}_{k,t}, & \text{if } \|\mathbf{p}_{k,t}\| \leq \mu, \\ \frac{1}{\|\mathbf{p}_{k,t}\|} \left(\sum_{i=1}^p \alpha_i f_2^{(i)}(t, \hat{\mathbf{x}}_t) \right)^\top \mathbf{p}_{k,t}, & \text{otherwise,} \end{cases}$$

and, for $k \in \mathcal{T}$, $i \in \{1, \dots, p\}$, the term $\nabla_{\mathbf{u}} F_\mu(\mathbf{H}_k^{(i)})$ is

$$\begin{cases} -\frac{1}{\mu} (f_2^{(i)}(t, \hat{\mathbf{x}}_t))^\top \mathbf{H}_k^{(i)}, & \text{if } \|\mathbf{H}_k^{(i)}\| \leq \mu, \\ -\frac{1}{\|\mathbf{H}_k^{(i)}\|} (f_2^{(i)}(t, \hat{\mathbf{x}}_t))^\top \mathbf{H}_k^{(i)}, & \text{otherwise.} \end{cases}$$

Problem 2 (Online resource allocation):

$$\nabla_{\mathbf{u}} G_\mu(t, \mathbf{u}) = \frac{1}{T} \sum_{k \in \mathcal{T}} \nabla_{\mathbf{u}} F_\mu^S(\langle \mathbf{u}, \mathbf{p}_{k,t} \rangle) + q_t \nabla_{\mathbf{u}} F_\mu(\mathbf{u}),$$

where

$$\nabla_{\mathbf{u}} F_\mu(\mathbf{u}) := \begin{cases} \frac{1}{\mu} \mathbf{u}, & \text{if } \|\mathbf{u}\| \leq \mu, \\ \frac{1}{\|\mathbf{u}\|} \mathbf{u}, & \text{otherwise,} \end{cases}$$

and, for each $k \in \mathcal{T}$, the gradient $\nabla_{\mathbf{u}} F_\mu^S(\langle \mathbf{u}, \mathbf{p}_{k,t} \rangle)$ is

$$\begin{cases} -\mathbf{p}_{k,t}, & \text{if } \langle \mathbf{u}, \mathbf{p}_{k,t} \rangle \leq 1 - \mu, \\ -\frac{1 - \langle \mathbf{u}, \mathbf{p}_{k,t} \rangle}{\mu} \mathbf{p}_{k,t}, & \text{if } 1 - \mu \leq \langle \mathbf{u}, \mathbf{p}_{k,t} \rangle < 1, \\ 0, & \text{if } \langle \mathbf{u}, \mathbf{p}_{k,t} \rangle \geq 1. \end{cases}$$

These explicit expressions provide ingredients for the solution system. With different selections of the norm, the expression varies accordingly.

B. Stability Analysis of the Solution System

Here, we adapt dissipativity theory to address the performance of the online solution system (5). This part of the work is an online-algorithmic extension of the existing Nesterov's accelerated-gradient method and its convergence analysis in [39]–[41]. Our extension (5) inherits from the work in [40], where the difference is that gradient computations

in (5) are from time-varying objective functions in (P2'). To simplify the discussion, the notation we used in this subsection is different from that in the main body of the paper. Consider the online problem, analogous to (P2'), defined as follows

$$\min_{\mathbf{x} \in \mathcal{X}} f_t(\mathbf{x}), \quad t = 0, 1, 2, \dots \quad (8)$$

where $f_t(\mathbf{x})$ is locally Lipschitz in t with the parameter $h(\mathbf{x})$ and, at each time t , the objective function f_t are m_t -strongly convex and L_t -smooth, with $m_t \geq 0$ and $L_t > 0$. The convex set $\mathcal{X} \subset \mathbb{R}^n$ is analogous to that in Assumption IV.4 on convex decision oracle. The solution system to (8), analogous to (5), is

$$\begin{aligned} \mathbf{x}_{t+1} &= \Pi(\mathbf{y}_t - \alpha_t \nabla f_t(\mathbf{y}_t)), \\ \mathbf{y}_{t+1} &= \mathbf{x}_{t+1} + \beta_{t+1} (\mathbf{x}_{t+1} - \mathbf{x}_t), \end{aligned} \quad (9)$$

with some $\mathbf{y}_0 = \mathbf{x}_0 \in \mathcal{X}$,

where $\alpha_t \leq 1/L_t$ and β_t is selected iteratively, following

$$\delta_{-1} = 1, \quad \delta_{t+1} := \frac{1 + \sqrt{1 + 4\delta_t^2}}{2}, \quad \beta_t := \frac{\delta_{t-1} - 1}{\delta_t}.$$

Note that $\delta_t^2 - \delta_t = \delta_{t-1}^2$, $t = 0, 1, 2, \dots$. The projection $\Pi(\mathbf{x})$ at each time t is equivalently written as

$$\Pi(\mathbf{x}) = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|^2 + \alpha_t \ell(\mathbf{z}),$$

with $\ell(\mathbf{z}) = 0$ if $\mathbf{z} \in \mathcal{X}$, otherwise $+\infty$. Note that the projection operation is a convex problem with the objective function being strongly convex. Thus, $\Pi(\mathbf{x})$ is a singleton (the unique minimizer) and satisfies the optimality condition [42]

$$\mathbf{x} - \Pi(\mathbf{x}) \in \alpha_t \partial \ell(\Pi(\mathbf{x})),$$

where the r.h.s. is the sub-differential set of ℓ at $\Pi(\mathbf{x})$. Equivalently, we write the above condition as

$$\Pi(\mathbf{x}) = \mathbf{x} - \alpha_t \partial \ell(\Pi(\mathbf{x})).$$

We apply this equivalent representation to the solution system (9), resulting in

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{y}_t - \alpha_t \nabla f_t(\mathbf{y}_t) - \alpha_t \partial \ell(\mathbf{w}_t), \\ \mathbf{y}_{t+1} &= \mathbf{x}_{t+1} + \beta_{t+1} (\mathbf{x}_{t+1} - \mathbf{x}_t), \\ \mathbf{w}_t &= \mathbf{x}_{t+1}. \end{aligned} \quad (10)$$

Note that (10) is not an explicit online algorithm, as the state \mathbf{x}_{t+1} is yet to be determined. However, we leverage this equivalent reformulation for the convergence analysis of solutions to (9) to a sequence of optimizers of (8), denoted by $\{\mathbf{x}_t^*\}$. To do this, let $\mathbf{z}_t := (\mathbf{x}_t - \mathbf{x}_t^*, \mathbf{x}_{t-1} - \mathbf{x}_{t-1}^*)$ denote the tracking error vector and represent (10) as the error dynamical system

$$\mathbf{z}_{t+1} = A_t \mathbf{z}_t + B_t^u \mathbf{u}_t + B_t^v \mathbf{v}_t, \quad (11)$$

with $\mathbf{z}_1 = (\mathbf{x}_1 - \mathbf{x}_1^*, \mathbf{x}_0 - \mathbf{x}_0^*)$,

with the gradient input $\mathbf{u}_t := \nabla f_t(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t)$, the reference signal $\mathbf{v}_t := (\mathbf{x}_t^* - \mathbf{x}_{t-1}^*, \mathbf{x}_{t+1}^* - \mathbf{x}_t^*)$, the matrices

$$A_t = \begin{bmatrix} 1 + \beta_t & -\beta_t \\ 1 & 0 \end{bmatrix}, \quad B_t^u = \begin{bmatrix} -\alpha_t \\ 0 \end{bmatrix}, \quad B_t^v = \begin{bmatrix} \beta_t & -1 \\ 0 & 0 \end{bmatrix},$$

and the auxiliary variables

$$\begin{aligned} \mathbf{y}_t - \mathbf{x}_t^* &= [1 + \beta_t \quad -\beta_t] \mathbf{z}_t + [\beta_t \quad 0] \mathbf{v}_t, \\ \mathbf{w}_t - \mathbf{x}_t^* &= [1 \quad 0] \mathbf{z}_{t+1} + [0 \quad 1] \mathbf{v}_t. \end{aligned}$$

We provide the following stability analysis of the system.

Theorem A.1 (Stability of (9)) Consider the solution algorithm (9), or equivalently (10).

(1) For each $t \geq 1$, we have the following

$$f_t(\mathbf{x}_t) - f_t(\mathbf{x}_{t+1}) \geq \boldsymbol{\xi}_t^\top X_{1,t} \boldsymbol{\xi}_t,$$

$$f_t(\mathbf{x}_t^*) - f_t(\mathbf{x}_{t+1}) \geq \boldsymbol{\xi}_t^\top X_{2,t} \boldsymbol{\xi}_t.$$

Here, $\boldsymbol{\xi}_t := (\mathbf{z}_t, \mathbf{u}_t, \mathbf{v}_t)$, and

$$X_{1,t} := \frac{1}{2} \begin{pmatrix} m\beta^2 & -m\beta^2 & -\beta & m\beta^2 & 0 \\ -m\beta^2 & m\beta^2 & \beta & -m\beta^2 & 0 \\ -\beta & \beta & \alpha(2-L\alpha) & -\beta & 0 \\ m\beta^2 & -m\beta^2 & -\beta & m\beta^2 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

$$X_{2,t} := \frac{1}{2} \begin{pmatrix} m(1+\beta)^2 & -\eta & -(1+\beta) & \eta & 0 \\ -\eta & m\beta^2 & \beta & -m\beta^2 & 0 \\ -(1+\beta) & \beta & \alpha(2-L\alpha) & -\beta & 0 \\ \eta & -m\beta^2 & -\beta & m\beta^2 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

with $\eta = m(1+\beta)\beta$ and the parameters (m, L, α, β) are a short-hand notation for $(m_t, L_t, \alpha_t, \beta_t)$.

(2) Given the horizon parameter $T_0 \in \mathbb{Z}_{>0}$ with $T = \min\{t-1, T_0\}$. Then, for any $t \geq 2$, the solution \mathbf{x}_t from (9) achieves

$$f_t(\mathbf{x}_t) - f_t(\mathbf{x}_t^*) \leq \frac{4G_t}{(t+2)^2} + TF_t + TK_t$$

$$+ \frac{4(t-T-1+\delta_0)^2}{(t+2)^2} (f_{t-T}(\mathbf{x}_{t-T}) - f_{t-T}(\mathbf{x}_{t-T}^*)).$$

where the time-dependent parameters G_t , F_t and K_t are determined by f_t , α_t and β_t .

Proof of Theorem A.1: (1) By the m -strong convexity and L -smoothness of f , we have

$$f(\mathbf{x}) - f(\mathbf{y}) \geq \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) + \frac{m}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad (12)$$

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \nabla f(\mathbf{y})^\top (\mathbf{y} - \mathbf{x}) - \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (13)$$

(1a) Consider (12) with $(\mathbf{x}, \mathbf{y}) \equiv (\mathbf{x}_t, \mathbf{y}_t)$. We leverage $\mathbf{y}_t = \mathbf{x}_t + \beta(\mathbf{x}_t - \mathbf{x}_{t-1})$ and the distributive law¹⁰ for

$$f(\mathbf{x}_t) - f(\mathbf{y}_t)$$

$$\geq \beta \nabla f(\mathbf{y}_t)^\top (\mathbf{x}_{t-1} - \mathbf{x}_t) + \frac{m\beta^2}{2} \|\mathbf{x}_{t-1} - \mathbf{x}_t\|^2,$$

$$= \beta (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t))^\top (\mathbf{x}_{t-1} - \mathbf{x}_t - \mathbf{x}_{t-1}^* + \mathbf{x}_t^*)$$

$$+ \frac{m\beta^2}{2} \|\mathbf{x}_{t-1} - \mathbf{x}_t - \mathbf{x}_{t-1}^* + \mathbf{x}_t^*\|^2$$

$$+ \beta (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t))^\top (\mathbf{x}_{t-1}^* - \mathbf{x}_t^*)$$

$$- \beta \partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_{t-1} - \mathbf{x}_t)$$

$$+ m\beta^2 (\mathbf{x}_{t-1} - \mathbf{x}_t - \mathbf{x}_{t-1}^* + \mathbf{x}_t^*)^\top (\mathbf{x}_{t-1}^* - \mathbf{x}_t^*)$$

$$+ \frac{m\beta^2}{2} \|\mathbf{x}_{t-1}^* - \mathbf{x}_t^*\|^2.$$

¹⁰Apply 1) $a^\top c = (a+b)^\top (c-d) + (a+b)^\top d - b^\top c$ and 2) $c^\top c = (c-d)^\top (c-d) + 2(c-d)^\top d + d^\top d$, with $a = \nabla f(\mathbf{y}_t)$, $b = \partial \ell(\mathbf{w}_t)$, $c = \mathbf{x}_{t-1} - \mathbf{x}_t$, $d = \mathbf{x}_{t-1}^* - \mathbf{x}_t^*$,

We re-organize the the right-hand-side into the matrix form as

$$\frac{1}{2} \boldsymbol{\delta}_t^\top \begin{pmatrix} m\beta^2 & -m\beta^2 & -\beta & m\beta^2 \\ -m\beta^2 & m\beta^2 & \beta & -m\beta^2 \\ -\beta & \beta & 0 & -\beta \\ m\beta^2 & -m\beta^2 & -\beta & m\beta^2 \end{pmatrix} \boldsymbol{\delta}_t - \beta \partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_{t-1} - \mathbf{x}_t),$$

with $\boldsymbol{\delta}_t^\top := (\mathbf{x}_t - \mathbf{x}_t^*, \mathbf{x}_{t-1} - \mathbf{x}_{t-1}^*, \nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t), \mathbf{x}_t^* - \mathbf{x}_{t-1}^*)$.

(1b) Consider (13) with $(\mathbf{x}, \mathbf{y}) \equiv (\mathbf{x}_{t+1}, \mathbf{y}_t)$. We leverage $\mathbf{x}_{t+1} = \mathbf{y}_t - \alpha \nabla f_t(\mathbf{y}_t) - \alpha \partial \ell(\mathbf{w}_t)$ and the distribution law, resulting in

$$f(\mathbf{y}_t) - f(\mathbf{x}_{t+1}) \geq \alpha \nabla f(\mathbf{y}_t)^\top (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t))$$

$$- \frac{L\alpha^2}{2} \|\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t)\|^2,$$

$$= \frac{\alpha(2-L\alpha)}{2} \|\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t)\|^2$$

$$- \alpha \partial \ell(\mathbf{w}_t)^\top (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t)).$$

Now, we sum the terms involving $\partial \ell(\mathbf{w}_t)$ in the r.h.s. of inequalities in (1a) and (1b), leverage (10), and then apply the convexity of ℓ , $\mathbf{x}_t \in \mathcal{X}$ and $\mathbf{w}_t = \mathbf{x}_{t+1} \in \mathcal{X}$, to obtain the following

$$- \beta \partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_{t-1} - \mathbf{x}_t) - \alpha \partial \ell(\mathbf{w}_t)^\top (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t))$$

$$= -\partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_t - \mathbf{w}_t) \geq \ell(\mathbf{w}_t) - \ell(\mathbf{x}_t) = 0,$$

which results in $f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) \geq \boldsymbol{\xi}_t^\top X_{1,t} \boldsymbol{\xi}_t$.

Note that we have identified (f, m, L, α, β) with $(f_t, m_t, L_t, \alpha_t, \beta_t)$, and note that $\nabla f_t(\mathbf{x}_t^*) + \partial \ell(\mathbf{x}_t^*) = 0$.

(1c) Similarly, consider (12) with $(\mathbf{x}, \mathbf{y}) \equiv (\mathbf{x}_t^*, \mathbf{y}_t)$. From $\mathbf{y}_t = \mathbf{x}_t + \beta(\mathbf{x}_t - \mathbf{x}_{t-1})$ and the distributive law,

$$f(\mathbf{x}_t^*) - f(\mathbf{y}_t)$$

$$\geq \nabla f(\mathbf{y}_t)^\top (\mathbf{x}_t^* - \mathbf{y}_t) + \frac{m}{2} \|\mathbf{x}_t^* - \mathbf{y}_t\|^2,$$

$$= (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t))^\top (\mathbf{x}_t^* - \mathbf{y}_t + \beta \mathbf{x}_t^* - \beta \mathbf{x}_{t-1}^*)$$

$$+ \frac{m}{2} \|\mathbf{x}_t^* - (1+\beta)(\mathbf{x}_t - \mathbf{x}_t^*) + \beta(\mathbf{x}_{t-1} - \mathbf{x}_{t-1}^*)\|^2$$

$$- \beta (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t))^\top (\mathbf{x}_t^* - \mathbf{x}_{t-1}^*) - \partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_t^* - \mathbf{y}_t)$$

$$- m\beta [-(1+\beta)(\mathbf{x}_t - \mathbf{x}_t^*) + \beta(\mathbf{x}_{t-1} - \mathbf{x}_{t-1}^*)]^\top (\mathbf{x}_t^* - \mathbf{x}_{t-1}^*)$$

$$+ \frac{m\beta^2}{2} \|\mathbf{x}_t^* - \mathbf{x}_{t-1}^*\|^2,$$

$$= \frac{1}{2} \boldsymbol{\delta}_t^\top \begin{pmatrix} m(1+\beta)^2 & -\eta & -(1+\beta) & \eta \\ -\eta & m\beta^2 & \beta & -m\beta^2 \\ -(1+\beta) & \beta & 0 & -\beta \\ \eta & -m\beta^2 & -\beta & m\beta^2 \end{pmatrix} \boldsymbol{\delta}_t$$

$$- \partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_t^* - \mathbf{y}_t),$$

with $\eta = m(1+\beta)\beta$. We add this inequality to that in (1b) and leverage

$$- \partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_t^* - \mathbf{y}_t) - \alpha \partial \ell(\mathbf{w}_t)^\top (\nabla f(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t))$$

$$= -\partial \ell(\mathbf{w}_t)^\top (\mathbf{x}_t^* - \mathbf{w}_t) \geq \ell(\mathbf{w}_t) - \ell(\mathbf{x}_t^*) = 0,$$

resulting in $f(\mathbf{x}_t^*) - f(\mathbf{x}_{t+1}) \geq \boldsymbol{\xi}_t^\top X_{2,t} \boldsymbol{\xi}_t$.

(2) Let us define the time varying function

$$V_t(\mathbf{z}_t) := \begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t^* - \mathbf{x}_{t-1}^* \end{bmatrix}^\top H_t \begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t^* - \mathbf{x}_{t-1}^* \end{bmatrix},$$

where we take

$$H_t := \frac{1}{2\alpha_{t-1}} \begin{bmatrix} \delta_{t-1} \\ 1 - \delta_{t-1} \\ \delta_{t-1} \end{bmatrix} [\delta_{t-1}, \quad 1 - \delta_{t-1}, \quad \delta_{t-1}],$$

with $\{\alpha_t\}_t$ those in the solution system (9) and $\{\delta_t\}_t$ the sequence of scalars which defines $\{\beta_t\}_t$. Now, verify

$$V_{t+1}(\mathbf{z}_{t+1}) - \frac{\alpha_{t-1}}{\alpha_t} V_t(\mathbf{z}_t) = \boldsymbol{\xi}_t^\top J_t \boldsymbol{\xi}_t,$$

where $\boldsymbol{\xi}_t := (\mathbf{z}_t, \mathbf{u}_t, \mathbf{v}_t)$, which are those define (11), resulting in $\boldsymbol{\xi}_t := (\mathbf{x}_t - \mathbf{x}_t^*, \mathbf{x}_{t-1} - \mathbf{x}_{t-1}^*, \nabla f_t(\mathbf{y}_t) + \partial \ell(\mathbf{w}_t), \mathbf{x}_t^* - \mathbf{x}_{t-1}^*, \mathbf{x}_{t+1}^* - \mathbf{x}_t^*)$ and

$$J_t = \frac{1}{2\alpha_t} \begin{pmatrix} 0 & 0 & -\alpha_t \delta_t \delta_{t-1} & -\delta_{t-1} & 0 \\ 0 & 0 & \alpha_t \beta_t \delta_t^2 & \beta_t \delta_t & 0 \\ -\alpha_t \delta_t \delta_{t-1} & \alpha_t \beta_t \delta_t^2 & \alpha_t^2 \delta_t^2 & -\alpha_t \beta_t \delta_t^2 & 0 \\ -\delta_{t-1} & \beta_t \delta_t & -\alpha_t \beta_t \delta_t^2 & 1 - 2\delta_{t-1} & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Let us compute

$$M_t := \delta_{t-1}^2 \mathbf{X}_{1,t} + \delta_t \mathbf{X}_{2,t} \\ = \frac{1}{2} \begin{pmatrix} m_t(\delta_t^2 - 1) & -m_t \beta_t \delta_t \delta_{t-1} & -\delta_t \delta_{t-1} & m_t \beta_t \delta_t \delta_{t-1} & 0 \\ -m_t \beta_t \delta_t \delta_{t-1} & m_t \beta_t^2 \delta_t^2 & \beta_t \delta_t^2 & -m_t \beta_t^2 \delta_t^2 & 0 \\ -\delta_t \delta_{t-1} & \beta_t \delta_t^2 & \alpha_t(2 - L_t \alpha_t) \delta_t^2 & -\beta_t \delta_t^2 & 0 \\ m_t \beta_t \delta_t \delta_{t-1} & -m_t \beta_t^2 \delta_t^2 & -\beta_t \delta_t^2 & m_t \beta_t^2 \delta_t^2 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

and then achieve

$$\boldsymbol{\xi}_t^\top (J_t - M_t) \boldsymbol{\xi}_t = \begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t^* - \mathbf{x}_{t-1}^* \end{bmatrix}^\top N_{1,t} \begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t^* - \mathbf{x}_{t-1}^* \end{bmatrix} \\ + \begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t^* - \mathbf{x}_{t-1}^* \end{bmatrix}^\top N_{2,t} \begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t^* - \mathbf{x}_{t-1}^* \end{bmatrix} \\ - \alpha_t (1 - L_t \alpha_t) \mathbf{u}_t^\top \mathbf{u}_t,$$

with, for each $t \geq 1$,

$$N_{1,t} := \frac{1}{2} \begin{pmatrix} -m_t(\delta_t^2 - 1) & m_t \beta_t \delta_t \delta_{t-1} & -m_t \beta_t \delta_t \delta_{t-1} \\ m_t \beta_t \delta_t \delta_{t-1} & -m_t \beta_t^2 \delta_t^2 & m_t \beta_t^2 \delta_t^2 \\ -m_t \beta_t \delta_t \delta_{t-1} & m_t \beta_t^2 \delta_t^2 & -m_t \beta_t^2 \delta_t^2 \end{pmatrix}, \\ \cong \frac{m_t}{2} \begin{pmatrix} -(\delta_t^2 - 1) & \beta_t \delta_t \delta_{t-1} & 0 \\ \beta_t \delta_t \delta_{t-1} & -\beta_t^2 \delta_t^2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \leq 0,$$

and, using the fact that $\delta_t > (t+1)/2, \forall t \geq 0$, we have

$$N_{2,t} := \frac{1}{2} \begin{pmatrix} 0 & 0 & -\delta_{t-1} \\ 0 & 0 & \beta_t \delta_t \\ -\delta_{t-1} & \beta_t \delta_t & 1 - 2\delta_{t-1} \end{pmatrix} \leq 0.$$

Then, if we select $\alpha_t \leq 1/L_t$, it results in

$$\boldsymbol{\xi}_t^\top (J_t - M_t) \boldsymbol{\xi}_t \leq 0.$$

We rewrite it as

$$V_{t+1}(\mathbf{z}_{t+1}) - \frac{\alpha_{t-1}}{\alpha_t} V_t(\mathbf{z}_t) \leq \boldsymbol{\xi}_t^\top M_t \boldsymbol{\xi}_t, \\ \leq \delta_{t-1}^2 (f_t(\mathbf{x}_t) - f_t(\mathbf{x}_{t+1})) + \delta_t (f_t(\mathbf{x}_t^*) - f_t(\mathbf{x}_{t+1})), \\ = -\delta_t^2 (f_t(\mathbf{x}_{t+1}) - f_t(\mathbf{x}_t^*)) + \delta_{t-1}^2 (f_t(\mathbf{x}_t) - f_t(\mathbf{x}_t^*)).$$

As f_t being locally Lipschitz in t , there exists a non-negative function $h(\mathbf{x})$ such that

$$f_{t+1}(\mathbf{x}_{t+1}) - f_t(\mathbf{x}_{t+1}) \leq h(\mathbf{x}_{t+1}),$$

resulting in

$$V_{t+1}(\mathbf{z}_{t+1}) - \frac{\alpha_{t-1}}{\alpha_t} V_t(\mathbf{z}_t) \\ \leq -\delta_t^2 (f_{t+1}(\mathbf{x}_{t+1}) - f_{t+1}(\mathbf{x}_{t+1}^*)) + \delta_{t-1}^2 (f_t(\mathbf{x}_t) - f_t(\mathbf{x}_t^*)) \\ - \delta_t^2 (f_{t+1}(\mathbf{x}_{t+1}^*) - f_t(\mathbf{x}_t^*)) + \delta_t^2 h(\mathbf{x}_{t+1}), \forall t$$

Summing up the above set of inequalities over the moving horizon window $t \in \mathcal{T} = \{t-1, \dots, t-T\}$, where $T = \min\{t-1, T_0\}$ with some $T_0 \in \mathbb{Z}_{>0}$, we obtain

$$V_t(\mathbf{z}_t) + \sum_{k \in \mathcal{T}} (1 - \frac{\alpha_{k-1}}{\alpha_k}) V_k(\mathbf{z}_k) - V_{t-T}(\mathbf{z}_{t-T}) \\ \leq -\delta_{t-1}^2 (f_t(\mathbf{x}_t) - f_t(\mathbf{x}_t^*)) \\ + \delta_{t-T-1}^2 (f_{t-T}(\mathbf{x}_{t-T}) - f_{t-T}(\mathbf{x}_{t-T}^*)) \\ - \sum_{k \in \mathcal{T}} \delta_k^2 (f_{k+1}(\mathbf{x}_{k+1}^*) - f_k(\mathbf{x}_k^*)) + \sum_{k \in \mathcal{T}} \delta_k^2 h(\mathbf{x}_{k+1}).$$

Let us denote by G_t , K_t , and F_t , respectively, the horizon accumulated potential, the bound of the locally Lipschitz function h , and the variation bound of the optimal objective values. That is,

$$G_t := V_{t-T}(\mathbf{z}_{t-T}) - V_t(\mathbf{z}_t) - \sum_{k \in \mathcal{T}} (1 - \frac{\alpha_{k-1}}{\alpha_k}) V_k(\mathbf{z}_k),$$

$$K_t := \max_{k \in \mathcal{T}} \{h(\mathbf{x}_{k+1})\},$$

$$F_t := \max_{k \in \mathcal{T}} \{|f_{k+1}(\mathbf{x}_{k+1}^*) - f_k(\mathbf{x}_k^*)|\}.$$

Then, using the fact that (1) $\delta_{t-1} \geq (t+2)/2$, for all $t \geq 0$; (2) $\delta_{t-T-1} \leq t-T-1 + \delta_0$ with $\delta_0 = (1 + \sqrt{5})/2$, and (3) δ_t is monotonically increasing, we have

$$f_t(\mathbf{x}_t) - f_t(\mathbf{x}_t^*) \leq \frac{4G_t}{(t+2)^2} + TF_t + TK_t \\ + \frac{4(t-T-1 + \delta_0)^2}{(t+2)^2} (f_{t-T}(\mathbf{x}_{t-T}) - f_{t-T}(\mathbf{x}_{t-T}^*)).$$

Note that, when $t \leq T_0 + 1$, we have $T = t - 1$. This gives

$$f_t(\mathbf{x}_t) - f_t(\mathbf{x}_t^*) \leq \frac{4G_t}{(t+2)^2} + (t-1)F_t + (t-1)K_t \\ + \frac{4\delta_0^2}{(t+2)^2} (f_1(\mathbf{x}_1) - f_1(\mathbf{x}_1^*)). \blacksquare$$

C. Proofs of lemmas and theorems

Proof of Lemma IV.1: By the definition of the ambiguity set, we have that, for any distribution $\mathbb{Q} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})$

$$d_W(\mathbb{Q}, \hat{\mathbb{P}}_{t+1|t}) \leq \hat{\epsilon},$$

which, by Kantorovich-Rubinstein Theorem, is equivalent to

$$\int_{\mathcal{Z}} h(\mathbf{x}) \mathbb{Q}(d\mathbf{x}) - \int_{\mathcal{Z}} h(\mathbf{x}) \hat{\mathbb{P}}_{t+1|t}(d\mathbf{x}) \leq \hat{\epsilon}, \quad \forall h \in \mathcal{L},$$

where \mathcal{L} is the set of functions with Lipschitz constant 1 and \mathcal{Z} is the support of the random variable \mathbf{x} . For a given \mathbf{u} , let us select h to be

$$h(\mathbf{x}) := \frac{\ell(\mathbf{u}, \mathbf{x})}{L(\mathbf{u})},$$

where L is the positive Lipschitz function as in Assumption IV.1. Substituting h to the above inequality, we have

$$\int_{\mathcal{Z}} \ell(\mathbf{u}, \mathbf{x}) \mathbb{Q}(d\mathbf{x}) - \int_{\mathcal{Z}} \ell(\mathbf{u}, \mathbf{x}) \hat{\mathbb{P}}_{t+1|t}(d\mathbf{x}) \leq \hat{\varepsilon}L(\mathbf{u}),$$

or equivalently

$$\mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{u}, \mathbf{x})] \leq \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] + \hat{\varepsilon}L(\mathbf{u}).$$

As the inequality holds for every $\mathbb{Q} \in \mathcal{P}_{t+1}$, therefore

$$\begin{aligned} \sup_{\mathbb{Q} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{u}, \mathbf{x})] \\ \leq \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] + \hat{\varepsilon}(t, T, \beta, \boldsymbol{\alpha}, \mathbf{u})L(\mathbf{u}). \quad \blacksquare \end{aligned}$$

Proof of Theorem IV.1: We show this by constructing a distribution in the ambiguity set. By Assumption IV.2 on convex and gradient-accessible functions, there exists an index $j \in \mathcal{T}$ such that the derivative $\nabla_{\mathbf{x}}\ell(\mathbf{u}, \mathbf{x})$ at $(\mathbf{u}, \bar{\mathbf{x}}^{(j)})$, $\bar{\mathbf{x}}^{(j)} := \sum_{i=1}^p \alpha_i \xi_j^{(i)}(\boldsymbol{\alpha}, \mathbf{u})$, satisfies

$$\|\nabla_{\mathbf{x}}\ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)})\| = L(\mathbf{u}).$$

Now using this index j , we construct a parameterized distribution as follows

$$\mathbb{Q}(\Delta\mathbf{x}) = \frac{1}{T} \sum_{k \in \mathcal{T}, k \neq j} \delta_{\{\sum_{i=1}^p \alpha_i \xi_k^{(i)}(\boldsymbol{\alpha}, \mathbf{u})\}} + \frac{1}{T} \delta_{\{\bar{\mathbf{x}}^{(j)} + \Delta\mathbf{x}\}},$$

where $\Delta\mathbf{x} \in \mathbb{R}^n$ with $\|\Delta\mathbf{x}\| \leq T\hat{\varepsilon}$. By the definition of the ambiguity set and, since the support of the distribution \mathbb{P} is $\Xi_{t+1} = \mathbb{R}^n$, we have $\mathbb{Q}(\Delta\mathbf{x}) \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})$.

Next, we quantify the lower bound of the following term

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}(\Delta\mathbf{x})}[\ell(\mathbf{u}, \mathbf{x})] - \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] \\ = \frac{1}{T} \left(\ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)} + \Delta\mathbf{x}) - \ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)}) \right). \end{aligned}$$

By Assumption IV.2 on the convexity of ℓ on \mathbf{x} , we have

$$\ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)} + \Delta\mathbf{x}) - \ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)}) \geq \nabla_{\mathbf{x}}\ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)})^\top \Delta\mathbf{x}.$$

Then, by selecting

$$\Delta\mathbf{x} := \frac{T\hat{\varepsilon}\nabla_{\mathbf{x}}\ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)})}{\|\nabla_{\mathbf{x}}\ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)})\|},$$

we have

$$\nabla_{\mathbf{x}}\ell(\mathbf{u}, \bar{\mathbf{x}}^{(j)})^\top \Delta\mathbf{x} = T\hat{\varepsilon}L(\mathbf{u}).$$

These bounds result in

$$\mathbb{E}_{\mathbb{Q}(\Delta\mathbf{x})}[\ell(\mathbf{u}, \mathbf{x})] - \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] \geq \hat{\varepsilon}L(\mathbf{u}).$$

As $\mathbb{Q}(\Delta\mathbf{x}) \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})$, therefore

$$\sup_{\mathbb{Q} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})} \mathbb{E}_{\mathbb{Q}}[\ell(\mathbf{u}, \mathbf{x})] \geq \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] + \hat{\varepsilon}L(\mathbf{u}).$$

Finally, with Assumption IV.1 on Lipschitz loss functions and Lemma IV.1 on an upper bound of (P1), we equivalently write Problem (P1) as

$$\inf_{\mathbf{u} \in \mathcal{U}} \mathbb{E}_{\hat{\mathbb{P}}_{t+1|t}(\boldsymbol{\alpha}, \mathbf{u})}[\ell(\mathbf{u}, \mathbf{x})] + \hat{\varepsilon}(t, T, \beta, \boldsymbol{\alpha}, \mathbf{u})L(\mathbf{u}),$$

which is the Problem (P2). \blacksquare

Proof of Lemma IV.2: This is the direct application of the definition of the local Lipschitz condition. \blacksquare

Proof of Lemma V.1: First, we have

$$F_{\mu}(\mathbf{u}) \leq F(\mathbf{u}) + \frac{1}{2\mu}\|\mathbf{u} - \mathbf{u}\|^2 = F(\mathbf{u}), \quad \forall \mathbf{u} \in \mathcal{U}.$$

Then, we compute

$$\begin{aligned} F(\mathbf{u}) - F_{\mu}(\mathbf{u}) &= \sup_{\mathbf{z} \in \mathcal{U}} \left\{ F(\mathbf{u}) - F(\mathbf{z}) - \frac{1}{2\mu}\|\mathbf{z} - \mathbf{u}\|^2 \right\}, \\ &\leq \sup_{\mathbf{z} \in \mathcal{U}} \left\{ \mathbf{g}(\mathbf{u})^\top (\mathbf{u} - \mathbf{z}) - \frac{1}{2\mu}\|\mathbf{z} - \mathbf{u}\|^2 \right\}, \\ &\leq \sup_{\mathbf{z}} \left\{ \mathbf{g}(\mathbf{u})^\top (\mathbf{u} - \mathbf{z}) - \frac{1}{2\mu}\|\mathbf{z} - \mathbf{u}\|^2 \right\}, \\ &\leq \frac{\mu}{2} \mathbf{g}(\mathbf{u})^\top \mathbf{g}(\mathbf{u}) \leq \frac{D}{2} \mu, \end{aligned}$$

where the equality comes from the definition of $F_{\mu}(\mathbf{u})$, the first inequality leverages the convexity of F , the second one relaxes the constraint set, the third one applies the achieved optimizer $\mathbf{z}^* = \mathbf{u} - \mu\mathbf{g}(\mathbf{u})$, and the last one is from the boundedness of subgradients.

Further, given F as described, it is well-known (see, e.g., [43, Proposition 12.15] for details) that F_{μ} is convex and continuously differentiable where its gradient ∇F_{μ} is Lipschitz continuous with constant $1/\mu$. In addition, the minimizer $\mathbf{z}^*(\mathbf{u})$ of F_{μ} is achievable and unique, resulting in an explicit gradient expression of F_{μ} as follows

$$\nabla F_{\mu}(\mathbf{u}) = \frac{1}{\mu}(\mathbf{u} - \mathbf{z}^*(\mathbf{u})).$$

In addition, we claim that, if F is M -strongly convex, F_{μ} is $M/(1 + \mu M)$ -strongly convex, following [44, Theorem 2.2]. Finally, we equivalently write the minimization problem as follows

$$\begin{aligned} \min_{\mathbf{u} \in \mathcal{U}} F_{\mu}(\mathbf{u}) &= \min_{\mathbf{u} \in \mathcal{U}} \min_{\mathbf{z} \in \mathcal{U}} \left\{ F(\mathbf{z}) + \frac{1}{2\mu}\|\mathbf{z} - \mathbf{u}\|^2 \right\} \\ &= \min_{\mathbf{z} \in \mathcal{U}} \min_{\mathbf{u} \in \mathcal{U}} \left\{ F(\mathbf{z}) + \frac{1}{2\mu}\|\mathbf{z} - \mathbf{u}\|^2 \right\} \\ &= \min_{\mathbf{z} \in \mathcal{U}} F(\mathbf{z}), \end{aligned}$$

where the first line applies the achievability of the minimizer of the problem that defines F_{μ} , the second switches the minimization operators, the third applies the fact that $\mathbf{u} = \mathbf{z}$ solves the inner problem. This concludes that any \mathbf{u} that minimizes F_{μ} also minimizes F , and vice versa. \blacksquare

Proof of Theorem V.1: Let us consider the solution system (5). At each time t , let us select $\varepsilon := \varepsilon_t = 1/\text{Lip}(G_{\mu})$, or equivalently, μ/b with $b = \max_{k \in \mathcal{T}} b_k$. Let η_t satisfy

$$\delta_{-1} = 1, \quad \delta_{t+1} := \frac{1 + \sqrt{1 + 4\delta_t^2}}{2}, \quad \eta_t := \frac{\delta_{t-1} - 1}{\delta_t}.$$

Then, by Theorem A.1 with $t \geq 2$, the following holds

$$G_{\mu}(t, \mathbf{u}_t) - G_{\mu}(t, \mathbf{u}_t^*) \leq \frac{4W_t}{(t+2)^2} + TF_t, \quad (14)$$

where \mathbf{u}_t^* is a solution to (P2'), $T = \min\{t-1, T_0\}$ with some horizon parameter $T_0 \in \mathbb{Z}_{>0}$. Notice that T_0 is the length of the used historical data whenever such data are available. The time-varying parameter W_t depends on the initial objective discrepancy and the accumulated energy storage in the considered time horizon \mathcal{T} , and F_t is the variation bound of the optimal objective values in \mathcal{T} . Specifically, we have

$$F_t = \max_{k \in \mathcal{T}} \{|G_\mu(k+1, \mathbf{u}_{k+1}^*) - G_\mu(k, \mathbf{u}_k^*)|\} + \bar{L},$$

with \bar{L} the variation bound of $G_\mu(t, \mathbf{u}_t)$ w.r.t. time t . Let us consider the storage function $V_t(\mathbf{z}_t) := \mathbf{z}_t^\top H_t \mathbf{z}_t$, where $\mathbf{z}_t := (\mathbf{u}_t - \mathbf{u}_t^*, \mathbf{u}_{t-1} - \mathbf{u}_{t-1}^*, \mathbf{u}_t^* - \mathbf{u}_{t-1}^*)$ and

$$H_t := \frac{1}{2\varepsilon_{t-1}} \begin{bmatrix} \delta_{t-1} \\ 1 - \delta_{t-1} \\ \delta_{t-1} \end{bmatrix} [\delta_{t-1} \quad 1 - \delta_{t-1} \quad \delta_{t-1}] \succeq 0.$$

Then we have

$$\begin{aligned} W_t &= V_{t-T}(\mathbf{z}_{t-T}) - V_t(\mathbf{z}_t) - \sum_{k \in \mathcal{T}} \left(1 - \frac{\varepsilon_{k-1}}{\varepsilon_k}\right) V_k(\mathbf{z}_k) \\ &\quad + (t - T - 1 + \delta_0)^2 (f_{t-T}(\mathbf{x}_{t-T}) - f_{t-T}(\mathbf{x}_{t-T}^*)), \end{aligned}$$

where the first two term is the energy decrease in the horizon \mathcal{T} ; the third sum term indicates the instantaneous energy change, which depends on the online, estimated Lipschitz constant; the last term depends on the goodness of the initial decision at the beginning of the current \mathcal{T} . Note how the selection of ε_t and T affect W_t (or G_t in Theorem A.1). In the most conservative scenario, we select $\varepsilon_t := \min\{\varepsilon_{t-1}, \mu/b_t\}$ and $T_0 = \infty$, which results in a constant upper bound of W_t as follows

$$W_t \leq V_1(\mathbf{z}_1) + \delta_0^2 (f_1(\mathbf{x}_1) - f_1(\mathbf{x}_1^*)),$$

therefore, in this case, the bound (14) essentially depends on the growing term $(t-1)F_t$. A less conservative way is to use moving horizon strategy, with $\varepsilon_t := \min\{\varepsilon_{t-1}, \mu/b_t\}$ but a finite T_0 . Then, as t is sufficiently large, we have

$$W_t \leq V_{t-T}(\mathbf{z}_{t-T}) + t^2 (f_{t-T}(\mathbf{x}_{t-T}) - f_{t-T}(\mathbf{x}_{t-T}^*)),$$

where, in this case, the bound (14) essentially depends on F_t and $f_{t-T}(\mathbf{x}_{t-T}) - f_{t-T}(\mathbf{x}_{t-T}^*)$.

Now, we consider for any $t \geq 2$. By Definition V.1, there exists a constant $a > 0$ such that

$$G(t, \mathbf{u}_t) - a\mu \leq G_\mu(t, \mathbf{u}_t),$$

and by Lemma V.1, we have that \mathbf{u}_t^* is a minimizer of (P2') if and only if it is that of (P2), and

$$G_\mu(t, \mathbf{u}_t^*) \equiv G(t, \mathbf{u}_t^*).$$

This results in

$$G(t, \mathbf{u}_t) - G(t, \mathbf{u}_t^*) \leq \frac{4W_t}{(t+2)^2} + TF_t + a\mu, \quad (15)$$

with an equivalent expression of F_t as

$$F_t = \max_{k \in \mathcal{T}} \{|G_{k+1}^* - G_k^*|\} + \bar{L},$$

where $G_k^* := G(k, \mathbf{u}_k^*)$ is the optimal objective value of (P2) or, later we see, equivalent to that of (P1).

Next, by Theorem IV.1 on the equivalence of (P1) and (P2), \mathbf{u}_t^* is a minimizer of (P2) if and only if it is also that of (P1), and

$$G(t, \mathbf{u}_t^*) \equiv \sup_{\mathbb{Q} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u}_t^*)} \mathbb{E}_{\mathbb{Q}} [\ell(\mathbf{u}_t^*, \mathbf{x})]. \quad (16)$$

Further, as in Section IV, we claim that Problem (P1) provides a probabilistic bound for the objective of (P), resulting in

$$\text{Prob}(\mathbb{P}_{t+1|t} \in \mathcal{P}_{t+1}) \geq \rho(t), \text{ or equivalently,} \quad (17)$$

$$\text{Prob}(\mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}_t, \mathbf{x})] \leq G(t, \mathbf{u}_t)) \geq \rho(t), \quad (18)$$

with $\rho(t)$ as in Theorem III.1. Then by (17), we know that $\mathbb{P}_{t+1|t} \in \mathcal{P}_{t+1}$ if and only if $d_W(\mathbb{P}_{t+1|t}, \hat{\mathbb{P}}_{t+1|t}) \leq \hat{\varepsilon}$ where $\hat{\varepsilon}$ is selected as in Theorem III.1. Further, since d_W is a metric, for any $\mathbb{Q} \in \mathcal{P}_{t+1}$, we claim

$$\begin{aligned} d_W(\mathbb{Q}, \mathbb{P}_{t+1|t}) &\leq d_W(\mathbb{Q}, \hat{\mathbb{P}}_{t+1|t}) + d_W(\mathbb{P}_{t+1|t}, \hat{\mathbb{P}}_{t+1|t}), \\ &\leq \hat{\varepsilon} + \hat{\varepsilon} \leq 2\hat{\varepsilon}. \end{aligned}$$

By Assumption IV.1 and the same proof procedure of Lemma IV.1 on the above inequality, we have, for every \mathbf{u} , the following:

$$\sup_{\mathbb{Q} \in \mathcal{P}_{t+1}(\boldsymbol{\alpha}, \mathbf{u})} \mathbb{E}_{\mathbb{Q}} [\ell(\mathbf{u}, \mathbf{x})] \leq \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}, \mathbf{x})] + 2L(\mathbf{u})\hat{\varepsilon}.$$

By taking $\mathbf{u} := \mathbf{u}_t^*$ and using (16), we have

$$G(t, \mathbf{u}_t^*) \leq \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}_t^*, \mathbf{x})] + 2L(\mathbf{u}_t^*)\hat{\varepsilon}. \quad (19)$$

We combine the inequality (15), (18) and (19), resulting in

$$\begin{aligned} \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}_t, \mathbf{x})] - \mathbb{E}_{\mathbb{P}_{t+1|t}} [\ell(\mathbf{u}_t^*, \mathbf{x})] \\ \leq \frac{4W_t}{(t+2)^2} + TF_t + a\mu + 2L(\mathbf{u}_t^*)\hat{\varepsilon}, \end{aligned}$$

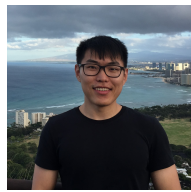
with the probability at least $\rho(t)$, holds for any $t \geq 2$. Furthermore, if all historical data are assimilated for the decision \mathbf{u}_t , i.e., we select $T_0 = \infty$ with $\varepsilon_t := \min\{\varepsilon_{t-1}, \mu/b_t\}$, then, the term W_t is upper bound by a constant and, the radius $\hat{\varepsilon}$ asymptotically goes to zero due to the selection as in [29, Section IV]. Consequently, this results in

$$\liminf_{t \rightarrow \infty} \text{Prob}(R_t \leq TF_t + a\mu) \geq 1 - \beta. \quad \blacksquare$$

REFERENCES

- [1] S. Shalev-Shwartz, "Online learning and online convex optimization," *Foundations and Trends in Machine Learning*, vol. 4, no. 2, pp. 107–194, 2012.
- [2] A. Rakhlin, K. Sridharan, and A. Tewari, "Online learning: Random averages, combinatorial parameters, and learnability," in *Advances in Neural Information Processing Systems*, pp. 1984–1992, 2010.
- [3] A. Simonetto, E. Dall'Anese, S. Paternain, G. Leus, and G. B. Giannakis, "Time-varying convex optimization: Time-structured algorithms and applications," *Proceedings of the IEEE*, vol. 108, pp. 2032–2048, November 2020.
- [4] M. Zinkevich, "Online convex programming and generalized infinitesimal gradient ascent," in *Int. Conf. on Machine Learning*, pp. 928–936, 2003.
- [5] E. Hazan, "Introduction to online convex optimization," *Foundations and Trends in Optimization*, vol. 2, no. 3-4, pp. 157–325, 2016.
- [6] A. Mokhtari, S. Shahrampour, A. Jadbabaie, and A. Ribeiro, "Online optimization in dynamic environments: Improved regret rates for strongly convex problems," in *IEEE Int. Conf. on Decision and Control*, pp. 7195–7201, 2016.
- [7] A. Jadbabaie, A. Rakhlin, S. Shahrampour, and K. Sridharan, "Online optimization: Competing with dynamic comparators," in *Artificial Intelligence and Statistics*, pp. 398–406, 2015.

- [8] A. Rakhlin and K. Sridharan, "Online learning with predictable sequences," *Journal of Machine Learning Research*, 2013.
- [9] E. Hall and R. Willett, "Online convex optimization in dynamic environments," *IEEE Journal of Selected Topics in Signal Processing*, vol. 9, no. 4, pp. 647–662, 2015.
- [10] N. Chen, A. Agarwal, A. Wierman, S. Barman, and L. Andrew, "Online convex optimization using predictions," in *ACM Int. Conf. on Measurement and Modeling of Computer Systems*, pp. 191–204, 2015.
- [11] Y. Li, G. Qu, and N. Li, "Online optimization with predictions and switching costs: Fast algorithms and the fundamental limit," *arXiv preprint arXiv:1801.07780*, 2018.
- [12] L. Ljung, *System identification*. Prentice Hall, 1999.
- [13] J. C. Willems, P. Rapisarda, I. Markovskiy, and B. D. Moor, "A note on persistency of excitation," *Systems and Control Letters*, vol. 54, no. 4, pp. 325–329, 2005.
- [14] T. Maupong and P. Rapisarda, "Data-driven control: A behavioral approach," *Systems and Control Letters*, vol. 101, pp. 37–43, 2017.
- [15] C. D. Persis and P. Tesi, "Formulas for data-driven control: Stabilization, optimality, and robustness," *IEEE Transactions on Automatic Control*, vol. 65, no. 3, pp. 909–924, 2019.
- [16] J. Coulson, J. Lygeros, and F. Dörfler, "Data-enabled predictive control: In the shallows of the DeePC," in *European Control Conference*, pp. 307–312, 2019.
- [17] J. Berberich, J. Köhler, M. Müller, and F. Allgower, "Data-driven model predictive control with stability and robustness guarantees," *IEEE Transactions on Automatic Control*, 2020.
- [18] A. Allibhoy and J. Cortés, "Data-based receding horizon control of linear network systems," *IEEE Control Systems Letters*, vol. 5, no. 4, pp. 1207–1212, 2020.
- [19] M. Nonhoff and M. A. Müller, "Online convex optimization for data-driven control of dynamical systems," vol. 1, pp. 180–193, 2022.
- [20] S. Oymak and N. Ozay, "Non-asymptotic identification of LTI systems from a single trajectory," in *American Control Conference*, pp. 5655–5661, 2019.
- [21] A. Tsiamis and G. J. Pappas, "Finite-sample analysis of stochastic system identification," in *IEEE Int. Conf. on Decision and Control*, pp. 3648–3654, 2019.
- [22] S. Fattahi, N. Matni, and S. Sojoudi, "Learning sparse dynamical systems from a single sample trajectory," in *IEEE Int. Conf. on Decision and Control*, pp. 2682–2689, 2019.
- [23] N. Fournier and A. Guillin, "On the rate of convergence in Wasserstein distance of the empirical measure," *Probability Theory and Related Fields*, vol. 162, no. 3–4, p. 707–738, 2015.
- [24] R. Gao and A. Kleywegt, "Distributionally robust stochastic optimization with Wasserstein distance," *arXiv preprint arXiv:1604.02199*, 2016.
- [25] P. Mohajerin Esfahani and D. Kuhn, "Data-driven distributionally robust optimization using the Wasserstein metric: performance guarantees and tractable reformulations," *Mathematical Programming*, vol. 171, no. 1–2, pp. 115–166, 2018.
- [26] S. Shafieezadeh-Abadeh, D. Kuhn, and P. Mohajerin Esfahani, "Regularization via mass transportation," *Journal of Machine Learning Research*, vol. 20, no. 103, pp. 1–68, 2019.
- [27] D. Boskos, J. Cortés, and S. Martínez, "Data-driven ambiguity sets with probabilistic guarantees for dynamic processes," *IEEE Transactions on Automatic Control*, vol. 66, no. 7, pp. 2991–3006, 2021.
- [28] D. Boskos, J. Cortés, and S. Martínez, "Data-driven ambiguity sets for linear systems under disturbances and noisy observations," in *American Control Conference*, (Denver, CO), pp. 4491–4496, July 2020.
- [29] D. Li, D. Fooladivanda, and S. Martínez, "Online learning of parameterized uncertain dynamical environments with finite-sample guarantees," *IEEE Control Systems Letters*, vol. 5, no. 4, pp. 1351–1356, 2021.
- [30] D. Dursvyatskiy and L. Xiao, "Stochastic optimization with decision-dependent distributions," *Mathematics of Operations Research*, 2022.
- [31] K. Wood, G. Bianchin, and E. Dall-Anese, "Online projected gradient descent for stochastic optimization with decision-dependent distributions," *IEEE Control Systems Letters*, vol. 6, pp. 1646–1651, 2022.
- [32] L. V. Kantorovich and G. S. Rubinstein, "On a space of completely additive functions," *Vestnik Leningrad. Univ.*, vol. 13, no. 7, p. 52–59, 1958.
- [33] D. Bertsekas, *Dynamic programming and optimal control: Volume I*, vol. 4. Athena scientific, 2012.
- [34] D. Li, D. Fooladivanda, and S. Martínez, "Online learning of parameterized uncertain dynamical environments with finite-sample guarantees," *ArXiv. Preprint arXiv:2009.02390*, 2020.
- [35] A. Beck and M. Teboulle, "Smoothing and first order methods: A unified framework," *SIAM Journal on Optimization*, vol. 22, no. 2, pp. 557–580, 2012.
- [36] Y. Nesterov, "Smooth minimization of non-smooth functions," *Mathematical Programming*, vol. 103, no. 1, pp. 127–152, 2005.
- [37] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*, vol. 87. Springer Science & Business Media, 2013.
- [38] S. M. LaValle, *Planning algorithms*. Cambridge University Press, 2006.
- [39] B. Hu and L. Lessard, "Dissipativity theory for Nesterov's accelerated method," *arXiv preprint arXiv:1706.04381*, 2017.
- [40] L. Lessard, B. Recht, and A. Packard, "Analysis and design of optimization algorithms via integral quadratic constraints," *SIAM Journal on Optimization*, vol. 26, no. 1, pp. 57–95, 2016.
- [41] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [42] R. T. Rockafellar and R. J.-B. Wets, *Variational analysis*. Springer, 1998.
- [43] H. H. Bauschke and P. L. Combettes, *Convex analysis and monotone operator theory in Hilbert spaces*, vol. 408. Springer, 2011.
- [44] C. Lemaréchal and C. Sagastizábal, "Practical aspects of the Moreau-Yosida regularization: Theoretical preliminaries," *SIAM Journal on Optimization*, vol. 7, no. 2, pp. 367–385, 1997.



Dan Li received the B.E. degree in automation from the Zhejiang University, Hangzhou, China, in 2013, the M.Sc. degree in chemical engineering from Queen's University, Kingston, Canada, in 2016. He is currently a Ph.D. student at University of California, San Diego, CA, USA. His current research interests include data-driven systems and optimization, dynamical systems and control, optimization algorithms, applied computational methods, and stochastic systems. He received Outstanding Student Award from Zhejiang University in 2012, Graduate Student Award from Queen's University in 2014, and Fellowship Award from University of California, San Diego, in 2016.



Dariush Fooladivanda received the Ph.D. degree from the University of Waterloo, in 2014, and a B.S. degree from the Isfahan University of Technology, all in electrical engineering. He is currently a Postdoctoral Research Associate in the Department of Electrical Engineering and Computer Sciences at the University of California Berkeley. His research interests include theory and applications of control and optimization in large scale dynamical systems.



Sonia Martínez is a Professor at the Department of Mechanical and Aerospace Engineering at the University of California, San Diego. She received her Ph.D. degree in Engineering Mathematics from the Universidad Carlos III de Madrid, Spain, in May 2002. Following a year as a Visiting Assistant Professor of Applied Mathematics at the Technical University of Catalonia, Spain, she obtained a Postdoctoral Fulbright Fellowship and held appointments at the Coordinated Science Laboratory of the University of Illinois, Urbana-Champaign during 2004, and at the Center for Control, Dynamical systems and Computation (CCDC) of the University of California, Santa Barbara during 2005.

Her research interests include networked control systems, multi-agent systems, and nonlinear control theory with applications to robotics and cyber-physical systems. For her work on the control of underactuated mechanical systems she received the Best Student Paper award at the 2002 IEEE Conference on Decision and Control. She co-authored with Jorge Cortés and Francesco Bullo "Motion coordination with Distributed Information" for which they received the 2008 Control Systems Magazine Outstanding Paper Award. She is a Senior Editor of the IEEE Transactions on Control of Networked Systems and an IEEE Fellow.