

Constrained Model-Free Reinforcement Learning for Process Optimization

Elton Pan^a, Panagiotis Petsagkourakis^{b,*}, Max Mowbray^c, Dongda Zhang^c, Ehecatl Antonio del Rio-Chanona^{a,*}

^aCentre for Process Systems Engineering, Department of Chemical Engineering, Imperial College London, UK

^bCentre for Process Systems Engineering, Department of Chemical Engineering, University College London, UK

^cDepartment of Chemical Engineering and Analytical Science, University of Manchester, UK

Abstract

Reinforcement learning (RL) is a control approach that can handle nonlinear stochastic optimal control problems. However, despite the promise exhibited, RL has yet to see marked translation to industrial practice primarily due to its inability to satisfy state constraints. In this work we aim to address this challenge. We propose an “oracle”-assisted constrained Q-learning algorithm that guarantees the satisfaction of joint chance constraints with a high probability, which is crucial for safety critical tasks. To achieve this, constraint tightening (backoffs) are introduced and adjusted using Broyden’s method, hence making the backoffs self-tuned. This results in a methodology that can be imbued into RL algorithms to ensure constraint satisfaction. We analyze the performance of the proposed approach and compare against nonlinear model predictive control (NMPC). The favorable performance of this algorithm signifies a step towards the incorporation of RL into real world optimization and control of engineering systems, where constraints are essential.

Keywords: Machine Learning, Batch Optimization, Process Control, Q-learning, Dynamic Systems, Data-Driven Optimization

1. Introduction

The optimization of nonlinear stochastic processes poses a challenge for conventional control schemes given the requirement of an accurate process model and a method to simultaneously handle process stochasticity and satisfy state and safety constraints. Recent works have explored the application of model-free Reinforcement Learning (RL) methods for dynamic optimization of batch processes within the chemical and biochemical industries [37, 47]. Many of these demonstrate the capability of RL algorithms to learn a control law independently from a nominal process model, but negate proper satisfaction of state and safety constraints. In this work, we propose constrained Q-learning, a model-free algorithm to meet the operational and safety requirements of constraint satisfaction with high probability.

*Corresponding authors

Email addresses: p.petsagkourakis@ucl.ac.uk (Panagiotis Petsagkourakis), a.del-rio-chanona@imperial.ac.uk (Ehecatl Antonio del Rio-Chanona)

1.1. Model-Free Reinforcement Learning

RL encompasses a subfield of machine learning, which aims to learn an optimal policy for a system that can be described as a Markov decision process (MDP). Importantly, MDPs assume the Markov property, such that the future transitions (dynamics) of the process are only dependent upon the current state and control action, and not upon the process history [52].

Dynamic programming (DP) methods provide exact solution to MDPs under knowledge of the exact process dynamics [4]. A subset of RL algorithms are known as action-value (or Q-learning) methods. These methods are closely related to DP, but instead, learn an approximate parameterization of the optimal action-value function independently of explicit knowledge of the exact dynamics [50]. This is achieved by sampling the response of the underlying Markov process, eliminating the requirement for explicit assumption regarding the stochastic nature of the system. This is a particularly powerful concept in the domain of process control and optimization, given the inherent uncertainties and (slow) non-stationary dynamics often characteristic of process systems [49, 37].

The field of model-free reinforcement learning extends beyond action-value methods; two other approaches exist in the form of policy optimization [24, 51] and actor-critic [16, 26] methods. The relationships between the different sub-classes of RL algorithms are expressed by Fig. 1. The following analysis proceeds with reference to that Figure. Action-value methods explicitly learn a parameterization of the action-value function, whereas policy optimization methods implicitly learn the value space and instead learn and parameterize a policy directly [46]. Typically, policy optimization approaches deploy a Monte Carlo method to gain estimate of the value-function corresponding to the current policy parameters. This provides a search direction for further policy improvement. However, policy optimization methods tend to follow *on-policy* learning rules, which means that data collected under a given policy may only be used for one learning update before being discarded [50]. On-policy learning, combined with dependency on the use of a Monte Carlo method for evaluation of the search direction, provides significant sample complexity. Actor-critic methods partially reduce the associated sample complexity by explicitly learning a parameterization of the action-value function (the critic) as well as a policy (the actor) - removing dependency on Monte Carlo sampling [52]. In essence, the parameterization of the action-value function enables conceptualisation of actor-critic methods as those algorithms at the intersection of policy optimisation and action-value methods. Via this analysis, the apparent benefit of action-value methods is of reduced sample-complexity. This is further improved when experience replay is integrated [41], enabling saving of data observed in previous simulations for future learning updates. As such, this directs focus in the analysis provided by the following section, where review of action-value function approximation in chemical engineering is presented.

1.2. Action-Value Function Approximation in Chemical Engineering

In recent years, there has been a growing interest in the development of RL controllers in the domain of chemical and biochemical processing and this is reflected by the rapid growth of literature in the area. The interest is primarily generated by two attractive properties - reduced sample complexity and the ability to integrate experience replay [7]. There are two classes of RL algorithms which leverage value function approximation: actor-critic, and action-value methods. The use of the two approaches is reviewed in the following discussion.

Actor-critic methods have been applied widely, given their explicit parameterisation of a policy (actor) and the associated sample efficiency induced via parameterisation of a value function (critic). In [15] an actor-critic method for pH control of wastewater from an industrial

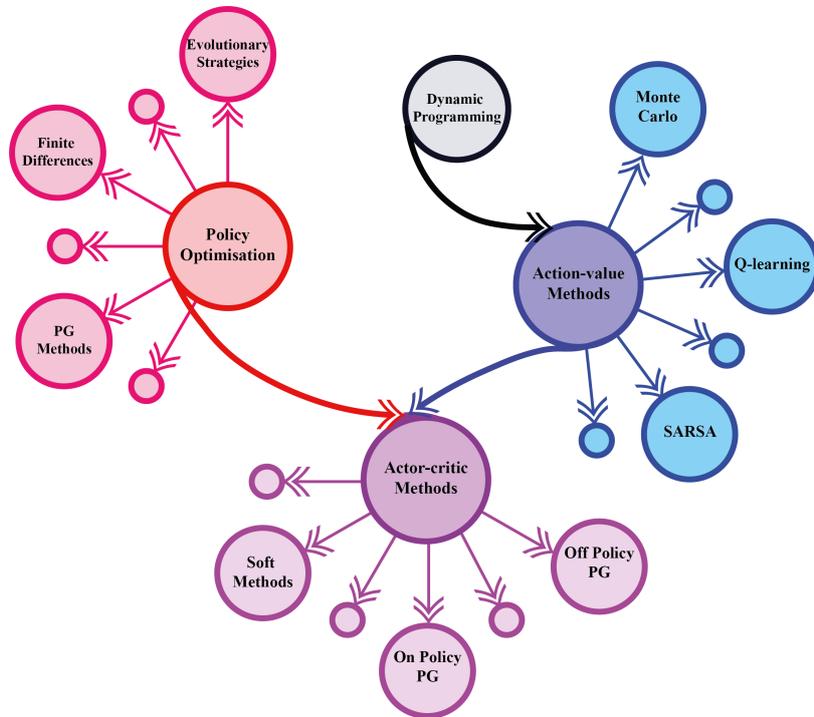


Figure 1: The landscape of model free reinforcement learning (RL). Model free RL can be broadly constituted by policy optimization and action-value methods with the intersection of the two characterised by actor-critic methods. The figure does not exhaustively detail the different algorithms, but rather broadly describes those approaches that are dominant within each respective class. The use of bubbles with no description denotes those algorithms not detailed for the sake of conciseness.

electroplating process is proposed and benchmarked against a proportional-integral-derivative (PID) controller. Similarly, in [21], the authors present an actor-critic based RL framework for optimal PID controller tuning, including a mechanism for *antiwindup*. In [59], the authors propose a framework to augment the actor-critic algorithm (Deep deterministic policy gradient, DDPG) with nonlinear model predictive control (NMPC). The framework is then demonstrated on two case studies, both of which highlight the framework’s ability to provide marked improvements in the efficacy of offline policy training, relative to a vanilla implementation of DDPG. In [49], the authors identify the potential benefits of a data-driven RL controller with respect to ease of system re-identification and demonstrate the application of an actor-critic algorithm to a number of case studies including control of a high purity distillation column.

For action-value methods to be deployed in continuous control spaces, they typically require augmentation with a further optimization method for determining the optimal control [43], and currently, very efficient optimization algorithms exist [56, 33] to facilitate this. Action-value methods are particularly attractive given they generally encompass the most sample efficient of algorithms in the RL toolbox. In [44], a Q-learning approach was applied for nanostructure and nanosurface design. Similarly, [60] proposes a Q-learning method for the purpose of molecular design. In [19], the authors integrate a Q-learning method with Aspen Plus for control of a downstream separation process, demonstrating improved performance relative to open loop

operation. In [22] an action-value method is proposed in the context of process scheduling and in [23] the authors demonstrate the application of action-value methods for the control of stochastic, nonlinear processes. The above have shown good success in chemical process optimization and control providing basis for further use of action-value methods and development of 'saving focused' learning strategies. However, for RL methods to be deployed in many instances of process engineering, they must satisfy constraints (with high probability in the stochastic setting). This is the challenge that the current work addresses.

1.3. Related Work

As highlighted in [25], there exists relative inertia in application of RL to industrial control problems. Specifically, in the chemical and biochemical process industries, the development of RL methods to guarantee safe process operation and constraint satisfaction would enhance prospective deployment of RL-based systems in the scope of optimization and control. The literature documents a number of approaches to safe constraint satisfaction that typically augment a pure model-free RL-based controller, with a separate system that has basis in direct optimal control. Such augmented systems are broadly constituted by barrier function [9, 53, 8] and safety filter methods [12, 55]. The deployment of these approaches dictates a nominal description of the process dynamics, method to handle process stochasticity and often impose nontrivial policy or value learning rules. The aspect of model dependency particularly dampens the initial attraction of an RL approach within the context of process control.

Other works have explored the development of methods for safe constraint satisfaction by leveraging the value framework provided by MDPs, preserving performance independent of a process model. These methods either add penalty to the original reward function (objective) for constraint violation [23, 54] or augment the original MDP to take the form of a constrained MDP (CMDP) [3]. If approached crudely, the former approach introduces a number of hyperparameters, which are typically chosen on the basis of heuristics and have bearing on policy optimality. This is also discussed in [2, 14]. The latter approach is underpinned by the learning of surrogate cost functions for each individual constraint combined with appropriate adaptation of the policy [2] or value learning rule. Both approaches ensure constraint satisfaction only in expectation [45], which is insufficient for control and optimization of (bio)chemical processes. As most engineering systems are safety critical, satisfaction of constraints with high probability is a necessity. In view of this problem, [42] present an approach which combines the CMDP concept with consideration of worst case realisations of process uncertainty, in an attempt to robustify the framework. However, the approach is only demonstrated empirically in a discrete domain (almost all industrial problems are continuous). In the same rationale, a Lyapunov-based approach is proposed in [10], where a Lyapunov function is found and the unconstrained policy is projected to a safety layer allowing the satisfaction of constraints in expectation. However this is not the optimal trade-off, as the closest action is not necessarily optimal. Additionally, the satisfaction of constraints in expectation is not sufficient for the most real-world applications. A similar approach is also presented in [18], where a Lyapunov function is identified, but for the purpose of safe exploration. This is not necessarily directly applicable to the process industries, where online exploratory control decisions pose high operational risk and the fact that operational constraints are not explicitly considered. More recently, work has been proposed in [34], based on the augmented Lagrangian. However, the method negates the inclusion of state and control dependency within the penalty term, instead appending the penalty in a crude form of credit assignment, therefore implicitly neglecting an aspect of the RL problem. It is likely that this approach leads to conservative control policies. The method presented in [36], provides a mechanism for the state and control dependency desired.

The work presented here, leverages similar concepts but in the context of action-value function learning, improving the data-efficiency of the algorithm alongside other advantages discussed herein.

1.4. Contribution

To our knowledge, no approach has been proposed which achieves constraint satisfaction with high probability for action-value methods. In this work, we propose an action-value method, which guarantees constraint satisfaction with high probability. Here, we first learn an unconstrained actor and surrogate constraint action-value functions. We then subsequently construct a constrained action-value function as a superimposition of the unconstrained actor with the surrogate constraints (in the form of a constrained optimization problem). The constrained actor is iteratively tuned, as learning proceeds, via localised backoffs [6, 28, 40] to penalize constraint violation. Conceptually, backoffs provide a policy variant shaping mechanism to ensure high probability satisfaction [32]. Tuning comprises a Monte Carlo method to estimate the probability of constraint violation under the policy combined with Broyden’s root finding method. Specifically, Broyden’s method is used to update the backoff values. Importantly, the dimensionality of the tuning problem is equivalent to the number of operational constraints imposed upon the problem providing a scale-able algorithm. Given the constrained actor action-value function and the fast inference associated with neural networks, efficient optimization strategies may be deployed for determination of the optimal control. Further, the algorithm proposed incorporates the use of an experience replay memory store, promoting interpretation of the method as a saving algorithm [7]. This is particularly important in the scope of continual learning and improvement of the policy once deployed to the real process [41]. The work is arranged as follows; the problem description is formalized in section 2, the methodology proposed in section 3 and demonstrated empirically in section 4 via two benchmark case studies.

2. Problem Statement

2.1. Reinforcement Learning in Process Engineering

Using reinforcement learning directly on an industrial plant to construct an accurate controller would require prohibitive amounts of data due to the random initialization of the policy. As such the initial training phase would require a model of the process dynamics, which could either be a data-driven or based on first principles. Additionally, the policy may be warm-started by techniques such as behavioral cloning or apprentice learning [31]. This would also ensure that safety violations do not occur in the real plant. The simplified workflow shown in Fig. 2 starts with either a randomly initialized policy or a policy that is warm-started by an existing controller and apprenticeship learning [1]. Preliminary training is performed using closed-loop simulations from the offline process model (notice that a stochastic model can be used). Here, the resulting control policy is a good approximation of the optimal policy of the real plant, which is subsequently deployed in the real plant for further training online. Importantly, system stochasticity is accounted for and the controller will continue to adapt and learn to better control and optimize the process, hence addressing plant-model mismatch [13, 57, 35]. It is thought a number of processes with uncertain, nonlinear dynamics would particularly benefit from the maturation of RL-based controllers [46]. These include biochemical reaction systems [37], multi-phase separations as well as those with complex rheological behaviours. However, it should be emphasised that essentially any process common to industry could benefit.

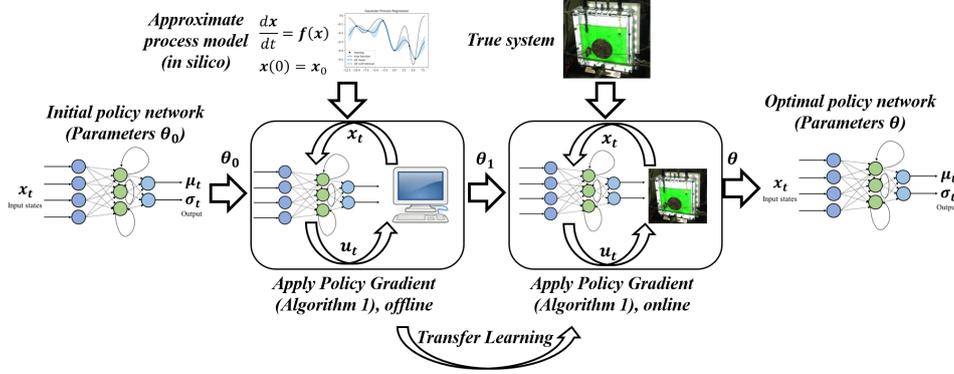


Figure 2: Schematic representation of RL for chemical process optimization (Adapted from [37]).

2.2. Stochastic Optimal Control Problem

We assume that the stochastic dynamic system in question follows a Markov process and transitions are given by

$$\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t), \quad (1)$$

where $p(\mathbf{x}_{t+1})$ is the probability density function of future state \mathbf{x}_{t+1} given a current state $\mathbf{x}_t \in \mathbb{R}^{n_x}$ and control $\mathbf{u}_t \in \mathbb{R}^{n_u}$ at discrete time t , and the initial state is given by $\mathbf{x}_0 \sim p_{\mathbf{x}_0}(\cdot)$. Without loss of generality we can write Eq. (1) as:

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t, \mathbf{d}_t, \mathbf{p}), \quad (2)$$

where $\mathbf{p} \in \mathbb{R}^{n_p}$ are the uncertain parameters of the system and $\mathbf{d}_t \in \mathbb{R}^{n_d}$ are the stochastic disturbances. In this work, the goal is to maximize a predefined economic metric via a policy π (a function that outputs a control \mathbf{u}_t given state \mathbf{x}_t) subject to constraints. Consequently, this problem can be framed as an optimal control problem:

$$\mathcal{P}(\pi(\cdot)) := \begin{cases} \max_{\pi(\cdot)} \mathbb{E} \{ J(\mathbf{x}_0, \dots, \mathbf{x}_{t_f}, \mathbf{u}_0, \dots, \mathbf{u}_{t_f}) \} \\ \text{s.t.} \\ \mathbf{x}_0 \sim p_{\mathbf{x}_0}(\mathbf{x}_0) \\ \mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \\ \mathbf{u}_t = \pi(\mathbf{x}_t) \\ \mathbf{u}_t \in \mathbb{U} \\ \mathbb{P} \left(\bigcap_{t=0}^{t_f} \{ \mathbf{x}_t \in \mathbb{X}_t \} \right) \geq 1 - \omega \\ \forall t \in \{0, \dots, t_f\} \end{cases} \quad (3)$$

where J is the objective function, \mathbb{U} is the set of hard constraints for the controls and \mathbb{X}_t denotes constraints for states that must be satisfied. In other words,

$$\mathbb{X}_t = \{ \mathbf{x}_t \in \mathbb{R}^{n_x} \mid g_{j,t}(\mathbf{x}_t) \leq 0, j = 1, \dots, n_g \}, \quad (4)$$

with n_g being the total number of constraints to be satisfied, and $g_{j,t}$ being the j th constraint to be satisfied at time t . Joint constraint satisfaction must occur with high probability $1 - \omega$ where $\omega \in [0, 1]$. Herein, we present a Q-learning algorithm that allows to obtain the optimal policy which satisfies joint chance constraints.

2.3. Q-learning

Q-learning is a model-free reinforcement learning algorithm which trains an agent to behave optimally in a Markov process [58]. The agent performs actions to maximize the expected sum of all future discounted rewards given an objective function $J(\cdot)$, which can be defined as

$$J = \sum_{t=0}^{t_f} \gamma^t R_t(\mathbf{x}_t, \mathbf{u}_t), \quad (5)$$

where $\gamma \in [0, 1]$ is the discount factor and R_t represents the reward at time t given values \mathbf{x}_t and \mathbf{u}_t . In the context of process control, the agent is the controller, which uses a policy $\pi(\cdot)$ to maximize the expected future reward through a feedback loop. Interaction between the agent (or controller) and system (in this case, a simulator) at each sampling time returns a value for the reward R that represents the performance of the policy.

In Q-learning, for a policy π an action-value function can be defined as

$$Q^\pi(\mathbf{x}_t, \mathbf{u}_t) = R_{t+1} + \gamma \max_{\mathbf{u}_{t+1}} Q^\pi(\mathbf{x}_{t+1}, \mathbf{u}_{t+1}) \quad (6)$$

with $Q^\pi(\mathbf{x}_{t+1}, \mathbf{u}_{t+1})$ being the expected sum of all the future rewards the agent receives in the resultant state \mathbf{x}_{t+1} . Importantly, the Q -value is the expected discounted reward for a given state and action, and therefore the optimal policy π^* can be found using iterative updates with the Bellman equation (Eq. (6)). Upon convergence, the optimal Q -value Q^* is defined as:

$$Q^*(\mathbf{x}_t, \mathbf{u}_t) = \mathbb{E}_{\mathbf{x}_{t+1} \sim p} \left[R_{t+1} + \gamma \max_{\mathbf{u}_{t+1}} Q^*(\mathbf{x}_{t+1}, \mathbf{u}_{t+1}) \mid \mathbf{x}_t, \mathbf{u}_t \right] \quad (7)$$

$Q(\mathbf{x}_t, \mathbf{u}_t)$ can be represented by function approximators such as neural networks, Gaussian process [11], tree-based regressors [38] amongst others. In this work, the Q -function is approximated with a deep Q-network (DQN) Q_θ parameterized by weights θ [29]. Here, the inputs specifically include the state \mathbf{x}_t , the corresponding time step t and control \mathbf{u}_t . The DQN is trained with the use of a replay buffer that addresses the issue of correlated sequential samples [27]. Huber loss is used as the error function [17, 30]. Initial exploration is encouraged using an ϵ -greedy policy starting with high ϵ values, which is decayed over the course of training to ensure eventual exploitation and convergence to the optimal policy.

3. Methodology

Our proposed approach can be found in Algorithm 1. Firstly, the Q-network (Q_θ) and the constraint networks ($G_{j,\phi}$) and their respective buffers were initialized (Steps 1 and 2). Here, each constraint network $G_{j,\phi}$ corresponds to a neural network that learns the oracle values for the j th constraint. Subsequently, MC roll-outs of the process (Step A) generated Q -values (Step b) and oracle values $\hat{g}_{j,t}$ (Step c), which were stored in the respective buffers \mathcal{D} and \mathcal{G}_j . The oracle $\hat{g}_{j,t}$ (Step c) is defined as the maximum level of violation to occur in all current and future time steps in the process realization as shown in Eq. (9), and will be further discussed in Section 3.1. In the MC roll-outs (Step a), actions are chosen as shown in Fig. 3 (a), by optimizing a fitness function (sub-problem in Algorithm 1) that comprises predictions by these neural networks to ensure that chosen actions satisfy constraints while maximizing reward. This can be framed as the following constrained optimization problem:

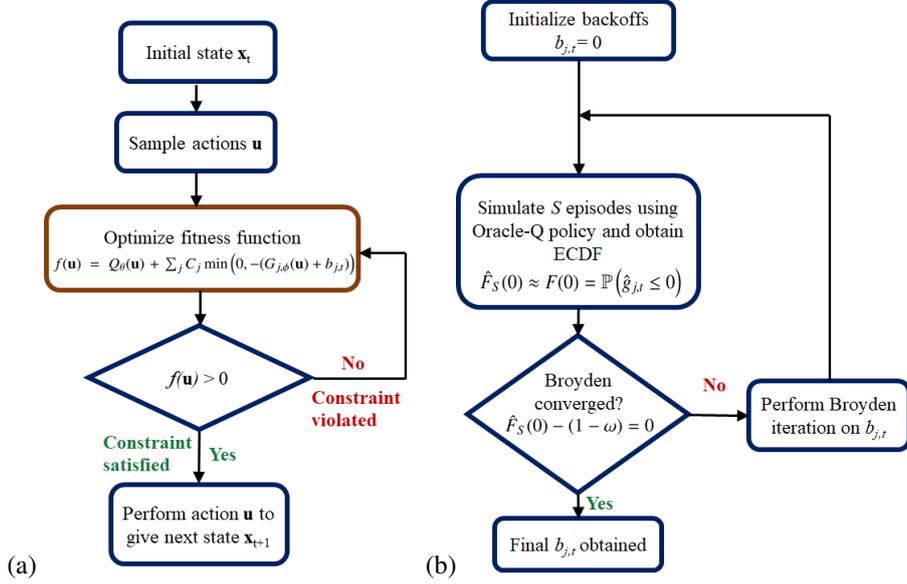


Figure 3: Flow chart for (a) choosing action \mathbf{u} to satisfy constraints while maximizing reward (b) adjustment of backoffs.

$$\begin{aligned}
 & \max_{\mathbf{u}} Q_\theta(\mathbf{x}, \mathbf{u}) \\
 & \text{s.t.} \\
 & G_{j,\phi}(\mathbf{x}, \mathbf{u}) + b_{j,t} \leq 0, j = 1, \dots, n_g
 \end{aligned} \tag{8}$$

where $b_{j,t}$ are the backoffs which tighten the former feasible set \mathbb{X}_t stated in Eq. (4). The nominal case with backoff $b_{j,t} = 0$ corresponds to the absence of tightening. This is further explained in Section 3.2 Neural networks (Q_θ and $G_{j,\phi}$) have been trained by random sampling from the replay/constraint buffers (Steps B and C in Algorithm 1), followed by performing a gradient optimization to learn weights θ and ϕ (Step D). After training using Algorithm 1, the optimal Q-network Q_θ^* and constraint networks $G_{j,\phi}$ are obtained.

Subsequently, we perform constraint tightening using backoffs as described in Fig. 3 (b) such that constraint satisfaction occurs with desired probability of $1 - \omega$ as in Eq. (3), which will be

discussed in Section 3.2.

Algorithm 1: Oracle-assisted constrained Q-learning

1. Initialize replay buffer \mathcal{D} of size $s_{\mathcal{D}}$ and constraint buffers \mathcal{G}_j of size $s_{\mathcal{G}}$, $j = 1, \dots, n_g$
2. Initialize Q-network Q_θ and constraint networks $G_{j,\phi}$ with random weights θ and ϕ ,
 $j = 1, \dots, n_g$
3. Initialize ϵ and backoffs $b_{j,t}$
4. **for** training iteration = $1, \dots, M$ **do**
 - A. **for** episode = $1, \dots, N$ **do**
 - Initialize state $\mathbf{x}_0 \sim p_{\mathbf{x}_0}(\mathbf{x}_0)$ and episode \mathcal{E}
 - a. **for** $t = 0, \dots, t_f$ **do**
 - With probability ϵ select random control \mathbf{u}_t
 - otherwise select $\mathbf{u}_t = \max_{\mathbf{u}} Q_\theta(\mathbf{x}_t, \mathbf{u}_t) \mid G_{j,\phi}(\mathbf{x}_t, \mathbf{u}_t) + b_{j,t} \leq 0, j = 1, \dots, n_g$
(Sub-problem^a)
 - Execute control \mathbf{u}_t and observe reward R_t and new state \mathbf{x}_{t+1}
 - Store transition $(\mathbf{x}_t, \mathbf{u}_t, R_t, \mathbf{x}_{t+1})$ in \mathcal{E}
 - b. Extract Q-values from \mathcal{E} and store datapoint $(\mathbf{x}_t, \mathbf{u}_t, Q_t)$ in \mathcal{D}
 - c. Extract oracle-constraint values from \mathcal{E} using:
 $\hat{g}_{j,t} = \max_{t' \geq t} [g_{j,t'}], j = 1, \dots, n_g$
 - d. Store datapoint $(\mathbf{x}_t, \mathbf{u}_t, \hat{g}_{j,t})$ in $\mathcal{G}_j, j = 1, \dots, n_g$
 - B. Sample random minibatch of datapoints of size G $(\mathbf{x}_t, \mathbf{u}_t, Q_t)$ from \mathcal{D}
 - C. Sample random minibatch of datapoints of size H_j $(\mathbf{x}_t, \mathbf{u}_t, \hat{g}_{j,t})$ from \mathcal{G}_j
 - D. Perform a gradient-descent type step (e.g. Adam optimizer^b) on Q_θ and $G_{j,\phi}$ and update weights θ and ϕ
 - E. Decay ϵ using $\epsilon = D_1 \epsilon$
 - F. Decay backoffs using $b_{j,t} = D_2 b_{j,t}$

Output: Optimal Q-network Q_θ^* and constraint networks $G_{j,\phi}, j = 1, \dots, n_g$

^aSub-problem: A derivative-free algorithm (e.g. evolutionary or Bayesian optimization) is used to optimize the nonconvex constrained Q-function using fitness function $f(\mathbf{u}) = Q_\theta(\mathbf{u}) + \sum_j C_j \min(0, -(G_{j,\phi}(\mathbf{u}) + b_{j,t}))$ where $g_{j,t}$ is the j th constraint violation at time t , and $b_{j,t}$ is the corresponding backoff. C_j are large values to ensure large negative fitness values for controls that lead to constraint violation.

^bAny other partial, or full optimization step can be used here.

3.1. Oracle-assisted Constrained Q-learning

Q-learning, when unconstrained, may offer little practical utility in process optimization due to unbounded exploration by the RL agent. For instance, an unconstrained policy may often result in a thermal runaway leading to a safety hazard in the process. As such, herein constraints $g_{j,t}$ are incorporated through the use of an oracle $\hat{g}_{j,t}$ which is formulated as

$$\hat{g}_{j,t} = \max_{t' \geq t} [g_{j,t'}] \quad (9)$$

with $g_{j,t}$ being the j th constraint to be satisfied at time t , and the oracle $\hat{g}_{j,t}$ is determined by the maximum level of violation to occur in all current and future time steps t' in the process realization.

The intuition behind this framework is as follows: Imagine a car (agent) accelerating towards the wall with the goal of minimizing the time it takes to reach some distance from the wall (objective) without actually crashing into the wall (constraint). Accelerating the car without

foresight causes it to go so fast that it cannot brake and stop in time, causing it to crash into the wall (constraint violated). As such, there is a need for foresight to ensure constraint satisfaction.

Effectively, the framework shown in Eq. (9) is akin to an oracle (or fortune-teller peeking into a crystal ball) advising the agent on the *worst* (or maximum) violation that a specific action can cause in the future given the current state. These values are easily obtained using Monte-Carlo simulations of the system. Analogous to the way a Q-function gives the sum of all future rewards, the oracle provides the worst violation in all future states if a certain action is taken by the agent, hence imbuing in the agent a sense of foresight to avoid future constraint violation.

Similar to the Q-function, constraint values are represented by neural networks $G_{j,\phi}$ with state and action as input features. However, the subtle difference between the two is that the state representation of the input for $G_{j,\phi}$ involves time-to-termination $t_f - t$ instead of time t for the case of batch processes.

3.2. Constraint Tightening

To satisfy the constraints with high probability, constraints are tightened with backoffs [6, 39] $b_{j,t}$ as:

$$\bar{\mathbb{X}}_t = \left\{ \mathbf{x}_t \in \mathbb{R}^{n_x} \mid g_{j,t}(\mathbf{x}_t) + b_{j,t} \leq 0, j = 1, \dots, n_g \right\} \quad (10)$$

where $b_{j,t}$ are the backoffs which tighten the former feasible set \mathbb{X}_t stated in Eq. (4). The result of this would be the reduction of the perceived feasible space by the agent, which consequently allows for the satisfaction of constraints. Notice that the value of the backoffs necessarily imply a trade-off: large backoff values ensure constraint satisfaction, but renders the policy over-conservative hence sacrificing performance. Conversely, smaller backoff values afford solutions with higher rewards, but may not guarantee constraint satisfaction. Therefore, the values of $b_{j,t}$ are the minimum value needed to guarantee satisfaction of constraints.

To determine the desired backoffs, the cumulative distribution function (CDF) F of the oracle $\hat{g}_{j,t}$ is approximated using sample average approximation (SAA) with S Monte Carlo (MC) simulations to give its empirical cumulative distribution function (ECDF) \hat{F}_S where

$$\hat{F}_S(0) \approx F(0) = \mathbb{P}(\hat{g}_{j,t} \leq 0) \quad (11)$$

hence $\hat{F}_S(0)$ is the approximate probability for a trajectory to satisfy a constraint. We can therefore pose a root-finding problem such that we adjust the backoffs $b_{j,t}$ to find:

$$\hat{F}_S(0) - (1 - \omega) = 0 \quad (12)$$

We solve Eq. (12) via the quasi-Newton Broyden's method [20] given its fast convergence near optimal solutions. Where ω is a tunable parameter depending on the case study, such that constraint satisfaction occurs with high probability $1 - \omega$ as shown in Eq. (3). Alternatively, the empirical lower bound of the ECDF can be forced to be $1 - \omega$, and guarantee with confidence $1 - \epsilon$ that $\mathbb{P}(\hat{g}_{j,t} \leq 0) \geq 1 - \omega$. Note that this constraint tightening takes place directly during the construction of $G_{j,\phi}$, $j = 1, \dots, n_g$

4. Case Studies

4.1. Case Study 1

This case study pertains to the photoproduction of phycocyanin synthesized by cyanobacterium *Arthrospira platensis*. Phycocyanin is a high-value bioproduct, and serves its biological role by

increasing the photosynthetic efficiency of cyanobacteria and red algae [6]. In addition, it is used as a natural colorant to substitute toxic synthetic pigments in cosmetic and food manufacturing. Moreover, it possesses antioxidant, and anti-inflammatory properties.

The dynamic system comprises a system of ODEs from [6] that describes the evolution of concentration (c) of biomass (x), nitrate (N) and product (q) under parametric uncertainty. The model is based on Monod kinetics, which describes the growth of microorganism in nutrient-sufficient cultures, where intracellular nutrient concentration is kept constant because of rapid replenishment. Here, a fixed volume fed-batch is assumed. The controls are light intensity ($u_1 = I$) and inflow rate ($u_2 = F_N$). The mass balance equations are as follows:

$$\begin{aligned}\frac{dc_x}{dt} &= u_m \frac{I}{I + k_s + I^2/k_i} \frac{c_x c_N}{c_N + K_N} - u_d c_x \\ \frac{dc_N}{dt} &= -Y_{N/X} \frac{u_m I}{I + k_s + I^2/k_i} \frac{c_x c_N}{c_N + K_N} + F_N \\ \frac{dc_q}{dt} &= \frac{k_m I}{I + k_{sq} + I^2/k_{iq}} c_x - \frac{k_d c_q}{c_N + K_{N_q}}\end{aligned}\quad (13)$$

This case study and parameter values are adopted from [6]. Uncertainty in the system is two-fold: First, the initial concentration adopts a Gaussian distribution, where $[c_{x,0}, c_{N,0}] \sim \mathcal{N}([1.0, 150.0], \text{diag}(10^{-3}, 22.5))$ and $c_q(0) = 0$. Second, parametric uncertainty is assumed to be: $\frac{k_s}{(\mu\text{mol}/\text{m}^2/\text{s})} \sim \mathcal{N}(178.9, \sigma_{k_s}^2)$, $\frac{k_i}{(\text{mg}/\text{L})} \sim \mathcal{N}(447.1, \sigma_{k_i}^2)$, $\frac{k_N}{(\mu\text{mol}/\text{m}^2/\text{s})} \sim \mathcal{N}(393.1, \sigma_{k_N}^2)$ where the variance $\sigma_i^2 = 10\%$ of its corresponding mean value. This type of uncertainty is common in engineering settings, as the parameters are experimentally determined, and therefore subject to confidence intervals after being calculated through nonlinear regression techniques. The objective function is to maximize the product concentration (c_q) at the end of the batch, hence the reward is defined as:

$$R_{t_f} = c_{q,t_f} \quad (14)$$

where t_f is the terminal time step. The two path constraints are as follows: Nitrate concentration (c_N) is to remain below 800 mg/L, and the ratio of bioproduct concentration (c_q) to biomass concentration (c_x) cannot exceed 11.0 mg/g for high density biomass cultivation. These constraints can be formulated as:

$$\begin{aligned}g_{1,t} &= c_N - 800 \leq 0 \quad \forall t \in \{0, \dots, t_f\} \\ g_{2,t} &= c_q - 0.011c_x \leq 0 \quad \forall t \in \{0, \dots, t_f\}\end{aligned}\quad (15)$$

The control inputs are subject to hard constraints to be in the interval $0 \leq F_N \leq 40$ and $120 \leq I \leq 400$. The time horizon was set to 12 with an overall batch time of 240 h, and hence giving a sampling time of 20 h. The Q-network Q_θ consists of 2 fully connected hidden layers, each consisting of 200 neurons with a leaky rectified linear unit (LeakyReLU) as activation function. The parameters used in Algorithm 1 for training the agent are: $\epsilon = 0.99$, $b_{1,t} = -500$, $b_{2,t} = -0.05$, $s_D = 3000$, $s_G = 30000$, $M = 2000$, $N = 100$, $G = 100$, $H_1 = 500$, $H_2 = 1000$, $D_1 = 0.99$ and $D_2 = 0.995$. Upon completion of training, validation was conducted via the trained policy with respect to the Q function. The policy is optimized through an evolutionary strategy [48] given its nonconvex nature, as discussed in Algorithm 1 and Fig. 3 (a).

After completion of training using Algorithm 1, the backoffs are adjusted to satisfy Eq. (12) with $S = 1000$ and $\omega = 0.01$, with backoffs at all time-steps t being constant. The constraint satisfaction is shown in Fig. 4, where the shaded areas represent the 99th to 1st percentiles. Here, we elucidate the importance of applying backoffs to the policy: As shown in Fig. 4 (a), even

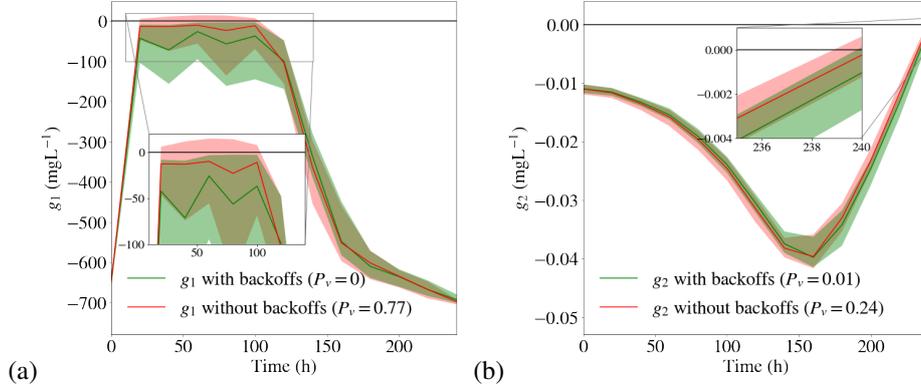


Figure 4: Case Study 1: Constraints $g_{1,t}$ (a) and $g_{2,t}$ (b) when backoffs are applied (green), and when they are absent (red) with probabilities of violation P_v within the parentheses. Inset: Zoomed-in region where violation of constraints occur. Solid lines represent the expected values. Shaded areas represent the 99th to 1st percentiles.

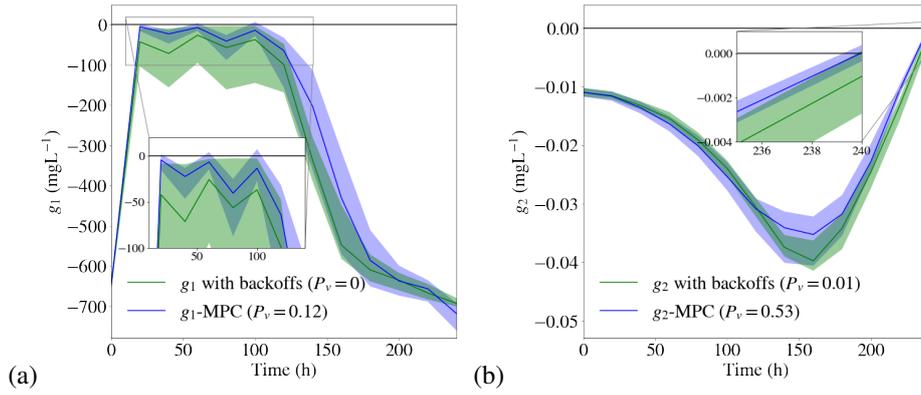


Figure 5: Case Study 1: Constraints $g_{1,t}$ (a) and $g_{2,t}$ (b) when backoffs are applied (green), and for NMPC (blue) with probabilities of violation P_v within the parentheses. Inset: Zoomed-in region where violation of constraints occur. Solid lines represent the expected values. Shaded areas represent the 99th to 1st percentiles.

Table 1: Case Study 1: Comparison of probabilities of constraint violation P_v and objective values of different algorithms

Algorithm	Violation probability P_v	Objective (c_{q,t_f})
Oracle Q-learning with backoffs	0.01	0.166
Oracle Q-learning without backoffs	0.82	0.169
NMPC	0.53	0.168

though it may seem at face value that $g_{1,t}$ values for both methods are similar, the zoomed-in region (in the inset) clearly shows that oracle Q-learning without backoffs (red) results in a high probability of constraint violation ($P_v = 0.77$), with parts of the red shaded regions exceeding zero. The violation probabilities P_v in Fig. 4 and 5 correspond to the fraction of 400 MC trajectories that violate a certain constraint. Gratifying, when backoffs are applied (green) in Fig. 4 (a), all

constraints are satisfied ($P_v = 0$), as shown by all the green shaded regions staying below zero.

In the same vein, in Fig. 4 (b), applying backoffs resulted in a drastic reduction of constraint violation from $P_v = 0.24$ to 0.01. This is expected since the backoffs are adjusted using $\omega = 0.01$ in Eq. (3). The objective value, represented by the final concentration of product c_q , are 0.166 and 0.169 for oracle Q-learning with and without backoffs, respectively. Consequently, this indicates that a small compromise in objective value can result in high probability of constraint satisfaction, where violation probability is reduced from 0.82 to 0.01 (in boldface) upon applying backoffs as shown in Table 1.

In addition, the performance of the oracle Q-learning algorithm with backoffs has been compared with that of nonlinear NMPC using the nominal parameters of the model, which is one of the main process control techniques used in chemical process optimization and hence serves as an important benchmark.

Although NMPC achieves a slightly higher objective value (Table 1), it fares poorly in terms of constraint satisfaction as shown in blue Fig. 5 (a) and (b) where probabilities of violation are 12 and 53% for g_1 and g_2 , respectively. This is unsurprising, since NMPC is only able to satisfy constraints in *expectation*, which means that in a stochastic system, loosely speaking, violation occurs 50% of the time. On the other hand, oracle Q-learning with backoffs violated a constraint only 1 % of the time (boldface in Table 1). Therefore, it is clear that this algorithm offers a more effective means of handling constraints compared to NMPC.

It is also noteworthy to contrast the proposed method with policy-based methods. Previous work in the group [36] involves policy optimization, which has its own benefits, such as avoiding an online optimization routine, this however, comes with the drawback of being unable to handle constraints naturally. On the other hand, value-based methods, as proposed in this paper, have an inner optimization loop, which allows them to handle constraints easily with mature algorithms from the numerical optimization community. Therefore, value-based methods can be better tailored to satisfy constraints.

4.2. Case Study 2

The second case study involves a challenging semi-batch reactor adopted from [5], with the following chemical reactions in the reactor catalyzed by H_2SO_4 :



Here, the reactions are first-order. Reactions (1) and (2) are exothermic and endothermic, respectively. The temperature is controlled by a cooling jacket. The controls are the flowrate of reactant A entering the reactor and the temperature of the cooling jacket T_0 . Therefore, the state is represented by the concentrations of A, B, and C in mol/L (c_A, c_B, c_C), reactor temperature in K (T), and the reactor volume in L (Vol). The model of the physical system can be found in [5].

The objective function is to maximize the amount of product ($c_C \cdot Vol$) at the end of the batch. Two path constraints exist. Firstly, the reactor temperature needs to be below 420 K due to safety reasons and secondly, the reactor volume is required to be below the maximum reactor capacity of 800 L and therefore:

$$\begin{aligned} g_{1,t} &= T - 420 \leq 0 \quad \forall t \in \{0, \dots, t_f\} \\ g_{2,t} &= Vol - 800 \leq 0 \quad \forall t \in \{0, \dots, t_f\} \end{aligned} \quad (17)$$

The time horizon is fixed to 10 with an overall batch time of 4 h, therefore the sampling time is 0.4 h. Parametric uncertainty is set as: $\theta_1 \sim \mathcal{N}(4, 0.1)$, $A_2 \sim \mathcal{N}(0.08, 1.6 \times 10^{-4})$, $\theta_4 \sim \mathcal{N}(100, 5)$. The

initial concentrations of A, B and C are set to zero. The initial reactor temperature and volume are 290 K and 100 L, respectively.

In this case study, due to its more challenging nature in terms of constraint satisfaction compared to the first case study, the backoffs have been adjusted to satisfy Eq. (12) using $\omega = 0.1$ in Eq. (3). We observe that backoffs again proved to be necessary to ensure high probability of constraint satisfaction. From the inset of Fig. 6 (a), we can see that without backoffs the policy violates g_1 41% of the time, and this probability is reduced to 9% when backoffs are applied. The same applies for g_2 in Fig. 6 (b) where P_v is completely eliminated from 3% to 0% using backoffs.

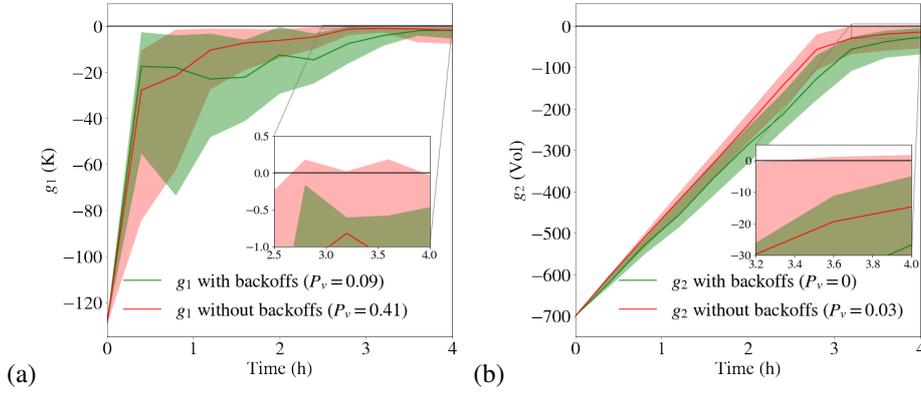


Figure 6: Case Study 2: Constraints $g_{1,t}$ (a) and $g_{2,t}$ (b) when backoffs are applied (green), and when they are absent (red) with probabilities of violation P_v within the parentheses. Inset: Zoomed-in region where violation of constraints occur. Solid lines represent the expected values. Shaded areas represent the 95th-5th percentiles for (a) and 99th-1st percentiles for (b).

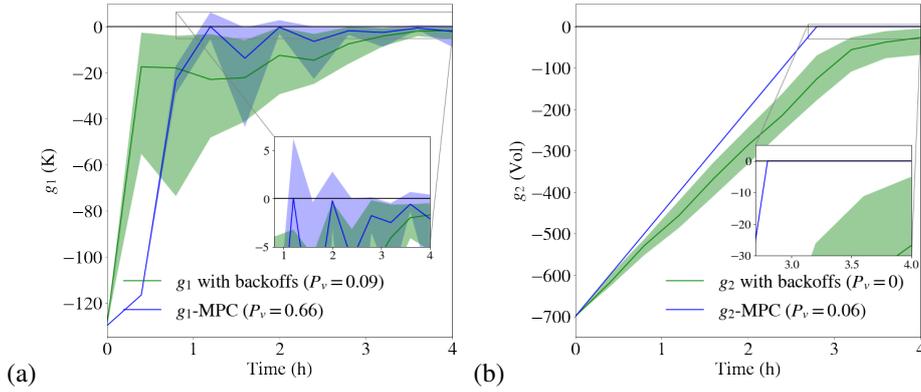


Figure 7: Case Study 2: Constraints $g_{1,t}$ (a) and $g_{2,t}$ (b) when backoffs are applied (green), and for NMPC (blue) with probabilities of violation P_v within the parentheses. Inset: Zoomed-in region where violation of constraints occur. Solid lines represent the expected values. Shaded areas represent the 95th-5th percentiles for (a) and 99th-1st percentiles for (b).

To compare the performance of NMPC with oracle Q-learning with backoffs in the context of this case study. The average runtime to solve for an individual control for the oracle Q-learning

Table 2: Case Study 2: Comparison of probabilities of constraint violation P_v and objective values of different algorithms

Algorithm	Violation probability (P_v)	Objective ($c_{C,t_f} \cdot Vol_{t_f}$)
Oracle Q-learning with backoffs	0.09	532
Oracle Q-learning without backoffs	0.44	680
NMPC	0.66	714

approach (0.03 ± 0.01 s) is faster than NMPC (0.4 ± 0.1 s). This can be attributed to the fact that MPC relies on a model-based optimization, whereas our approach has already learnt this offline. Moreover, it is noteworthy that chemical systems are rarely deterministic in nature, hence limiting the applicability of NMPC. In a stochastic system, NMPC often struggles in terms of constraint handling. This can be clearly seen in Fig. 7 (a), where the MPC trajectories only satisfy g_1 in expectation (blue line), hence resulting in high levels of violations (66%). Intriguingly, for g_2 the NMPC trajectory in Fig. 7 displayed little variation, resulting in only small probability of violation (6%).

In terms of objective values, unlike the first case study, oracle Q-learning with backoffs saw a significant decrease in objective value in Table 2 after applying backoffs. This is expected because we further restrict the feasible space of the controller leading to a more conservative solution, hence exhibiting a trade-off between constraint satisfaction and objective value.

This trade-off is justified as the NMPC solution results in 66% probability of constraint violation (Table 2). In the context of a chemical plant, the NMPC solution is infeasible due to the high risk. Process operation in such industries necessitates that these probabilities are minimized as safety is of utmost importance in chemical engineering. This provides basis for the use of RL in the process industries.

On the other hand, for oracle Q-learning, it can be seen that the probability of constraint violation has been significantly improved from 66% (for NMPC) to 9% (boldface in Table 2). Clearly, oracle Q-learning offers an effective means of not only satisfying constraints in expectation (green lines in Fig. 6), but more importantly with high probability (all green shaded areas below zero).

However, it is worth noting that this algorithm is based on Q-learning, which is expected to take longer time to train than MPC, particularly because it requires backoffs to be tuned. This is a direct consequence of shifting the computation time from online to offline. Indeed, such a tradeoff can be justified as this guarantees robust constraint satisfaction *online* with fast computation time, which is crucial in many safety critical engineering applications.

5. Conclusions

In this paper we propose a new reinforcement learning methodology for finding a controller policy that can satisfy constraints with high probability in stochastic and complex process systems. The proposed algorithm - oracle-assisted constrained Q-learning - uses constraint tightening by applying backoffs to the original feasible set. Backoffs restrict the perceived feasible space by the controller, hence allowing guarantees on the satisfaction of chance constraints. Here, we find the smallest backoffs (least conservative) that still guarantee the desired probability of satisfaction by solving a root-finding problem using Broyden’s method. Results show that our proposed methodology compares favorably to nonlinear model predictive control (NMPC), a benchmark

control technique commonly used in the industry, in terms of constraint handling. This is expected since NMPC guarantees constraint satisfaction only in *expectation* (loosely speaking constraints are satisfied only 50% of the time), while our algorithm ensures constraint satisfaction with probabilities as high as 99% as shown in the case studies. Being able to solve constraint policy optimization problems with high probability constraint satisfaction has been one of the main hurdles of the widespread use of RL in engineering applications. The promising performance of this algorithm is an encouraging step towards applying RL to real-world industrial chemical processes, where constraints on policies are absolutely critical due to safety reasons.

Acknowledgment

This project has received funding from the EPSRC projects (EP/R032807/1) and (EP/P016650/1).

References

- [1] P. Abbeel, A. Y. Ng, 2004. Apprenticeship learning via inverse reinforcement learning. In: Proceedings of the twenty-first international conference on Machine learning. p. 1.
- [2] J. Achiam, D. Held, A. Tamar, P. Abbeel, 2017. Constrained policy optimization.
- [3] E. Altman, 1999. Constrained Markov decision processes. Vol. 7. CRC Press.
- [4] D. P. Bertsekas, D. P. Bertsekas, D. P. Bertsekas, D. P. Bertsekas, 1995. Dynamic programming and optimal control. Vol. 1. Athena scientific Belmont, MA.
- [5] E. Bradford, L. Imsland, 2018. Economic stochastic model predictive control using the unscented kalman filter. IFAC-PapersOnLine 51 (18), 417–422.
- [6] E. Bradford, L. Imsland, D. Zhang, E. A. del Rio Chanona, 2020. Stochastic data-driven model predictive control using gaussian processes. Computers & Chemical Engineering 139, 106844.
- [7] J. Buckman, 2021. How to think about replay memory.
URL <https://jacobbuckman.com/2021-02-13-how-to-think-about-replay-memory/>
- [8] R. Cheng, G. Orosz, R. M. Murray, J. W. Burdick, 2019. End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks. In: Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. pp. 3387–3395.
- [9] J. Choi, F. Castañeda, C. J. Tomlin, K. Sreenath, 2020. Reinforcement learning for safety-critical control under model uncertainty, using control lyapunov functions and control barrier functions. arXiv preprint arXiv:2004.07584.
- [10] Y. Chow, O. Nachum, A. Faust, E. Duenez-Guzman, M. Ghavamzadeh, 2019. Lyapunov-based safe policy optimization for continuous control.
URL <http://arxiv.org/abs/1901.10031>
- [11] G. Chowdhary, M. Liu, R. Grande, T. Walsh, J. How, L. Carin, 2014. Off-policy reinforcement learning with gaussian processes. IEEE/CAA Journal of Automatica Sinica 1 (3), 227–238.
- [12] G. Dalal, K. Dvijotham, M. Vecerik, T. Hester, C. Paduraru, Y. Tassa, 2018. Safe exploration in continuous action spaces. CoRR abs/1801.08757.
URL <http://arxiv.org/abs/1801.08757>
- [13] E. A. del Rio Chanona, P. Petsagkourakis, E. Bradford, J. E. A. Graciano, B. Chachuat, 2021. Real-time optimization meets bayesian optimization and derivative-free optimization: A tale of modifier adaptation. Computers & Chemical Engineering 147, 107249.
- [14] L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, A. Madry, 2020. Implementation matters in deep policy gradients: A case study on ppo and trpo. arXiv preprint arXiv:2005.12729.
- [15] D. A. Goulart, R. D. Pereira, 2020. Autonomous ph control by reinforcement learning for electroplating industry wastewater. Computers & Chemical Engineering, 106909.
- [16] T. Haarnoja, A. Zhou, P. Abbeel, S. Levine, 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. arXiv preprint arXiv:1801.01290.
- [17] P. J. Huber, 1964. Robust estimation of a location parameter: Annals mathematics statistics, 35.
- [18] S. Huh, I. Yang, 2020. Safe reinforcement learning for probabilistic reachability and safety specifications: A lyapunov-based approach. arXiv preprint arXiv:2002.10126.
- [19] S. Hwangbo, G. Sin, 2020. Design of control framework based on deep reinforcement learning and monte-carlo sampling in downstream separation. Computers & Chemical Engineering, 106910.
- [20] C. T. Kelley, 1995. Iterative methods for linear and nonlinear equations. SIAM.

- [21] N. P. Lawrence, G. E. Stewart, P. D. Loewen, M. G. Forbes, J. U. Backstrom, R. B. Gopaluni, 2020. Optimal pid and antiwindup control design as a reinforcement learning problem. arXiv preprint arXiv:2005.04539.
- [22] J. H. Lee, J. M. Lee, 2006. Approximate dynamic programming based approach to process control and scheduling. *Computers & chemical engineering* 30 (10-12), 1603–1618.
- [23] J. M. Lee, J. H. Lee, 2005. Approximate dynamic programming-based approaches for input-output data-driven control of nonlinear processes. *Automatica* 41 (7), 1281–1288.
- [24] J. Lehman, J. Chen, J. Clune, K. O. Stanley, 2018. Es is more than just a traditional finite-difference approximator. In: *Proceedings of the Genetic and Evolutionary Computation Conference*. pp. 450–457.
- [25] E. Leurent, D. Efimov, O.-A. Maillard, 2020. Robust-adaptive control of linear systems: beyond quadratic costs.
- [26] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, D. Wierstra, 2015. Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971.
- [27] L.-J. Lin, 1993. Reinforcement learning for robots using neural networks. Tech. rep., Carnegie-Mellon Univ Pittsburgh PA School of Computer Science.
- [28] S. Mehta, L. A. Ricardez-Sandoval, 2016. Integration of design and control of dynamic systems under uncertainty: A new back-off approach. *Industrial & Engineering Chemistry Research* 55 (2), 485–498.
- [29] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, M. Riedmiller, 2013. Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602.
- [30] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al., 2015. Human-level control through deep reinforcement learning. *nature* 518 (7540), 529–533.
- [31] M. Mowbray, R. Smith, E. A. Del Rio-Chanona, D. Zhang, 2021. Using process data to generate an optimal control policy via apprenticeship and reinforcement learning. Submitted to Journal
- [32] A. Y. Ng, D. Harada, S. Russell, 1999. Policy invariance under reward transformations: Theory and application to reward shaping. In: *ICML*. Vol. 99. pp. 278–287.
- [33] J. Nocedal, S. Wright, 2006. Numerical optimization. Springer Science & Business Media.
- [34] B. Peng, Y. Mu, J. Duan, Y. Guan, S. E. Li, J. Chen, 2021. Separated proportional-integral lagrangian for chance constrained reinforcement learning.
- [35] P. Petsagkourakis, F. Galvanin, 2020. Safe model-based design of experiments using gaussian processes.
- [36] P. Petsagkourakis, I. O. Sandoval, E. Bradford, F. Galvanin, D. Zhang, E. A. del Rio-Chanona, 2020. Chance constrained policy optimization for process control and optimization. arXiv preprint arXiv:2008.00030.
- [37] P. Petsagkourakis, I. O. Sandoval, E. Bradford, D. Zhang, E. A. del Rio-Chanona, 2020. Reinforcement learning for batch bioprocess optimization. *Computers & Chemical Engineering* 133, 106649.
- [38] L. D. Pyeatt, A. E. Howe, et al., 2001. Decision tree function approximation in reinforcement learning. In: *Proceedings of the third international symposium on adaptive systems: evolutionary computation and probabilistic graphical models*. Vol. 2. Cuba, pp. 70–77.
- [39] M. Rafiei, L. A. Ricardez-Sandoval, 2018. Stochastic back-off approach for integration of design and control under uncertainty. *Industrial & Engineering Chemistry Research* 57 (12), 4351–4365.
- [40] M. Rafiei, L. A. Ricardez-Sandoval, 2020. A trust-region framework for integration of design and control. *AIChE Journal* 66 (5), e16922.
- [41] D. Rolnick, A. Ahuja, J. Schwarz, T. P. Lillicrap, G. Wayne, 2019. Experience replay for continual learning.
- [42] R. H. Russel, M. Benosman, J. V. Baar, 2020. Robust constrained-mdps: Soft-constrained robust policy optimization under model uncertainty.
- [43] M. Ryu, Y. Chow, R. Anderson, C. Tjandraatmadja, C. Boutilier, 2019. Caql: Continuous action q-learning. arXiv preprint arXiv:1909.12397.
- [44] I. Sajedian, T. Badloe, J. Rho, 2019. Optimisation of colour generation from dielectric nanostructures using reinforcement learning. *Optics express* 27 (4), 5874–5883.
- [45] H. Satija, P. Amortila, J. Pineau, 2020. Constrained markov decision processes via backward value functions. arXiv preprint arXiv:2008.11811.
- [46] J. Shin, T. A. Badgwell, K.-H. Liu, J. H. Lee, 2019. Reinforcement learning—overview of recent progress and implications for process control. *Computers & Chemical Engineering* 127, 282–294.
- [47] V. Singh, H. Kodamana, 2020. Reinforcement learning based control of batch polymerisation processes. *IFAC-PapersOnLine* 53 (1), 667–672.
- [48] A. Slowik, H. Kwasnicka, 2020. Evolutionary algorithms and their applications to engineering problems. *Neural Computing and Applications*, 1–17.
- [49] S. Spielberg, A. Tulsyan, N. P. Lawrence, P. D. Loewen, R. Bhushan Gopaluni, 2019. Toward self-driving processes: A deep reinforcement learning approach to control. *AIChE Journal* 65 (10), e16689.
- [50] R. S. Sutton, A. G. Barto, 2018. Reinforcement learning: An introduction 2nd ed.
- [51] R. S. Sutton, D. A. McAllester, S. P. Singh, Y. Mansour, 2000. Policy gradient methods for reinforcement learning with function approximation. In: *Advances in neural information processing systems*. pp. 1057–1063.

- [52] C. Szepesvári, 2010. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning* 4 (1), 1–103.
- [53] A. Taylor, A. Singletary, Y. Yue, A. Ames, 2020. Learning for safety-critical control with control barrier functions. In: *Learning for Dynamics and Control*. PMLR, pp. 708–717.
- [54] C. Tessler, D. J. Mankowitz, S. Mannor, 2018. Reward constrained policy optimization. *arXiv preprint arXiv:1805.11074*.
- [55] K. P. Wabersich, M. N. Zeilinger, 2018. Safe exploration of nonlinear dynamical systems: A predictive safety filter for reinforcement learning. *arXiv preprint arXiv:1812.05506*.
- [56] A. Wächter, 2002. An interior point algorithm for large-scale nonlinear optimization with applications in process engineering. Ph.D. thesis, PhD thesis, Carnegie Mellon University.
- [57] Z. Wang, H.-X. Li, C. Chen, 2019. Incremental reinforcement learning in continuous spaces via policy relaxation and importance weighting. *IEEE Transactions on Neural Networks and Learning Systems*.
- [58] C. J. Watkins, P. Dayan, 1992. Q-learning. *Machine learning* 8 (3-4), 279–292.
- [59] H. Xie, X. Xu, Y. Li, W. Hong, J. Shi, 2020. Model predictive control guided reinforcement learning control scheme. In: *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, pp. 1–8.
- [60] Z. Zhou, S. Kearnes, L. Li, R. N. Zare, P. Riley, 2019. Optimization of molecules via deep reinforcement learning. *Scientific reports* 9 (1), 1–10.