

# On PI Controllers for Updating Lagrange Multipliers in Constrained Optimization

Motahareh Sohrabi<sup>\*†</sup> Juan Ramirez<sup>\*†</sup> Tianyue H. Zhang<sup>†</sup> Simon Lacoste-Julien<sup>†‡</sup> Jose Gallego-Posada<sup>†</sup>

## Abstract

Constrained optimization offers a powerful framework to prescribe desired behaviors in neural network models. Typically, constrained problems are solved via their min-max Lagrangian formulations, which exhibit unstable oscillatory dynamics when optimized using gradient descent-ascent. The adoption of constrained optimization techniques in the machine learning community is currently limited by the lack of reliable, general-purpose update schemes for the Lagrange multipliers. This paper proposes the  $\nu$ PI algorithm and contributes an optimization perspective on Lagrange multiplier updates based on PI controllers, extending the work of [Stooke et al. \(2020\)](#). We provide theoretical and empirical insights explaining the inability of momentum methods to address the shortcomings of gradient descent-ascent, and contrast this with the empirical success of our proposed  $\nu$ PI controller. Moreover, we prove that  $\nu$ PI generalizes popular momentum methods for single-objective minimization. Our experiments demonstrate that  $\nu$ PI reliably stabilizes the multiplier dynamics and its hyperparameters enjoy robust and predictable behavior.

## 1. Introduction

The need to enforce complex behaviors in neural network models has reinvigorated the interest of the machine learning community in constrained optimization techniques. Recent applications include fairness ([Cotter et al., 2019](#); [Zafar et al., 2019](#); [Fioretto et al., 2020](#); [Hashemizadeh et al., 2024](#)), sparsity ([Gallego-Posada et al., 2022](#)), active learning ([Elenter et al., 2022](#)), reinforcement learning ([Stooke et al., 2020](#); [Farahmand & Ghavamzadeh, 2021](#)) and model quantization ([Hounie et al., 2023](#)).

<sup>\*</sup> Equal contribution. <sup>†</sup> Mila—Quebec AI Institute and DIRO, Université de Montréal. <sup>‡</sup> Canada CIFAR AI Chair. Correspondence to: Juan Ramirez <[juan.ramirez@mila.quebec](mailto:juan.ramirez@mila.quebec)>.

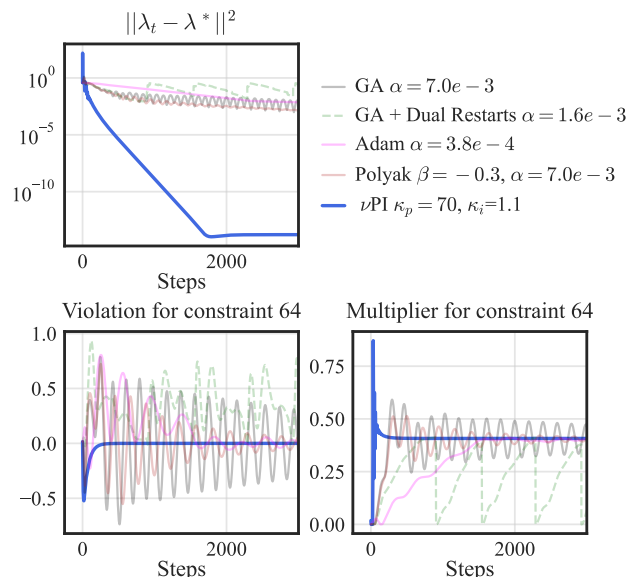


Figure 1: Dynamics for different dual optimizers on a hard-margin SVM problem (Eq. (12)). Amongst the tested methods,  **$\nu$ PI is the only method to successfully converge to the optimal dual variables**. Each optimizer uses the best hyperparameters found after a grid-search aiming to minimize the distance to the optimal  $\lambda^*$  after 5,000 steps. For improved readability, the plot shows the first 3,000 steps. Constraint 64 corresponds to a support vector. All methods achieved perfect training accuracy.

Algorithmic approaches based on the Lagrangian min-max representation of the original constrained optimization problem ([Boyd & Vandenberghe, 2004](#), §5) are commonly preferred in the context of neural networks since (i) they are amenable to inexact, gradient-based optimization ([Bertsekas, 2016](#), §5.2), (ii) making it easy to incorporate constraints into existing pipelines for unconstrained optimization ([Cotter et al., 2019](#); [Gallego-Posada & Ramirez, 2022](#)), and (iii) they do not require special structure in the objective or constraint functions (such as convexity or efficient projection onto the feasible set ([Nocedal & Wright, 2006](#))).

Despite their wider applicability, solving Lagrangian problems involving neural networks is challenging as it simultaneously entails the difficulties of nonconvex optimization on large-scale models ([Bottou et al., 2018](#)), and the potential for instability and oscillations due to the adversarial min-max nature of the Lagrangian ([Stooke et al., 2020](#)).

Lagrangian problems are commonly optimized using some variant of gradient-descent ascent (GDA) (Arrow et al., 1958). Despite local convergence results in idealized settings (Lin et al., 2020; Zhang et al., 2022), the optimization dynamics of GDA typically exhibit instabilities, overshoot or oscillations (Platt & Barr, 1987; Gidel et al., 2019a; Stooke et al., 2020; Gallego-Posada et al., 2022).

Alleviating the shortcomings of GDA on Lagrangian problems is an important step towards widespread adoption constrained optimization in deep learning. Recently, Stooke et al. (2020) proposed a solution based on a PID controller (Åström & Hägglund, 1995) for updating the Lagrange multipliers in safety-constrained reinforcement learning problems. Our manuscript expands on their work by providing an optimization-oriented analysis of  $\nu$ PI (Algo. 1), a related PI controller that incorporates an exponential moving average on the error signal.

Fig. 1 illustrates how our proposed  $\nu$ PI controller successfully dampens the oscillations on a hard-margin SVM task, achieving fast convergence to the optimal Lagrange multipliers. In contrast, a wide range of popular methods for single-objective minimization exhibit unstable, oscillatory dynamics and fail to converge in this task. See §5.1 for further details on this experiment.

**Contributions:** ① We introduce the  $\nu$ PI algorithm (§4) and prove that  $\nu$ PI generalizes popular momentum methods like POLYAK and NESTEROV (Thm. 1), as well as traditional PI controllers. ② We provide conceptual insights explaining how  $\nu$ PI improves the dynamics of the Lagrange multipliers: §4.3 presents a qualitative analysis of the updates executed by the  $\nu$ PI algorithm in contrast to gradient ascent; in §4.4 we study the spectral properties of the continuous-time system. ③ In §4.5, we provide a heuristic to tune the new hyperparameter  $\kappa_p$  of the  $\nu$ PI algorithm; we also demonstrate that it has a monotonic effect in the damping of oscillations. ④ Our experiments on hard-margin SVMs, sparsity tasks using ResNets, and algorithmic fairness demonstrate that  $\nu$ PI leads to improved stability and convergence.

**Code:** [github.com/motahareh-sohrabi/nuPI](https://github.com/motahareh-sohrabi/nuPI)

**Scope:** Due to the highly specialized techniques used for training neural networks (Dahl et al., 2023), in this work we concentrate on iterative schemes that do *not* modify the optimization protocol used on the model parameters. In other words, we restrict our attention to update schemes on the Lagrange multipliers only, which allows us to reuse the same optimizer choices for the (primal) model parameters as used in the unconstrained setting.

## 2. Related Works

**Constrained optimization.** We are interested in Lagrangian methods (Arrow et al., 1958) that allow tackling general

(nonconvex) constrained optimization problems with differentiable objective and constraints. Classical constrained optimization (Nocedal & Wright, 2006; Bertsekas, 2016) techniques include projection methods (Bertsekas, 1976), barrier methods (Dikin, 1967), and methods of feasible directions (Frank & Wolfe, 1956; Zoutendijk, 1960). These approaches usually make assumptions on the structure of the problem, such as convexity of the objective or constraints, the existence of an efficient projection operator onto the feasible set, or access to a linear minimization oracle. Such assumptions restrict their applicability to deep learning tasks. Other popular techniques such as penalty methods (Nocedal & Wright, 2006) and the method of multipliers (Bertsekas, 1975), apply to general nonconvex problems, but are outside the scope of this work.

**Min-max optimization.** The Lagrangian formulation of a nonconvex constrained optimization problem leads to a nonconvex concave min-max problem. Under idealized assumptions, gradient descent-ascent has local convergence guarantees for said problems (Lin et al., 2020), but may exhibit oscillations (Platt & Barr, 1987; Gidel et al., 2019b). Under stronger assumptions, extragradient (Korpelevich, 1976) and the optimistic gradient method (Popov, 1980) converge at a nearly optimal rate (Mokhtari et al., 2020a). These methods, as well as POLYAK with negative momentum (Gidel et al., 2019a) and PID controllers (Stooke et al., 2020), have been shown to dampen the oscillations of GDA. However, negative momentum may be suboptimal for strongly convex-strongly concave min-max problems (Zhang & Wang, 2021).

Our work focuses on the *dynamics* of Lagrangian games. We provide insights on why popular techniques for minimization may exacerbate oscillations and overshoot, and why PI controllers can be effective at damping oscillations. Our proposed method  $\nu$ PI is a generalization of both (negative) momentum and the optimistic gradient method.

**PID controllers and optimization.** An et al. (2018) studied PID control for training machine learning models by considering the negative loss gradient as the error signal to the controller. PID controllers have been shown to generalize gradient descent (Hu & Lessard, 2017) and momentum (Recht, 2018). Stooke et al. (2020); Casti et al. (2023) have highlighted the effectiveness of controllers at optimizing constrained optimization tasks.

In this work, we propose a PI-like update rule for the dual variables in a Lagrangian min-max game. We prove our algorithm generalizes momentum methods and we provide conceptual insights to support the empirical effectiveness of PI controllers in reducing oscillations and overshoot in the constrained optimization dynamics. In Appx. A, we elaborate on the distinctions between our work and existing research on PID controllers for optimization.

### 3. Lagrangian Optimization

Consider a constrained optimization problem with  $m$  inequality and  $n$  equality constraints, represented by functions  $g : \mathcal{X} \rightarrow \mathbb{R}^m$  and  $h : \mathcal{X} \rightarrow \mathbb{R}^n$ , respectively:

$$\min_{\mathbf{x}} f(\mathbf{x}) \text{ subject to } \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \text{ and } \mathbf{h}(\mathbf{x}) = \mathbf{0}. \quad (1)$$

We do not make any assumptions on the functions  $f$ ,  $g$ , and  $h$  beyond almost-everywhere differentiability. We refer to the values of  $g$  and  $h$  as the *constraint violations*. In particular, we are interested in optimization problems where  $\mathbf{x}$  corresponds to the parameters of a neural network, leading to objective and constraint functions that may be nonconvex. This typically precludes the use of “classical” constrained optimization methods, as those discussed in §2.

The Lagrangian min-max problem associated with the constrained optimization problem in Eq. (1) is given by:

$$\min_{\mathbf{x}} \max_{\lambda \geq 0, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) \triangleq f(\mathbf{x}) + \lambda^\top \mathbf{g}(\mathbf{x}) + \mu^\top \mathbf{h}(\mathbf{x}), \quad (2)$$

where  $\lambda$  and  $\mu$  are vectors of *Lagrange multipliers* associated with the inequality and equality constraints, respectively. Eq. (2) constitutes a nonconvex-concave zero-sum game between  $\mathbf{x}$  (known as the *primal player*) and  $\{\lambda, \mu\}$  (known as the *dual player*). We are interested in algorithmic approaches that identify saddle points of the Lagrangian  $\mathcal{L}(\mathbf{x}, \lambda, \mu)$  as these correspond to constrained optima.

In general, Lagrangian-based approaches do not constitute *feasible methods* (i.e. visiting only feasible iterates). We judge a method’s success based on its asymptotic feasibility, or at the end of a pre-determined optimization budget.

**Simultaneous updates.** The simplest algorithm to solve the problem in Eq. (2) is simultaneous gradient descent-ascent (GDA) (Arrow et al., 1958):

$$\begin{cases} \mu_{t+1} \leftarrow \mu_t + \eta_{\text{dual}} \nabla_{\mu} \mathcal{L}(\mathbf{x}_t, \lambda_t, \mu_t) = \mu_t + \eta_{\text{dual}} \mathbf{h}(\mathbf{x}_t) \\ \hat{\lambda}_{t+1} \leftarrow \lambda_t + \eta_{\text{dual}} \nabla_{\lambda} \mathcal{L}(\mathbf{x}_t, \lambda_t, \mu_t) = \lambda_t + \eta_{\text{dual}} \mathbf{g}(\mathbf{x}_t) \\ \lambda_{t+1} \leftarrow \Pi_{\mathbb{R}_+^m}(\hat{\lambda}_{t+1}) = \max\left(0, \hat{\lambda}_{t+1}\right) \\ \mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta_{\text{primal}} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t, \lambda_t, \mu_t), \end{cases}$$

where the middle two equations execute a projected gradient-ascent step enforcing the non-negativity of the multipliers  $\lambda$ .

To simplify notation, we will group the dual variables as  $\theta = [\lambda, \mu]^\top$  and the constraints  $\mathbf{c}(\mathbf{x}) = [\mathbf{g}(\mathbf{x}), \mathbf{h}(\mathbf{x})]^\top$  which yields the concise Lagrangian problem:

$$\min_{\mathbf{x}} \max_{\theta \in \mathbb{R}_+^m \times \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \theta) \triangleq f(\mathbf{x}) + \theta^\top \mathbf{c}(\mathbf{x}) \quad (3)$$

Note that the primal update direction  $\nabla_{\mathbf{x}} \mathcal{L}$  is a linear combination of the objective and constraint gradients—which can be efficiently computed using automatic differentiation,

without storing  $\nabla f$  and  $\mathcal{J}c^1$  separately. On the other hand,  $\nabla_{\theta} \mathcal{L} = \mathbf{c}(\mathbf{x})$ , and thus the GDA update on the multipliers corresponds to the *integration* (i.e. accumulation) of the constraint violations over time. We highlight that the cost of updating the Lagrange multipliers is typically negligible relative to the cost of computing  $f$  and  $\mathbf{c}$ .

**Alternating updates.** Prior work has demonstrated the advantages of alternating updates in min-max optimization: Zhang et al. (2022) established that alternating GDA achieves a near-optimal local convergence rate for strongly concave-strongly convex problems (strictly better than simultaneous GDA); Gidel et al. (2019b) showed that alternating GDA leads to bounded iterates on smooth bilinear games, as opposed to divergence for simultaneous updates. Besides the improved convergence and stability benefits, alternating updates are particularly suitable for Lagrangian games from a computational standpoint due to the linear structure of the Lagrangian with respect to the dual variables. Concretely, consider the *alternating* update scheme:

$$\begin{cases} \hat{\theta}_{t+1} \leftarrow \theta_t + \eta_{\text{dual}} \nabla_{\theta} \mathcal{L}(\mathbf{x}_t, \theta_t) = \theta_t + \eta_{\text{dual}} \mathbf{c}(\mathbf{x}_t) \\ \theta_{t+1} \leftarrow \Pi_{\mathbb{R}_+^m \times \mathbb{R}^n}(\hat{\theta}_{t+1}) \\ \mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta_{\text{primal}} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t, \theta_{t+1}) \\ \quad = \mathbf{x}_t - \eta_{\text{primal}} (\nabla f(\mathbf{x}_t) + \mathcal{J}c(\mathbf{x}_t) \theta_t) \end{cases} \quad (4)$$

The alternating updates in Eq. (4), only require computing  $f(\mathbf{x}_t)$  and  $\mathbf{c}(\mathbf{x}_t)$  once, just as when performing simultaneous updates. In a general zero-sum game, where  $\mathcal{L}(\mathbf{x}_t, \theta_t)$  does not decouple as in the Lagrangian case, the second part of the alternation might require re-evaluating  $\mathcal{L}(\mathbf{x}_t, \theta_{t+1})$  entirely. However, note that thanks to the affine structure of  $\mathcal{L}$  with respect to  $\theta$ , the update on  $\mathbf{x}$  can be calculated efficiently without having to re-evaluate  $f$  or  $\mathbf{c}$ .

These theoretical and practical advantages motivate our decision to concentrate on alternating update schemes like Eq. (4) for solving the problem in Eq. (3) in what follows.

**Practical remarks.** In practice, updates on the primal variables require more sophisticated methods (with intricate hyperparameter tuning) than the plain gradient descent update presented in Eq. (4) to achieve good performance, including any number of highly specialized procedures developed for training neural networks (Dahl et al., 2023).

Moreover, for certain applications, a training pipeline designed to minimize a single, *unconstrained* objective might be in place. In these cases, it is desirable to develop update schemes for the Lagrange multipliers that allow for seamlessly incorporating constraints into the model development pipeline *without having to engineer from scratch a new recipe for training the model*.

<sup>1</sup> $\mathcal{J}f \triangleq [\nabla f_1 \quad \dots \quad \nabla f_p] \in \mathbb{R}^{d \times p}$  denotes the (transpose) Jacobian matrix of a function  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^p$ .

In this paper, we concentrate on different update schemes for the Lagrange multipliers and assume that a well-tuned optimizer for the model parameters is available.

**Shortcomings of gradient ascent.** As mentioned previously, gradient ascent (GA) on the Lagrange multipliers corresponds to accumulating the observed constraint violations over time. For simplicity, let us concentrate on a single inequality constraint  $c(\mathbf{x})$ . Whenever the constraint is being violated (resp. satisfied), the violation is positive  $c(\mathbf{x}) > 0$  (resp. negative) and thus the value of the corresponding multiplier is increased (resp. decreased) by  $\eta_{\text{dual}} c(\mathbf{x})$ . Recall that the projection step ensures that the inequality multipliers remain non-negative.

Therefore, the value of the multiplier depends on the entire optimization trajectory through the value of the observed violations. In particular, after a long period of infeasibility, the value of the multiplier will be large, biasing the gradient  $\nabla_{\mathbf{x}} \mathcal{L}$  towards reducing the violation and thus improving the feasibility of the model.

An insufficient increase of the multiplier will cause the constraint to be *ignored*, while an excessively large value of the multiplier will lead the constraint to be enforced *beyond* the prescribed constraint level. The latter behavior can also occur if the multiplier fails to decrease sufficiently fast once the constraint is satisfied. Repeated cycles of insufficient or excessive change in the multiplier manifest in ignoring or overshooting, thus forming oscillations. See Figs. 1 and 2 for illustrations of these behaviors.

— GA    — GA + Dual Rest.    — Polyak  $\beta = 0.3$     —  $\nu$ PI  $\kappa_p = 16.0$

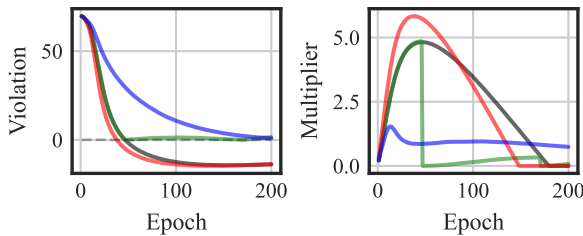


Figure 2: Constraint dynamics for GA, POLYAK and  $\nu$ PI in a sparsity task (§5.3). Constrained optimal solutions for this problem lie at the boundary of the feasible set. The excessive growth in the value of the multiplier for GA causes the constraint to overshoot into the interior of the feasible set. **The improved multiplier updates of the  $\nu$ PI algorithm remove the overshoot in the constraint and multiplier.**

In short, an ideal update rule for the multiplier would behave *adaptively*, based on the observed violations throughout the execution of the optimization. This begs the question of whether existing adaptive optimization such as POLYAK, NESTEROV and ADAM would reliably resolve these issues. Sections 4 and 5 provide a *negative* answer to this question.

**Dual restarts.** Gallego-Posada et al. (2022) proposed an

approach to mitigate the overshoot in inequality constraints called *dual restarts*: once a constraint is *strictly* satisfied, its associated dual variable is reset to zero. This corresponds to a best response (in game-theoretic terms) of the dual player. Dual restarts prevent excessive enforcement of constraints, which can degrade the achieved objective function value.

However, dual restarts are not suitable for general constrained optimization problems since they rely on determining the satisfaction of the constraint *exactly*. Constraint estimates may (wrongly) indicate strict feasibility due to (i) stochasticity in their estimation, (ii) numerical precision errors making active constraints appear strictly feasible, or (iii) a “temporary” strict satisfaction of the constraint. Fig. 1 illustrates the undesirable dynamics caused by dual restarts when applied to an SVM task in which the support vectors correspond to strictly active inequality constraints.

In §4, we show that  $\nu$ PI mitigates the overshoot of inequality constraints, with additional benefits: (i) controllable degree of overshoot (governed by the  $\kappa_p$  hyperparameter), (ii) compatibility with equality (and strictly feasible inequality) constraints, and (iii) damping of multiplier oscillations.

## 4. $\nu$ PI Control for Constrained Optimization

Following Stooke et al. (2020), we consider the learning of an optimal feasible model solving Eq. (1) as a dynamical system. Thus, we can think of the update rule for the multipliers as a control algorithm that aims to *steer the system toward feasibility*. We emphasize that we are not trying to control general dynamical systems, but rather systems that arise from partial, inexact minimization (e.g. gradient-based updates) on a min-max Lagrangian game. In other words, we assume that  $\mathbf{x}_{t-1} \mapsto \mathbf{x}_t$  is updated so as to minimize the current Lagrangian  $\mathcal{L}(\cdot, \boldsymbol{\theta}_t)$ . Figure 3 illustrates the control pipeline we consider throughout this work.

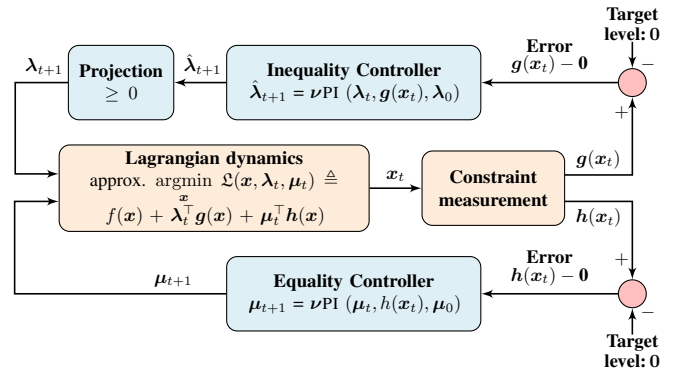


Figure 3:  $\nu$ PI control pipeline for updating the Lagrange multipliers in a constrained optimization problem. We consider the update on the primal variables as a black-box procedure that receives the multipliers  $\lambda_t$  and primal variables  $\mathbf{x}_{t-1}$  as input, and returns an updated  $\mathbf{x}_t$ . The multiplier update is executed by the controller, using the constraint violations as the error signal.

An assumption entailed by the control perspective in Fig. 3 is that an increase (resp. decrease) in the control variable (the Lagrange multipliers  $\theta_t$ ) leads to a decrease (resp. increase) in the controlled quantity (the constraint violations  $c(x_t)$ ). This assumption holds for constrained optimization problems since an increase in the multipliers leads the primal minimization of the Lagrangian to focus on reducing the value of the constraints (as mentioned during the discussion of gradient descent-ascent dynamics in §3).

Note that our black-box assumption on the nature of the primal update allows for an arbitrary choice of optimizer for minimizing  $\mathcal{L}(\cdot, \theta_t)$ . After obtaining an updated primal iterate  $x_t$ , the new constraint violations  $g(x_t)$  and  $h(x_t)$  are measured and used as the error signals for the inequality- and equality-constraint controllers, yielding updated multipliers  $\theta_{t+1}$ . The projection block ensures the non-negativity of the multipliers for inequality constraints.

#### 4.1. $\nu$ PI algorithm

Our main algorithmic contribution is the multiplier update scheme presented in Algo. 1. This is a simple generalization of a PI controller (i.e. a PID controller (Åström & Hägglund, 1995) with  $\kappa_d = 0$ ) by including an exponential moving average (of the error signal) in the proportional term. Indeed, the traditional PI controller is recovered when  $\nu = 0$ .

##### Algorithm 1 $\nu$ PI update

**Args:** EMA coefficient  $\nu$ , proportional ( $\kappa_p$ ) and integral ( $\kappa_i$ ) gains; initial conditions  $\xi_0$  and  $\theta_0$ .

- 1: Measure current system error  $e_t$
- 2:  $\xi_t \leftarrow \nu \xi_{t-1} + (1 - \nu)e_t$  ▷ for  $t \geq 1$
- 3:  $\theta_{t+1} \leftarrow \theta_0 + \kappa_p \xi_t + \kappa_i \sum_{\tau=0}^t e_\tau$

The  $\nu$ PI update can be equivalently expressed in terms of a recursive update (see Lemma 2 in Appx. B) as:

$$\theta_1 = \theta_0 + \kappa_i e_0 + \kappa_p \xi_0 \quad (5)$$

$$\theta_{t+1} = \theta_t + \kappa_i e_t + \kappa_p (\xi_t - \xi_{t-1}) \quad \text{for } t \geq 1. \quad (6)$$

#### 4.2. Connections to optimization methods

When the error signal corresponds to the negative gradient of a cost function  $e_t = -\nabla f_t$ , Algo. 1 has straightforward equivalences with common minimization methods (see Appx. B). For example,  $\nu$ PI ( $\nu = 0, \kappa_p = 0, \kappa_i$ ) is equivalent to GD ( $\alpha = \kappa_i$ ) (Stoake et al., 2020; Lessard et al., 2016; An et al., 2018). When  $\nu = 0$  and  $\kappa_p = \kappa_i = \alpha$ ,  $\nu$ PI recovers a single-player version of the OPTIMISTICGRADIENT (OG) method (Popov, 1980), with step-size  $\alpha$ . When  $\nu = 0$ , but  $\kappa_p$  and  $\kappa_i$  are allowed to differ,  $\nu$ PI coincides with the generalized OG studied by Mokhtari et al. (2020b). Since we use  $\nu$ PI for updating the multipliers, we phrase the updates in Algo. 1 based on a maximization convention.

Moreover, our proposed algorithm  $\nu$ PI generalizes popular momentum methods such as POLYAK—also known as HEAVYBALL—(Polyak, 1964) and NESTEROV (Nesterov, 1983).<sup>2</sup> This connection, stated formally in Thm. 1, will allow us to understand (§4.3) why traditional momentum methods are *insufficient* to address the shortcomings of gradient ascent for Lagrangian optimization.

We take advantage of the UNIFIEDMOMENTUM ( $\alpha, \beta, \gamma$ ) framework introduced by Shen et al. (2018) to concisely develop a joint analysis of POLYAK( $\alpha, \beta$ ) = UM( $\alpha, \beta, \gamma = 0$ ) and NESTEROV( $\alpha, \beta$ ) = UM( $\alpha, \beta, \gamma = 1$ ).

##### Algorithm 2 UNIFIEDMOMENTUM update (Shen et al., 2018)

**Args:** step-size  $\alpha$ , momentum coefficient  $\beta$ , interpolation factor  $\gamma \in [0, \frac{1}{1-\beta}]$ ; initial conditions  $\phi_0 = \mathbf{0}$  and  $\theta_0$ .

- 1: Measure current system error  $e_t$
- 2:  $\phi_{t+1} \leftarrow \beta \phi_t + \alpha e_t$
- 3:  $\theta_{t+1} \leftarrow \theta_t + \phi_{t+1} + \beta \gamma (\phi_{t+1} - \phi_t)$

**Theorem 1.** [Proof in Appx. B.] *Under the same initialization  $\theta_0$ , UNIFIEDMOMENTUM( $\alpha, \beta \neq 1, \gamma$ ) is a special case of the  $\nu$ PI algorithm with the hyperparameter choices:*

$$\nu \leftarrow \beta \quad \xi_0 \leftarrow (1 - \beta)e_0 \quad (7)$$

$$\kappa_i \leftarrow \frac{\alpha}{1 - \beta} \quad \kappa_p \leftarrow -\frac{\alpha\beta}{(1 - \beta)^2} [1 - \gamma(1 - \beta)]. \quad (8)$$

Table 2 in Appx. B summarizes the connections we have established between  $\nu$ PI and existing methods. We emphasize that the exponential moving average in  $\nu$ PI is a crucial component to obtain the generalization of momentum methods.

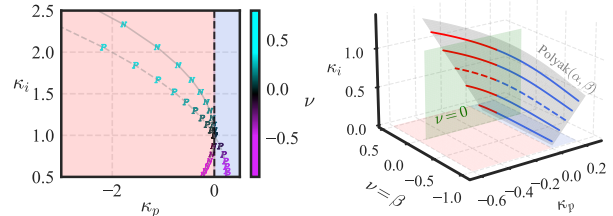


Figure 4: **Left:** Hyperparameter choices from Thm. 1 for which  $\nu$ PI ( $\nu, \kappa_p, \kappa_i$ ) realizes POLYAK( $\alpha = \frac{1}{2}, \beta$ ) and NESTEROV( $\alpha = \frac{1}{2}, \beta$ ). **Right:** The right plot zooms on the range  $-1 \leq \beta \leq 0.25$ . POLYAK comprises a limited surface in the  $(\nu, \kappa_p, \kappa_i)$  space, leaving configurations outside this surface unexplored. Note how positive (resp. negative) values of  $\beta$  result in negative (resp. positive) values of  $\kappa_p$ , colored in red (resp. blue). Colored paths correspond to different values of  $\alpha$ . The dashed curves match between both plots.

Fig. 4 visually emphasizes the greater generality of  $\nu$ PI compared to POLYAK and NESTEROV, presented in Thm. 1. Note

<sup>2</sup>We consider a variant of the Nesterov method that uses a constant momentum coefficient.

that the  $\kappa_p$  coefficient changes between POLYAK and NESTEROV, while the  $\kappa_i$  coefficient coincides. Formally,

$$\kappa_p^{\text{POLYAK}} = -\frac{\alpha\beta}{(1-\beta)^2}, \quad \kappa_p^{\text{NESTEROV}} = -\frac{\alpha\beta^2}{(1-\beta)^2} \leq 0. \quad (9)$$

Moreover,  $\kappa_p^{\text{NESTEROV}}$  is *non-positive*, regardless of  $\beta$ . In contrast, a negative momentum value  $\beta < 0$  induces a *positive*  $\kappa_p^{\text{POLYAK}}$ . This observation is in line with the benefits of using a negative POLYAK momentum coefficient (for both players) in adversarial games presented by Gidel et al. (2019a).

### 4.3. Interpreting the updates of $\nu$ PI

Consider the execution of  $\nu$ PI ( $\nu, \kappa_p, \kappa_i$ ) and GA ( $\alpha = \kappa_i$ ) at time  $t$ .<sup>3</sup> The relative size between these updates is:

$$\frac{\Delta\nu\text{PI}}{\Delta\text{GA}} \triangleq \frac{\theta_{t+1}^{\nu\text{PI}} - \theta_t}{\theta_{t+1}^{\text{GA}} - \theta_t} = \frac{1}{1-\psi} \left[ 1 - \frac{\psi\xi_{t-1}}{e_t} \right], \quad (10)$$

where  $\psi \triangleq \frac{\kappa_p(1-\nu)}{\kappa_i + \kappa_p(1-\nu)}$ . Fig. 5 illustrates the behavior of the relative size of updates of  $\nu$ PI compared to GA. The left plot displays  $\nu$ PI with  $\kappa_p > 0$  and  $\nu = 0$ . The right plot shows the  $\nu$ PI-equivalent of POLYAK with *positive* momentum.<sup>4</sup>

Consider the colored regions present in the *left* plot of Fig. 5:

**Mode A** When  $\xi_{t-1} < e_t$ , the current violation is greater than the historical violation average (right region).  $\nu$ PI algorithm increases the multiplier *faster* than GA. When  $e_t < 0$  (left region), the primal iterate is feasible and the  $\nu$ PI algorithm agrees with GA in decreasing the multiplier, but does so *much faster* (with a factor above  $\frac{1}{1-\psi}$ ).

**Mode B** When  $e_t \in [\psi\xi_{t-1}, \xi_{t-1}]$ , the constraint violation has improved compared to the historical average but is still infeasible. In this case,  $\nu$ PI increases the multiplier *more slowly* than GA, consistent with the perceived improvement in the violation.

**Mode C** When  $e_t \in [0, \kappa\xi_{t-1}]$ , the primal iterate is still infeasible. However, the  $\nu$ PI algorithm determines that the constraint improvement is large enough to warrant a *decrease* in the multiplier. Note that in this case, GA would have continued increasing the multiplier.

In all of these cases, the  $\nu$ PI optimizer can be seen as executing proactively by considering how the current constraint violation compares to the historical estimates. This proactive behavior allows the method to *increase the multiplier faster than GA* when the constraint satisfaction is degrading, *and reduce the multiplier faster than GA* whenever sufficient improvement has been made.

<sup>3</sup>It is sufficient to consider a single scalar multiplier since the updates of both algorithms decouple across constraints/multipliers.

<sup>4</sup>The case of POLYAK with negative momentum resembles the left plot of Fig. 5 See Appx. C for further details.

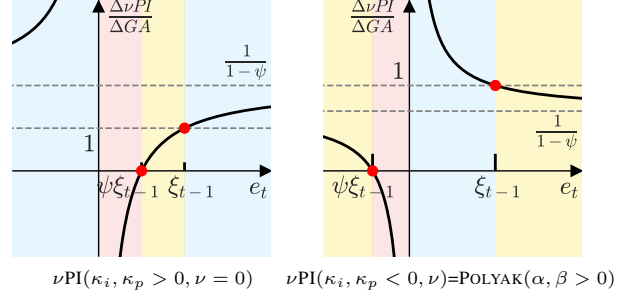


Figure 5: Comparing the update of  $\nu$ PI relative to GA.  $\nu$ PI increases the multipliers faster than GA when the constraint violation is large, enhancing convergence speed; and proactively decreases them near the feasible set, preventing overshoot. The blue, yellow, and red regions correspond to cases in which the updates performed by the  $\nu$ PI algorithm are faster, slower, or in the opposite direction than those of GA, respectively. This plot illustrates the case  $\xi_{t-1} > 0$ .

In stark contrast, Fig. 5 (right) shows a setting in which  $\kappa_i$  and  $\kappa_p$  have been chosen according to Thm. 1 for  $\beta = \nu = 0.3$ , i.e. using *positive* Polyak momentum. In this case, the algorithm would produce *stronger increases* of the multiplier whenever feasibility is improved, while *weaker increases* are executed whenever feasibility worsens. This counter-intuitive behavior may be the cause of oscillations and overshoot underlying the failure of positive momentum methods in Lagrangian games.

### 4.4. Oscillator dynamics

The continuous-time dynamics of gradient-descent/ $\nu$ PI-ascend on an equality-constrained problem can be characterized by the second-order differential equations (see Thm. 4):

$$\begin{cases} \ddot{\mathbf{x}} = -\left(\nabla^2 f + \sum_{c'} \mu_{c'} \nabla^2 h_{c'}\right) \dot{\mathbf{x}} - \mathcal{J} \mathbf{h} \dot{\boldsymbol{\mu}} & (11a) \\ \ddot{\boldsymbol{\mu}} = \kappa_i \mathcal{J} \mathbf{h}^\top \dot{\mathbf{x}} + \kappa_p \mathcal{J} \mathbf{h}^\top \ddot{\mathbf{x}} + \kappa_p \boldsymbol{\Xi}, & (11b) \end{cases}$$

where  $\boldsymbol{\Xi} = [\dot{\mathbf{x}}^\top \nabla^2 h_1 \dot{\mathbf{x}}, \dots, \dot{\mathbf{x}}^\top \nabla^2 h_c \dot{\mathbf{x}}]^\top \in \mathbb{R}^c$ .

In Appx. D we present the spectral analysis for the Lagrangian system associated with an equality-constrained quadratic program. In particular, we demonstrate how the continuous-time  $\nu$ PI algorithm can modify the eigenvalues of the system and transition between divergent, oscillatory, *critically damped* and overdamped behaviors. We show how these regime changes are controlled by the  $\kappa_p$  hyperparameter. Moreover, critical damping may require a non-zero value of  $\kappa_p$ , and is thus not achievable by GA.

### 4.5. Practical remarks

In practice, we suggest the initial condition  $\boldsymbol{\xi}_0 = \mathbf{e}_0$ , as it ensures that the first step of  $\nu$ PI matches that of gradient

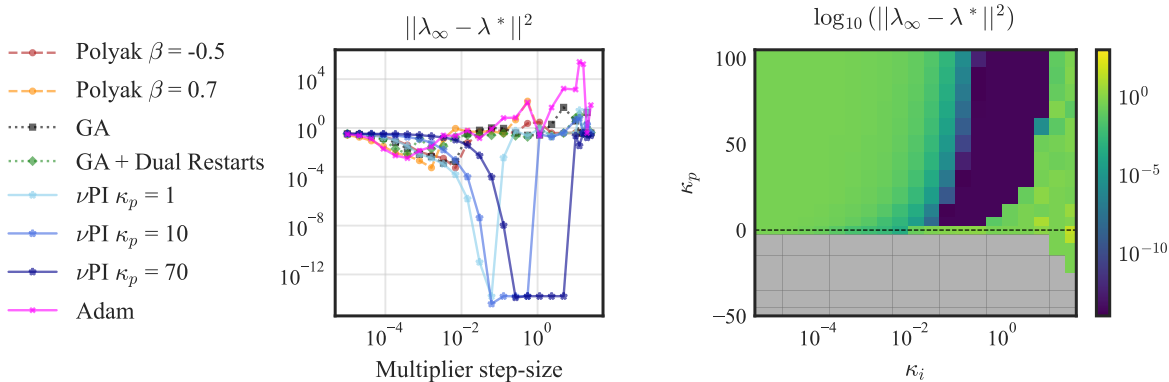


Figure 6: SVM experiment (§5.1). **Left:** Step-size sensitivity plot. **Right:**  $\kappa_p$  and  $\kappa_i$  sensitivity plot for  $\nu$ PI. **Without a  $\kappa_p$  of at least one, none of the methods converge to the optimal dual variable. Higher  $\kappa_p$  values allow for choosing higher and broader range of  $\kappa_i$ 's.** The x-axis of the left plot represents  $\kappa_i$  for the  $\nu$ PI parameter and  $\alpha$ , the step-size, for the other optimizers. In the right plot, the gray color shows the runs exceeding a distance of  $10^3$  to  $\lambda^*$ .

ascent. In cases where the constraints can be evaluated without noise, we suggest a default value of  $\nu = 0$ . This leaves only the additional hyperparameter  $\kappa_p$  to be tuned (besides the “step-size”  $\kappa_i$ ). We highlight that the main benefits of the  $\nu$ PI algorithm remain even when  $\nu = 0$ . However,  $\nu$  can be useful for filtering noise in the constraint measurement, as shown in our fairness experiments in §5.2.

There is a predictable monotonic behavior of the damping of the system as the  $\kappa_p$  coefficient increases. This is illustrated in Fig. 9 in §5.3 for a sparsity task. As a side effect, higher values of  $\kappa_p$  make the tuning of the  $\kappa_i$  coefficient easier, as seen in Fig. 6 in §5.1. As a heuristic to tune  $\nu$ PI, we suggest considering a large initial  $\kappa_p$  value (so that its influence on the optimization dynamics is significant), and then try a grid of  $\kappa_i$  values. A good starting place is a grid of  $\kappa_i$  values around a suitable step-size for gradient ascent.

## 5. Experiments

In this section, we present an empirical comparison between  $\nu$ PI and a series of baseline optimization methods popular for minimization. We consider gradient ascent, gradient ascent with positive (Polyak, 1964; Nesterov, 1983) and negative (Gidel et al., 2019a) momentum, and ADAM (Kingma & Ba, 2014). The goal of our experiments is to highlight the flexibility of  $\nu$ PI and its ability to mitigate oscillations and overshoot when used to optimize Lagrange multipliers.

Our implementations use PyTorch (Paszke et al., 2019) and the Cooper library for Lagrangian constrained optimization (Gallego-Posada & Ramirez, 2022).

### 5.1. Hard-margin SVMs

We consider solving a *hard-margin* linear SVM problem via its associated Lagrangian formulation. While specialized QP solvers exist to find solutions for this task, we consider

the Lagrangian formulation in order to illustrate the dynamics of the multipliers in a simple machine learning task. These experiments show how standard methods for minimization produce oscillations on the multipliers, which have detrimental effects on convergence. Consider

$$\min_{\mathbf{w}} \|\mathbf{w}\|^2/2 \quad \text{s.t.} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad \text{for } i \in [m], \quad (12)$$

where  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$  are labeled training datapoints, and  $\mathbf{w}$  and  $b$  are the parameters of the SVM classifier.

We perform binary classification on two linearly separable classes from the Iris dataset (Fisher, 1988). We apply alternating GDA updates on the Lagrangian associated with Eq. (12), with a fixed optimizer for the primal variables. For details on our SVM experiments, see Appx. F.1.

**Multiplier dynamics.** Fig. 1 shows the oscillations on the multiplier in all of the baselines. In these tasks, all of the methods that do not diverge achieve perfect training and validation accuracy. However, among the methods we experiment with, the only method capable of achieving zero constraint violation is the  $\nu$ PI algorithm. In contrast to all baselines,  $\nu$ PI dampens the oscillations and converges to the optimal multiplier value.

**Sensitivity analysis.** Fig. 6 (left) illustrates the robustness of  $\nu$ PI to the choice of  $\kappa_i$ . The considered baselines fail to converge to the ground truth multiplier value, across a wide range of step-sizes. For these baselines, small step-size choices avoid divergence but do not lead to recovering the optimal Lagrange multipliers, while large step-sizes increase the oscillations. In contrast, introducing a  $\kappa_p$  term of more than 1 results in convergence for some step-sizes within the selected range (see the  $\nu$ PI curves). Moreover, increasing  $\kappa_p$  to a higher value broadens the range of step-sizes that lead to convergence, and enables the use of bigger step-size values that converge. This behavior can be observed more extensively in the heatmap of Fig. 6 (right).

## 5.2. Fairness

We consider a classification task under *statistical parity* constraints, as described in Cotter et al. (2019). This leads to the following constrained optimization problem:

$$\min_{\mathbf{w}} L(\mathbf{w}) \quad \text{s.t.} \quad \mathbb{P}(\hat{y} = 1 | g) = \mathbb{P}(\hat{y} = 1), \quad \forall g \in G \quad (13)$$

where  $L(\mathbf{w})$  is the loss of model  $\mathbf{w}$ ,  $\hat{y}$  is the model prediction, and  $G$  represents the set of protected groups in the dataset. The constraints require the probability of positive prediction to be equal across all groups.

**Model and data.** We train binary classifiers on the Adult dataset (Becker & Kohavi, 1996). Groups correspond to the intersection of race (2 values) and gender (5 values), leading to 10 constraints. We use an MLP with two 100-dimensional hidden layers. Our experimental setup is similar to those of Zafar et al. (2019) and Cotter et al. (2019). However, in our setting, non-convexity precludes the use of specialized solvers (as done by Zafar et al. (2019)) and requires iterative optimization approaches.

**Optimization configuration.** We train the model using ADAM ( $\alpha = 10^{-2}$ ) with a batch size of 512. To mitigate the noise in the estimation of the constraint satisfaction, we update the multipliers once every epoch, using the exact constraint measurement over the entire training set.

**Results.** Figure 7 includes training curves for experiments with GA and  $\nu$ PI applied to the Lagrange multipliers. We report two of the multipliers, the model accuracy, and the maximum constraint violation (in absolute value).

For this task, gradient ascent is a strong baseline as it successfully reduces the violation of the constraints. Both GA and  $\nu$ PI ( $\nu = 0.99$ ) significantly improve compared to an *unconstrained* baseline which achieves a maximum violation of 20% (not shown in Fig. 7 for readability).

$\nu$ PI ( $\nu = 0$ ) runs exhibit unstable multiplier dynamics as the noise of the constraints is amplified by the proportional term. During our experiments, we observed that when  $\nu = 0$ , larger  $\kappa_p$  values lead to noisier multipliers and unstable optimization. In contrast,  $\nu$ PI ( $\nu = 0.99$ ) **reduces the maximum violation faster and achieves better training accuracy (92.4% vs 89%)**.

All experiments reach a final maximum violation of around 1.7%. We hypothesize that it is not possible to decrease this value further (while carrying out stochastic updates on the primal variables) since the constraint gradients may be misaligned across mini-batches.

**Multiplier dynamics.** As can be seen in the evolution of multipliers 2 and 7 shown in Fig. 7,  $\nu$ PI yields multipliers that stabilize at their limiting values faster than those of GA.

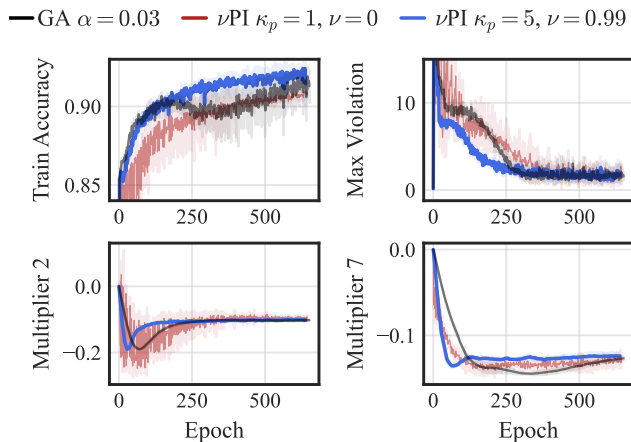


Figure 7: Dynamics of  $\nu$ PI compared to GA for the fairness task in Eq. (13).  $\nu$ PI has faster convergence in the multiplier value and achieves a better training accuracy than GA. All dual optimizers use a step-size ( $\kappa_i$  for  $\nu$ PI) of 0.03. Results are aggregated across five seeds.

## 5.3. Sparsity

We consider the problem of learning models under inequality  $L_0$ -sparsity constraints (Louizos et al., 2018; Gallego-Posada et al., 2022). See Appx. F.2 for further background.

$$\min_{\mathbf{w}, \phi \in \mathbb{R}^d} \mathbb{E}_{\mathbf{z}|\phi} [L(\mathbf{w} \odot \mathbf{z} | \mathcal{D})] \quad \text{s.t.} \quad \frac{\mathbb{E}_{\mathbf{z}|\phi} [\|\mathbf{z}\|_0]}{\#(\mathbf{w})} \leq \epsilon \quad (14)$$

When using GA updates for the multipliers, Gallego-Posada et al. (2022) observe a tendency of the model to “overshoot” into the feasible region and become significantly less dense than the prescribed level. Since a reduction in model density corresponds to a reduction in capacity, this overshoot may have a detrimental effect on the performance of the model.

Our experiments explore the effect of  $\nu$ PI on the sparsity-constrained task, and compare it with dual restarts (Gallego-Posada et al., 2022, §3). Our results show that  $\nu$ PI allows for fine-grained control over overshoot, thus enabling the sparse model to retain as much performance as possible.

**Experiment configuration and hyperparameters.** We consider classifying CIFAR-10 (Krizhevsky, 2009) images with ResNet-18 (He et al., 2016) models. To highlight the ease-of-use of  $\nu$ PI, our setup remains as close as possible to Gallego-Posada et al. (2022): we apply output channel sparsity on the first layer of each residual block in the model, and re-use the authors’ choice of optimizer and step-size for  $\phi$ . Our sparsity experiments consider  $\nu = 0$ .

**Global and layer-wise settings.** We present sparsity experiments with either ① one global constraint on the sparsity of the entire model, or ② multiple constraints, each prescribing a maximum density per layer.



The metrics reported in this section are aggregated across 5 seeds. Experimental details for this task can be found in Appx. F.2. For comprehensive experimental results across multiple sparsity levels, and ablations on the use of momentum and ADAM for updating the multipliers, see Appx. G.1.

**Results.** Fig. 8 shows how gradient ascent and positive and negative momentum values consistently yield runs that overshoot into becoming overly sparse. The extra reduction in capacity results in a loss in performance. In contrast,  $\nu$ PI consistently recovers feasible solutions, with minimal overshoot. While dual restarts do not incur in overshoot, they produce slightly infeasible solutions.

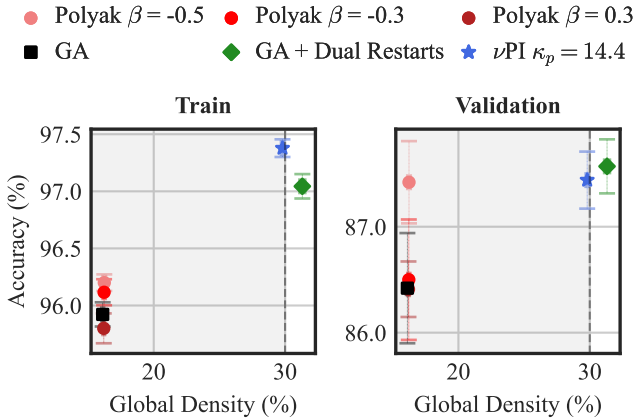


Figure 8: CIFAR10 trade-off plot for *global* sparsity under a 30% density target.  $\nu$ PI successfully achieves the desired sparsity while achieving the highest train accuracy. The shaded region is the feasible set. As higher density correlates to higher train accuracy, overshooting to a lower density is undesirable. All optimizers use the same step-size.

Fig. 9 consists on an ablation on the  $\kappa_p$  value. We observe that by increasing the hyper-parameter, overshoot is reduced, eventually turning into undershoot (which leads to infeasible solutions). Since the density of the model is monotonically tied to the choice of  $\kappa_p$ , tuning  $\nu$ PI for this task can be done via bisection search, without the need to consider a grid (which is usually required for tuning the step-size).

Table 1 shows sparsity experiments with layer-wise sparsity targets. Gradient ascent and momentum methods overshoot and the degree of overshoot differs significantly across layers. In contrast, GA with dual restarts and  $\nu$ PI mitigate overshoot and produce constraints spanning a narrow range of values. This highlights the robustness of  $\nu$ PI as the  $\kappa_p$  coefficient did not need to be tuned separately per constraint.

**Multiplier dynamics.** Figure 2 shows the training dynamics for a global sparsity constraint and its multiplier under a 30% density target. We observe that GA and POLYAK quickly lead to overshoot into the feasible set, but manage to regain some model capacity as training progresses. GA with dual

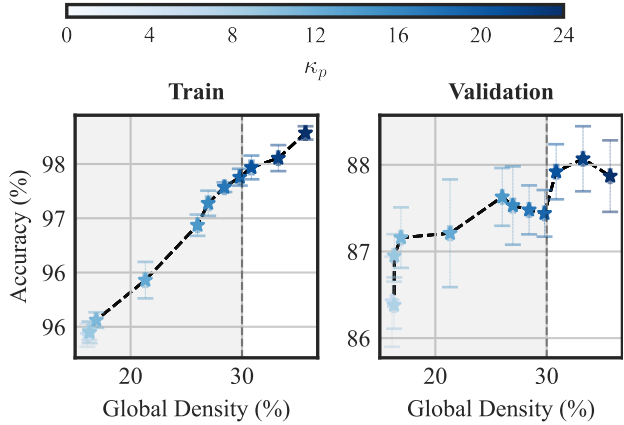


Figure 9: Ablation of  $\kappa_p$  values for  $\nu$ PI on CIFAR10. An increasing  $\kappa_p$  leads to more damping and less overshoot. Target density is 30%. The shaded region is the feasible set.

restarts sets the value of the multiplier to zero as soon as feasibility is achieved, thus preventing an incursion of the constraint into the feasible set.  $\nu$ PI produces well-behaved multipliers and successfully avoids constraint overshoot.

Table 1: CIFAR10 results for *layer-wise* sparsity under a 30% density target. GA and momentum methods overshoot to *different* values for each constraint.  $\nu$ PI achieves the desired sparsity on all layers while achieving the highest train accuracy. All dual optimizers use the same step-size.

Method	Accuracy		Violation		
	Train	Test	Min	Max	Range
Polyak $\beta = -0.5$	91.9	83.6	-26.5	-7.9	18.9
Polyak $\beta = -0.3$	92.1	83.4	-27.1	-6.7	20.6
Polyak $\beta = 0.3$	91.9	82.5	-26.3	-2.3	24.0
GA	92.0	84.1	-27.8	-5.2	22.0
GA + Dual Restarts	95.0	85.3	-0.0	1.2	1.2
Ours - $\nu$ PI $\kappa_p = 8.0$	95.1	86.2	-1.7	0.1	1.8

## 6. Conclusion

Previous work has highlighted that employing PID controllers on the multipliers in Lagrangian constrained optimization problems reduces oscillation and overshoot. In this paper, we consider  $\nu$ PI, a variant of a PI controller that generalizes various popular methods for optimization. We complement previous work by providing insights justifying why PI controllers are desirable for Lagrangian optimization. Moreover, we highlight some intuitions as to why momentum methods fail in this context. While we focus our efforts on constrained optimization, our results indicate that  $\nu$ PI may improve the dynamics of linear players in general min-max games. Investigating the behavior of  $\nu$ PI on non-linear players is left as a direction of future work.

## Impact Statement

Constrained optimization offers tools for reliably enforcing properties on machine learning models. It is, therefore, applicable for enhancing safety, robustness, and fairness in AI models. By integrating constraints into the model development process, rather than retrofitting safety measures as afterthoughts, we advocate for a paradigm shift towards building models that are inherently secure “by design.” We intend our fairness experiments as a conceptual illustration of the potential for positive impact of constrained approaches in the development of machine learning models.

Our paper presents insights into the robustness of algorithms for constrained optimization, and highlights  $\nu$ PI as a reliable tool for training models with constraints. Thus, our work lays the groundwork for practitioners to adopt and implement constrained approaches confidently in diverse real-world applications.

## Acknowledgements

This research was partially supported by an IVADO PhD Excellence Scholarship, the Canada CIFAR AI Chair program (Mila), the NSERC Discovery Grant RGPIN2017-06936 and by Samsung Electronics Co., Ltd. Simon Lacoste-Julien is a CIFAR Associate Fellow in the Learning in Machines & Brains program.

This research was enabled in part by compute resources, software, and technical help provided by Mila.

We would like to thank Ioannis Mitliagkas for useful discussions during the development of this work.

## Reproducibility Statement

We provide our code,<sup>5</sup> including scripts to replicate the experiments in this paper. §4.5 presents some considerations when using the  $\nu$ PI algorithm in practice. Experimental details, as well as the hyper-parameters used in our experiments, are included in Appx. F. Our implementations use the open-source libraries PyTorch (Paszke et al., 2019) and Cooper (Gallego-Posada & Ramirez, 2022).

## References

An, W., Wang, H., Sun, Q., Xu, J., Dai, Q., and Zhang, L. A PID Controller Approach for Stochastic Optimization of Deep Networks. In *CVPR*, 2018. (Cit. on p. 2, 5, 14)

Arrow, K., Hurwicz, L., and Uzawa, H. *Studies in Linear and Non-linear Programming*. Stanford University Press, 1958. (Cit. on p. 2, 3)

<sup>5</sup><https://github.com/motahareh-sohrabi/nuPI>

Åström, K. and Hägglund, T. *PID Controllers*. International Society for Measurement and Control, 1995. (Cit. on p. 2, 5)

Becker, B. and Kohavi, R. Adult. UCI Machine Learning Repository, 1996. DOI: <https://doi.org/10.24432/C5XW20>. (Cit. on p. 8, 24)

Bertsekas, D. On the Method of Multipliers for Convex Programming. *IEEE transactions on automatic control*, 1975. (Cit. on p. 2)

Bertsekas, D. *Nonlinear Programming*. Athena Scientific, 2016. (Cit. on p. 1, 2)

Bertsekas, D. P. On the Goldstein-Levitin-Polyak Gradient Projection Method. *IEEE Transactions on automatic control*, 1976. (Cit. on p. 2)

Bottou, L., Curtis, F. E., and Nocedal, J. Optimization Methods for Large-Scale Machine Learning. *SIAM Review*, 60(2):223–311, 2018. (Cit. on p. 1)

Boyd, S. Stanford ENGR108: Introduction to applied linear algebra: 2020: Lecture 53-VMLS CSTRD nonlinear LS, 2021. URL [https://youtu.be/SM\\_ZieyKicU?si=PWNMr7vxMQkhFBbf&t=815](https://youtu.be/SM_ZieyKicU?si=PWNMr7vxMQkhFBbf&t=815). (Cit. on p. 21)

Boyd, S. and Vandenberghe, L. *Convex Optimization*. Cambridge University Press, 2004. (Cit. on p. 1)

Casti, U., Bastianello, N., Carli, R., and Zampieri, S. A Control Theoretical Approach to Online Constrained Optimization. *arXiv preprint arXiv:2309.15498*, 2023. (Cit. on p. 2, 14)

Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., and Huq, A. Algorithmic decision making and the cost of fairness. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017. (Cit. on p. 24)

Cotter, A., Jiang, H., Gupta, M. R., Wang, S., Narayan, T., You, S., and Sridharan, K. Optimization with Non-Differentiable Constraints with Applications to Fairness, Recall, Churn, and Other Goals. *JMLR*, 2019. (Cit. on p. 1, 8, 24)

Dahl, G. E., Schneider, F., Nado, Z., Agarwal, N., Sastry, C. S., Hennig, P., Medapati, S., Eschenhagen, R., Kasimbeg, P., Suo, D., Bae, J., Gilmer, J., Peirson, A. L., Khan, B., Anil, R., Rabbat, M., Krishnan, S., Snider, D., Amid, E., Chen, K., Maddison, C. J., Vasudev, R., Badura, M., Garg, A., and Mattson, P. Benchmarking Neural Network Training Algorithms. *arXiv:2306.07179*, 2023. (Cit. on p. 2, 3)

- Dikin, I. I. Iterative Solution of Problems of Linear and Quadratic Programming. In *Doklady Akademii Nauk*. Russian Academy of Sciences, 1967. (Cit. on p. 2)
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness Through Awareness. In *Innovations in Theoretical Computer Science*, 2012. (Cit. on p. 24)
- Elenter, J., NaderiAlizadeh, N., and Ribeiro, A. A Lagrangian Duality Approach to Active Learning. In *NeurIPS*, 2022. (Cit. on p. 1)
- Farahmand, A.-M. and Ghavamzadeh, M. PID Accelerated Value Iteration Algorithm. In *ICML*, 2021. (Cit. on p. 1)
- Fioretto, F., Van Hentenryck, P., Mak, T. W. K., Tran, C., Baldo, F., and Lombardi, M. Lagrangian Duality for Constrained Deep Learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2020. (Cit. on p. 1)
- Fisher, R. A. Iris. UCI Machine Learning Repository, 1988. DOI: <https://doi.org/10.24432/C56C76>. (Cit. on p. 7, 22)
- Frank, M. and Wolfe, P. An Algorithm for Quadratic Programming. *Naval Research Logistics Quarterly*, 1956. (Cit. on p. 2)
- Gale, T., Elsen, E., and Hooker, S. The State of Sparsity in Deep Neural Networks. *arXiv:1902.09574*, 2019. (Cit. on p. 23)
- Gallego-Posada, J. and Ramirez, J. Cooper: a toolkit for Lagrangian-based constrained optimization. <https://github.com/cooper-org/cooper>, 2022. (Cit. on p. 1, 7, 10, 22)
- Gallego-Posada, J., Ramirez, J., Erraqabi, A., Bengio, Y., and Lacoste-Julien, S. Controlled Sparsity via Constrained Optimization or: *How I Learned to Stop Tuning Penalties and Love Constraints*. In *NeurIPS*, 2022. (Cit. on p. 1, 2, 4, 8, 23, 24)
- Gidel, G., Askari, R., Pezeshki, M., LePriol, R., Huang, G., Lacoste-Julien, S., and Mitliagkas, I. Negative Momentum for Improved Game Dynamics. In *AISTATS*, 2019a. (Cit. on p. 2, 6, 7, 19, 20, 23)
- Gidel, G., Berard, H., Vignoud, G., Vincent, P., and Lacoste-Julien, S. A Variational Inequality Perspective on Generative Adversarial Networks. In *ICLR*, 2019b. (Cit. on p. 2, 3)
- Hashemizadeh, M., Ramirez, J., Sukumaran, R., Farnadi, G., Lacoste-Julien, S., and Gallego-Posada, J. Balancing Act: Constraining Disparate Impact in Sparse Models. In *ICLR*, 2024. (Cit. on p. 1)
- He, K., Zhang, X., Ren, S., and Sun, J. Deep Residual Learning for Image Recognition. In *CVPR*, 2016. (Cit. on p. 8, 23)
- Hounie, I., Elenter, J., and Ribeiro, A. Neural Networks with Quantization Constraints. In *ICASSP*, 2023. (Cit. on p. 1)
- Hu, B. and Lessard, L. Control Interpretations for First-Order Optimization Methods. In *American Control Conference*, 2017. (Cit. on p. 2, 14)
- Kingma, D. and Ba, J. Adam: A Method for Stochastic Optimization. In *ICLR*, 2014. (Cit. on p. 7, 23)
- Korpelevich, G. M. The extragradient method for finding saddle points and other problems. *Matecon*, 1976. (Cit. on p. 2)
- Krizhevsky, A. Learning Multiple Layers of Features from Tiny Images. Technical report, University of Toronto, Toronto, Ontario, 2009. (Cit. on p. 8)
- Lessard, L., Recht, B., and Packard, A. Analysis and Design of Optimization Algorithms via Integral Quadratic Constraints. *SIAM Journal on Optimization*, 26(1):57–95, 2016. (Cit. on p. 5)
- Lin, T., Jin, C., and Jordan, M. On Gradient Descent Ascent for Nonconvex-Concave Minimax Problems. In *ICML*, 2020. (Cit. on p. 2)
- Loshchilov, I. and Hutter, F. SGDR: Stochastic Gradient Descent with Warm Restarts. In *ICLR*, 2017. (Cit. on p. 23)
- Louizos, C., Welling, M., and Kingma, D. P. Learning Sparse Neural Networks through  $L_0$  Regularization. In *ICLR*, 2018. (Cit. on p. 8, 23, 28)
- Maddison, C. J., Mnih, A., and Teh, Y. W. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. In *ICLR*, 2017. (Cit. on p. 23)
- Mokhtari, A., Ozdaglar, A. E., and Pattathil, S. Convergence Rate of  $O(1/k)$  for Optimistic Gradient and Extragradient Methods in Smooth Convex-Concave Saddle Point Problems. *SIAM Journal on Optimization*, 2020a. (Cit. on p. 2)
- Mokhtari, A., Ozdaglar, A. E., and Pattathil, S. A Unified Analysis of Extra-gradient and Optimistic Gradient Methods for Saddle Point Problems: Proximal Point Approach. In *AISTATS*, 2020b. (Cit. on p. 5)
- Nesterov, Y. E. A method of solving a convex programming problem with convergence rate  $o(1/k^2)$ . In *Doklady Akademii Nauk*, volume 269, pp. 543–547. Russian Academy of Sciences, 1983. (Cit. on p. 5, 7)

- Nocedal, J. and Wright, S. J. *Numerical Optimization*. Springer, 2006. (Cit. on p. 1, 2)
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS*, 2019. (Cit. on p. 7, 10, 22)
- Platt, J. C. and Barr, A. H. Constrained Differential Optimization. In *NeurIPS*, 1987. (Cit. on p. 2)
- Polyak, B. T. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964. (Cit. on p. 5, 7)
- Popov, L. D. A modification of the Arrow-Hurwicz method for search of saddle points. *Mathematical notes of the Academy of Sciences of the USSR*, 28:845–848, 1980. (Cit. on p. 2, 5)
- Recht, B. The Best Things in Life Are Model Free. arg min, blog: <https://archives.argmin.net/2018/04/19/pid/>, 2018. (Cit. on p. 2)
- Shen, L., Chen, C., Zou, F., Jie, Z., Sun, J., and Liu, W. A Unified Analysis of AdaGrad with Weighted Aggregation and Momentum Acceleration. In *IEEE TNNLS*, 2018. (Cit. on p. 5)
- Stooke, A., Achiam, J., and Abbeel, P. Responsive Safety in Reinforcement Learning by PID Lagrangian Methods. In *ICML*, 2020. (Cit. on p. 1, 2, 4, 5, 14, 18)
- Zafar, M. B., Valera, I., Gomez-Rodriguez, M., and Gum-madi, K. P. Fairness Constraints: A Flexible Approach for Fair Classification. *JMLR*, 2019. (Cit. on p. 1, 8, 24)
- Zhang, G. and Wang, Y. On the Suboptimality of Negative Momentum for Minimax Optimization. In *AISTATS*, 2021. (Cit. on p. 2)
- Zhang, G., Wang, Y., Lessard, L., and Grosse, R. B. Near-optimal Local Convergence of Alternating Gradient Descent-Ascent for Minimax Optimization. In *AISTATS*, 2022. (Cit. on p. 2, 3)
- Zoutendijk, G. *Methods of Feasible Directions: A Study in Linear and Non-linear Programming*. Elsevier Publishing Company, 1960. (Cit. on p. 2)

# Appendix

## Table of Contents

---

<b>A</b>	<b>Further discussion on prior works using PID controls in optimization</b>	<b>14</b>
<b>B</b>	<b>Connections between <math>\nu</math>PI and momentum methods</b>	<b>14</b>
<b>C</b>	<b>Interpreting the updates of <math>\nu</math>PI</b>	<b>16</b>
<b>D</b>	<b>Analysis of continuous-time <math>\nu</math>PI dynamics as an oscillator</b>	<b>18</b>
	D.1 Oscillator dynamics of GD/ $\nu$ PI flow . . . . .	18
	D.2 Dynamics of GD/ $\nu$ PI flow for a constrained quadratic program . . . . .	19
<b>E</b>	<b>Illustrative 2D nonconvex problem</b>	<b>21</b>
<b>F</b>	<b>Experimental details</b>	<b>22</b>
	F.1 Linear SVM experiments . . . . .	22
	F.2 Sparsity experiments . . . . .	23
	F.3 Fairness experiments . . . . .	24
<b>G</b>	<b>Comprehensive results on the sparsity task</b>	<b>24</b>
	G.1 Global . . . . .	24
	G.2 Layer-wise . . . . .	24
<b>H</b>	<b>Additional Experiments</b>	<b>27</b>
	H.1 Dynamics . . . . .	27
	H.2 Ablation on the value of $\kappa_p$ . . . . .	29
	H.3 ADAM . . . . .	30
	H.4 Momentum . . . . .	31

---

## A. Further discussion on prior works using PID controls in optimization

- In [Stooke et al. \(2020\)](#), the authors focus almost exclusively on applying PID control to constrained reinforcement learning. The authors do not explore the optimization aspects of PID-based updates for Lagrange multipliers, which are the main focus of our work. Our key theoretical contribution (Thm. 1) shows that  $\nu$ PI provides a generalization of momentum-based optimization techniques. Note that the controller considered by [Stooke et al. \(2020\)](#) is unable to generalize momentum methods. Thanks to the unifying framework provided by Thm. 1, we provide insights to understand why momentum fails at Lagrangian optimization tasks (Fig. 5). Moreover, our experiments encompass linear SVMs, sparsity, and fairness tasks, and are not restricted to reinforcement learning.
- [An et al. \(2018\)](#) propose directly updating the parameters of a neural network using a PID controller (for unconstrained minimization only). Their approach has not been widely adopted by the deep learning community, possibly due to the highly specialized training procedures that have been developed for training neural networks. Although connected due to their use of PID control, this paper is not directly relevant to our work as we limit our scope to not modifying the optimization protocol for the (primal) model parameters.
- The work of [Casti et al. \(2023\)](#) focuses on the theoretical aspects of using PID control for problems with linear constraints. Their analysis is not directly applicable to our setting since we are interested in general machine learning applications involving nonconvex constraints.
- [Hu & Lessard \(2017\)](#) present control interpretations of first-order optimization methods and show how worst-case convergence rates of optimization algorithms can be derived from a control theoretical perspective. The idea of examining a possible connection between our  $\nu$ PI algorithm and other momentum methods was inspired by this work.

## B. Connections between $\nu$ PI and momentum methods

Table 2: Classical optimization methods as instances of  $\nu$ PI ( $\nu, \kappa_p, \kappa_i; \xi_0$ ).

Algorithm	$\xi_0$	$\kappa_p$	$\kappa_i$	$\nu$
UNIFIEDMOMENTUM( $\alpha, \beta, \gamma$ )	$(1 - \beta)\mathbf{e}_0$	$-\frac{\alpha\beta}{(1 - \beta)^2} [1 - \gamma(1 - \beta)]$	$\frac{\alpha}{1 - \beta}$	$\beta$
POLYAK( $\alpha, \beta$ )	$(1 - \beta)\mathbf{e}_0$	$-\frac{\alpha\beta}{(1 - \beta)^2}$	$\frac{\alpha}{1 - \beta}$	$\beta$
NESTEROV( $\alpha, \beta$ )	$(1 - \beta)\mathbf{e}_0$	$-\frac{\alpha\beta^2}{(1 - \beta)^2}$	$\frac{\alpha}{1 - \beta}$	$\beta$
PI	$\mathbf{e}_0$	$\kappa_p$	$\kappa_i$	0
OPTIMISTICGRADIENTASCENT( $\alpha$ )	$\mathbf{e}_0$	$\alpha$	$\alpha$	0
$\nu$ PI( $\nu, \kappa_p, \kappa_i$ ) in practice	$\mathbf{0}$	$\kappa_i$	$\kappa_p$	$\nu$
GRADIENTASCENT( $\alpha$ )	–	0	$\alpha$	0

**Lemma 2.** *The  $\nu$ PI ( $\nu, \kappa_p, \kappa_i; \xi_0$ ) algorithm can be equivalently expressed as the recursion:*

$$\boldsymbol{\theta}_1 = \boldsymbol{\theta}_0 + \kappa_i \mathbf{e}_0 + \kappa_p \boldsymbol{\xi}_0 \quad (15a)$$

$$\boldsymbol{\xi}_t = \nu \boldsymbol{\xi}_{t-1} + (1 - \nu) \mathbf{e}_t \quad (15b)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \kappa_i \mathbf{e}_t + \kappa_p (1 - \nu) (\mathbf{e}_t - \boldsymbol{\xi}_{t-1}) \text{ for } t \geq 1 \quad (15c)$$

*Proof of Lemma 2.* For  $t \geq 1$ , we have:

$$\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \stackrel{(\nu\text{PI3})}{=} \kappa_p \boldsymbol{\xi}_t + \kappa_i \sum_{\tau=0}^t \mathbf{e}_\tau - \kappa_p \boldsymbol{\xi}_{t-1} - \kappa_i \sum_{\tau=0}^{t-1} \mathbf{e}_\tau = \kappa_i \mathbf{e}_t + \kappa_p (\boldsymbol{\xi}_t - \boldsymbol{\xi}_{t-1}) \stackrel{(\nu\text{PI2})}{=} \kappa_i \mathbf{e}_t + \kappa_p (1 - \nu) (\mathbf{e}_t - \boldsymbol{\xi}_{t-1}) \quad (16)$$

□

**Lemma 3.** The UNIFIEDMOMENTUM( $\alpha, \beta, \gamma; \phi_0 = \mathbf{0}$ ) algorithm can be expressed as the single-parameter recurrence:

$$\boldsymbol{\theta}_1 = \boldsymbol{\theta}_0 + \alpha(1 + \beta\gamma)\mathbf{e}_0 \quad (17a)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha\mathbf{e}_t + \beta(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}) + \alpha\beta\gamma(\mathbf{e}_t - \mathbf{e}_{t-1}) \text{ for } t \geq 1. \quad (17b)$$

**Proof of Lemma 3.**

$$\phi_1 \stackrel{(UM_2)}{=} \cancel{\beta\phi_0} + \alpha\mathbf{e}_0 = \alpha\mathbf{e}_0 \quad (18a)$$

$$\boldsymbol{\theta}_1 \stackrel{(UM_3)}{=} \boldsymbol{\theta}_0 + \phi_1 + \beta\gamma(\phi_1 - \cancel{\phi_0}) = \boldsymbol{\theta}_0 + \alpha(1 + \beta\gamma)\mathbf{e}_0. \quad (18b)$$

$$\boldsymbol{\theta}_{t+1} \stackrel{(UM_3)}{=} \boldsymbol{\theta}_t + \phi_{t+1} + \beta\gamma(\phi_{t+1} - \phi_t) \quad (19a)$$

$$\stackrel{(UM_2)}{=} \boldsymbol{\theta}_t + \beta\phi_t + \alpha\mathbf{e}_t + \beta\gamma(\beta\phi_t + \alpha\mathbf{e}_t - \beta\phi_{t-1} - \alpha\mathbf{e}_{t-1}) \quad (19b)$$

$$= \boldsymbol{\theta}_t + \beta\phi_t + \alpha\mathbf{e}_t + \beta\gamma(\beta(\phi_t - \phi_{t-1}) + \alpha(\mathbf{e}_t - \mathbf{e}_{t-1})) \quad (19c)$$

$$= \boldsymbol{\theta}_t + \alpha\mathbf{e}_t + \beta[\phi_t + \gamma\beta(\phi_t - \phi_{t-1})] + \alpha\beta\gamma(\mathbf{e}_t - \mathbf{e}_{t-1}) \quad (19d)$$

$$\stackrel{(UM_3)}{=} \boldsymbol{\theta}_t + \alpha\mathbf{e}_t + \beta(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}) + \alpha\beta\gamma(\mathbf{e}_t - \mathbf{e}_{t-1}). \quad (19e)$$

□

**Theorem. 1** Under the same initialization  $\boldsymbol{\theta}_0$ , UNIFIEDMOMENTUM( $\alpha, \beta \neq 1, \gamma; \phi_0 = \mathbf{0}$ ) is a special case of  $\nu$ PI ( $\nu, \kappa_p, \kappa_i; \boldsymbol{\xi}_0$ ) with the following hyperparameter choices:

$$\nu = \beta \quad \kappa_p = -\frac{\alpha\beta}{(1-\beta)^2} [1 - \gamma(1-\beta)] \quad \kappa_i = \frac{\alpha}{1-\beta} \quad \boldsymbol{\xi}_0 = (1-\beta)\mathbf{e}_0 \quad (20)$$

**Proof of Thm. 1.** We want to find values of  $\nu, \kappa_p, \kappa_i$  and  $\boldsymbol{\xi}_0$  such that the sequence of iterates produced by  $\nu$ PI ( $\nu, \kappa_p, \kappa_i; \boldsymbol{\xi}_0$ ) satisfies Eq. (17b). For  $t \geq 2$  we have:

$$\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \stackrel{(17b)}{=} \alpha\mathbf{e}_t + \beta(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}) + \alpha\beta\gamma(\mathbf{e}_t - \mathbf{e}_{t-1}) \quad (21a)$$

$$\kappa_i\mathbf{e}_t + \kappa_p(\boldsymbol{\xi}_t - \boldsymbol{\xi}_{t-1}) \stackrel{(15c)}{=} \alpha\mathbf{e}_t + \beta(\kappa_i\mathbf{e}_{t-1} + \kappa_p(\boldsymbol{\xi}_{t-1} - \boldsymbol{\xi}_{t-2})) + \alpha\beta\gamma(\mathbf{e}_t - \mathbf{e}_{t-1}) \quad (21b)$$

$$\mathbf{e}_t(\kappa_i - \alpha - \alpha\beta\gamma) + \mathbf{e}_{t-1}(-\beta\kappa_i + \alpha\beta\gamma) + \kappa_p[\boldsymbol{\xi}_t - \boldsymbol{\xi}_{t-1} - \beta(\boldsymbol{\xi}_{t-1} - \boldsymbol{\xi}_{t-2})] = 0 \quad (22)$$

Several applications of the definition of  $\boldsymbol{\xi}_t$  (line 2 in Algo. 1) give:

$$\boldsymbol{\xi}_t - \boldsymbol{\xi}_{t-1} - \beta(\boldsymbol{\xi}_{t-1} - \boldsymbol{\xi}_{t-2}) = (1-\nu)[\mathbf{e}_t - \boldsymbol{\xi}_{t-1}] - \beta\boldsymbol{\xi}_{t-1} + \beta\boldsymbol{\xi}_{t-2} \quad (23)$$

$$= (1-\nu)\mathbf{e}_t - (1+\beta-\nu)\boldsymbol{\xi}_{t-1} + \beta\boldsymbol{\xi}_{t-2} \quad (24)$$

$$= (1-\nu)\mathbf{e}_t - (1+\beta-\nu)[\nu\boldsymbol{\xi}_{t-2} + (1-\nu)\mathbf{e}_{t-1}] + \beta\boldsymbol{\xi}_{t-2} \quad (25)$$

$$= (1-\nu)\mathbf{e}_t - (1-\nu)(1+\beta-\nu)\mathbf{e}_{t-1} + (1-\nu)(\beta-\nu)\boldsymbol{\xi}_{t-2} \quad (26)$$

Thus we can re-arrange to get:

$$[\mathbf{e}_t \quad \mathbf{e}_{t-1} \quad \boldsymbol{\xi}_{t-2}] \begin{bmatrix} \kappa_i + (1-\nu)\kappa_p - \alpha(1+\beta\gamma) \\ -\beta\kappa_i - (1-\nu)(1+\beta-\nu)\kappa_p + \alpha\beta\gamma \\ (1-\nu)(\beta-\nu)\kappa_p \end{bmatrix} = 0 \quad (27)$$

Therefore, from both algorithms coincide when the following system of equations is satisfied:

$$\kappa_i + (1-\nu)\kappa_p = \alpha(1+\beta\gamma) \quad (28a)$$

$$\beta\kappa_i + (1-\nu)(1-\nu+\beta)\kappa_p = \alpha\beta\gamma \quad (28b)$$

$$(1-\nu)(\beta-\nu)\kappa_p = 0 \quad (28c)$$

For  $\beta \neq 1$ , the solution to this system is given by:

$$\nu \leftarrow \beta \quad \kappa_i \leftarrow \frac{\alpha}{1-\beta} \quad \kappa_p \leftarrow -\frac{\alpha\beta}{(1-\beta)^2} [1-\gamma(1-\beta)] \quad (29)$$

Finally, we choose the initial condition  $\xi_0$  that ensures that the first two steps of the algorithms match (at  $t = 0$  and  $t = 1$ ). The first iterate of  $\nu$ PI is given by  $\theta_1 = \theta_0 + \kappa_i e_0 + \kappa_p \xi_0$  as per Eq. (15a). Meanwhile, the first iterate of UNIFIEDMOMENTUM is given by:

$$\theta_1 \stackrel{(17a)}{=} \theta_0 + \alpha(1+\beta\gamma)e_0 \stackrel{(29)}{=} \theta_0 + (\kappa_i + (1-\beta)\kappa_p)e_0 = \theta_0 + \kappa_i e_0 + (1-\beta)\kappa_p e_0 \quad (30)$$

Therefore, setting  $\xi_0 \leftarrow (1-\beta)e_0$  makes both algorithms match in their first step at  $t = 0$ .

The second iterate from UNIFIEDMOMENTUM is  $\theta_2 = \theta_1 + \alpha\beta [1-\gamma(1-\beta)] e_0 + \alpha [1-\gamma(1-\beta)] e_1$ . On the other hand, the second iterate of  $\nu$ PI is  $\theta_2 = \theta_1 + (\kappa_i + \kappa_p(1-\nu))e_1 - \kappa_p(1-\beta)\xi_0$ . It is easy to see that, given the hyperparameter choices outlined above, both algorithms match at  $t = 1$ .

An induction argument yields the equivalence between the algorithms.  $\square$

### C. Interpreting the updates of $\nu$ PI

Consider the execution of the algorithms  $\nu$ PI ( $\nu, \kappa_p, \kappa_i$ ) and GA ( $\alpha = \kappa_i$ ) at time  $t$ , with updates given by:

$$\theta_{t+1}^{\nu\text{PI}} = \theta_t + \kappa_i e_t + \kappa_p(1-\nu)(e_t - \xi_{t-1}) \quad (31)$$

$$\theta_{t+1}^{\text{GA}} = \theta_t + \kappa_i e_t \quad (32)$$

Let  $\psi = \frac{\kappa_p(1-\nu)}{\kappa_i + \kappa_p(1-\nu)}$ . Note that whenever  $\kappa_p$  and  $\kappa_i$  are non-negative,  $\psi \in [0, 1]$ . The ratio between these updates is:

$$\frac{\Delta\nu\text{PI}}{\Delta\text{GA}} = \frac{\theta_{t+1}^{\nu\text{PI}} - \theta_t}{\theta_{t+1}^{\text{GA}} - \theta_t} = \frac{\kappa_i e_t + \kappa_p(1-\nu)(e_t - \xi_{t-1})}{\kappa_i e_t} = 1 + \frac{\kappa_p(1-\nu)}{\kappa_i} - \frac{\kappa_p(1-\nu)}{\kappa_i} \frac{\xi_{t-1}}{e_t} = \frac{1}{1-\psi} \left[ 1 - \frac{\psi\xi_{t-1}}{e_t} \right]. \quad (33)$$

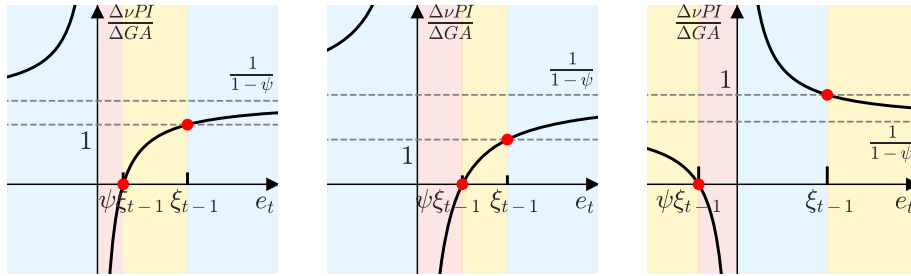


Figure 10: Comparing the update of  $\nu$ PI relative to GA, for different hyper-parameter configurations of  $\nu$ PI. **Left:**  $\nu$ PI is configured to recover POLYAK ( $\beta = -0.4$ ). Updates exhibit dampening similar to that of  $\nu$ PI ( $\nu = 0, \kappa_p = 1, \kappa_i = 1$ ). **Middle:**  $\nu$ PI ( $\nu = 0, \kappa_p = 1, \kappa_i = 1$ ) corresponding to a PI controller.  **$\nu$ PI increases the multipliers faster than GA when the constraint violation is large, enhancing convergence speed; and proactively decreases them near the feasible set, preventing overshoot.** **Right:**  $\nu$ PI is configured to recover POLYAK ( $\beta = 0.3$ ). We observe an increased eagerness to increase the multipliers *as progress toward feasibility occurs*. This increases the chances of overshoot and subsequent oscillations. The blue, yellow, and red regions correspond to cases in which the updates performed by the  $\nu$ PI algorithm are faster, slower, or in the opposite direction than those of GA, respectively. This plot illustrates the case  $\xi_{t-1} > 0$ . The middle and right figures presented here are the same as those in Fig. 5. We include them here for the reader's convenience.

**PI (Fig. 10, middle).** We consider  $\nu$ PI ( $\nu = 0, \kappa_p = 1, \kappa_i = 1$ ), which recovers PI ( $\kappa_p = 1, \kappa_i = 1$ ). The update of the PI optimizer relative to GA is as follows:



1. **Mode A** When either  $e_t \geq \xi_{t-1}$  or  $e_t \leq 0$ , the relative update exceeds one and thus the PI controller update can be seen as *eager* compared to gradient ascent.
  - (a) When  $e_t \geq \xi_{t-1}$ , the constraint has historically been infeasible and the current violation indicates an *increase in infeasibility*. In this case, PI not only increases the value of the multiplier but does so more strongly than GA. This proactive behavior serves to counteract the infeasibility increase.
  - (b) When  $e_t \leq 0$ , the constraint has been satisfied despite historical infeasibility ( $\xi_{t-1} > 0$ ). Here, the PI controller decreases the multiplier much more than GA. This serves to prevent overshoot into the feasible region.
2. **Mode B** In the range  $0 < \psi\xi_{t-1} \leq e_t \leq \xi_{t-1}$ , the constraint at step  $t$  (i) is not satisfied, (ii) it is smaller than the historical EMA of violations  $\xi_{t-1}$ , but not significantly (not beyond a factor of  $\psi$ ). In this case, the PI controller proactively exerts *friction* by having a smaller update than GA. This reduces the risk of overshoot under the assumption that the primal variables continue to make progress toward feasibility.
3. **Mode C** In the *optimistic* phase, where  $0 \leq e_t \leq \kappa\xi_{t-1}$ , the PI controller's update goes in the opposite direction to that of GA:  $\frac{\Delta\nu\text{PI}}{\Delta\text{GA}} \leq 0$ . This corresponds to a scenario where the constraint made significant progress toward feasibility relative to the historic violation EMA. While GA would increase the multiplier in this case (since  $g_t > 0$ ), PI *decreases* the value of the multiplier. This is useful to prevent overshoot since significant progress toward feasibility is an indicator that the multiplier is already exerting sufficient pressure for the constraints to be satisfied.

**Negative momentum (Fig. 10, left).** We consider POLYAK ( $\beta = -0.4$ ) as a realization of  $\nu\text{PI}$ , following Thm. 1. We observe similar behavior to that of  $\nu\text{PI}$  ( $\nu = 0, \kappa_p = 1, \kappa_i = 1$ ), in the middle figure of Fig. 10. Note that the current illustration assumes an equal value of the “optimizer state”  $\xi_{t-1}$  between the momentum and non-momentum cases. However, the value of  $\xi_t$  will be different depending on the momentum coefficient as  $\beta = \nu$  also influences the update of  $\xi$  (see Algo. 1).

**Positive momentum (Fig. 10, right).** The right plot of Fig. 10 considers POLYAK ( $\beta = 0.3$ ) as a realization of  $\nu\text{PI}$ , following Thm. 1. We observe significantly different behavior compared to the left and middle plots.

1. **Mode A** When infeasibility is reduced, the algorithm is *eager* to increase the multiplier more than GA. This is a counter-intuitive operation of the algorithm considering that the current value of the multiplier can apply sufficient pressure to improve the constraint satisfaction. Increasing the multiplier further can lead to a higher risk of overshoot.
2. **Mode B** Consider the cases in which infeasibility increases ( $e_t \geq \xi_{t-1}$ ), or the constraints suddenly become (sufficiently) strictly feasible  $e_t \leq \psi\xi_{t-1} \leq 0$ . These cases induce *frictioned* updates with the same sign as GA, but of smaller magnitude.
3. **Mode C** When the primal player is feasible, positive momentum would result in an *increase* of the multiplier; going against the update of GA, which would *decrease* the multiplier. In this context, increasing the multiplier is unreasonable since the current value of the multiplier is already sufficient to achieve feasibility.

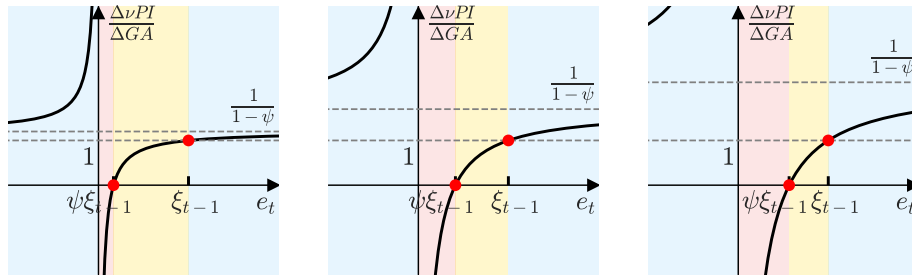


Figure 11: Effect of  $\kappa_p$  in the update of  $\nu\text{PI}$  relative to GA. When  $\kappa_p$  approaches 0,  $\nu\text{PI}$  recovers GA (a constant function  $y = 1$  for the relative update). **A larger  $\kappa_p$  leads to a wider “optimistic” region (in red) where  $\nu\text{PI}$  decreases the multiplier to prevent overshooting despite the constraint being violated.** We use  $\kappa_i = 1$  and  $\nu = 0$  and  $\kappa_p$  of 0.2, 0.7 and 1.3, respectively.

**Ablation on the influence of  $\kappa_p$ .** Figure 11 presents three configurations of  $\kappa_p$  for a  $\nu$ PI ( $\nu = 0, \kappa_p, \kappa_i = 1$ ) optimizer. We display  $\kappa_p$  at 0.2, 0.7 and 1.3, respectively. As  $\kappa_p \rightarrow 0$ ,  $\nu$ PI is equivalent to GA. This is confirmed by the relative updates between  $\nu$ PI and GA converging to a constant function  $y = 1$ . As  $\kappa_p$  increases in the middle and right plots of Fig. 11, the asymptote at  $1/(1 - \psi)$  moves further away from 1, and the width of the “optimistic” region (**Mode C**) increases. In other words, as  $\kappa_p$  grows, the threshold for “sufficient improvement” is relaxed and the optimizer is more prone to decrease the multipliers upon improvements in constraint violation. This leads to a more “cautious” behavior from the algorithm: the multiplier is decreased earlier when the problem approaches the feasible region, which prevents overshooting but with potentially slower convergence. One can monotonically control for the convergence and overshoot behaviors by adjusting the  $\kappa_p$  value, see Fig. 9 in §5.3.

## D. Analysis of continuous-time $\nu$ PI dynamics as an oscillator

In this section, we examine the spectral properties of the gradient-descent/ $\nu$ PI flow dynamics presented in Algo. 3. We extend the analysis of [Stooke et al. \(2020\)](#) (which only considers the dynamics of  $\mathbf{x}$ ) to also consider  $\boldsymbol{\lambda}$ .

Consider a constrained optimization problem with equality constraints  $\mathbf{h}$ . The GD/ $\nu$ PI flow corresponds to a *continuous-time* dynamical system in which the primal player implements gradient descent on the Lagrangian, and the dual player implements  $\nu$ PI ascent. This is formalized in Algo. 3.

---

### Algorithm 3 Continuous-time gradient descent/ $\nu$ PI

---

**Args:** proportional ( $\kappa_p$ ) and integral ( $\kappa_i$ ) gains for  $\nu$ PI flow

1:  $\dot{\mathbf{x}} = -\nabla f(\mathbf{x}) - \mathcal{J}\mathbf{h}(\mathbf{x})\boldsymbol{\mu}$

2:  $\dot{\boldsymbol{\mu}} = \kappa_i \mathbf{h}(\mathbf{x}) + \kappa_p \dot{\mathbf{h}}(\mathbf{x})$

---

Thm. 4 characterizes the GD/ $\nu$ PI flow in Algo. 3 in terms of a second-order dynamical system. Note that this relationship holds for any constrained problem where the objective and constraints have second derivatives. Appx. D.2 analyzes the resulting dynamical system for a quadratic program with linear equality constraints.

### D.1. Oscillator dynamics of GD/ $\nu$ PI flow

**Theorem 4.** *The dynamics of Algo. 3 can be characterized by the following system of second-order differential equations, with initial conditions  $\mathbf{x}(0) = \mathbf{x}_0$ ,  $\boldsymbol{\mu}(0) = \boldsymbol{\mu}_0$ ,  $\dot{\mathbf{x}}(0) = -\nabla f(\mathbf{x}_0) - \mathcal{J}\mathbf{h}(\mathbf{x}_0)\boldsymbol{\mu}_0$ , and  $\dot{\boldsymbol{\mu}}(0) = \kappa_i \mathbf{h}(\mathbf{x}_0) + \mathcal{J}\mathbf{h}(\mathbf{x}_0)\dot{\mathbf{x}}(0)$ :*

$$\left\{ \begin{array}{l} \ddot{\mathbf{x}} = - \underbrace{\left( \nabla^2 f + \sum_{c'} \mu_{c'} \nabla^2 \mathbf{h}_{c'} \right)}_{\boldsymbol{\Phi}} \dot{\mathbf{x}} - \mathcal{J}\mathbf{h}\dot{\boldsymbol{\mu}} \\ \ddot{\boldsymbol{\mu}} = \kappa_i \mathcal{J}\mathbf{h}^\top \dot{\mathbf{x}} + \kappa_p \mathcal{J}\mathbf{h}^\top \ddot{\mathbf{x}} + \kappa_p \boldsymbol{\Xi} \end{array} \right. \quad (34a)$$

$$\left\{ \begin{array}{l} \ddot{\mathbf{x}} = - \underbrace{\left( \nabla^2 f + \sum_{c'} \mu_{c'} \nabla^2 \mathbf{h}_{c'} \right)}_{\boldsymbol{\Phi}} \dot{\mathbf{x}} - \mathcal{J}\mathbf{h}\dot{\boldsymbol{\mu}} \\ \ddot{\boldsymbol{\mu}} = \kappa_i \mathcal{J}\mathbf{h}^\top \dot{\mathbf{x}} + \kappa_p \mathcal{J}\mathbf{h}^\top \ddot{\mathbf{x}} + \kappa_p \boldsymbol{\Xi} \end{array} \right. \quad (34b)$$

where  $\boldsymbol{\Xi} = [\dot{\mathbf{x}}^\top \nabla^2 h_1 \dot{\mathbf{x}}, \dots, \dot{\mathbf{x}}^\top \nabla^2 h_c \dot{\mathbf{x}}]^\top \in \mathbb{R}^c$ .

This can be concisely represented in matrix form as:

$$\begin{bmatrix} \mathbf{I}_{n \times n} & \mathbf{0}_{n \times c} \\ -\kappa_p \mathcal{J}\mathbf{h}^\top & \mathbf{I}_{c \times c} \end{bmatrix} \begin{bmatrix} \ddot{\mathbf{x}} \\ \ddot{\boldsymbol{\mu}} \end{bmatrix} + \begin{bmatrix} \boldsymbol{\Phi} & \mathcal{J}\mathbf{h} \\ -\kappa_i \mathcal{J}\mathbf{h}^\top & \mathbf{0}_{c \times c} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ -\beta \boldsymbol{\Xi} \end{bmatrix} = \mathbf{0}. \quad (35)$$

Or, equivalently:

$$\begin{bmatrix} \ddot{\mathbf{x}} \\ \ddot{\boldsymbol{\mu}} \end{bmatrix} + \begin{bmatrix} \boldsymbol{\Phi} & \mathcal{J}\mathbf{h} \\ \mathcal{J}\mathbf{h}^\top (\kappa_p \boldsymbol{\Phi} - \kappa_i \mathbf{I}) & \kappa_p \mathcal{J}\mathbf{h}^\top \mathcal{J}\mathbf{h} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ -\beta \boldsymbol{\Xi} \end{bmatrix} = \mathbf{0}. \quad (36)$$

**Proof of Thm. 4.** We start by computing the time derivatives of the objective gradient and constraint Jacobian:

$$\frac{d}{dt} [\nabla f] = \left[ \sum_j \frac{\partial(\nabla f)}{\partial x_j} \frac{dx_j}{dt} \right]_i = \nabla^2 f \dot{\mathbf{x}} \quad \frac{d}{dt} [\mathcal{J}\mathbf{h}] = [\nabla^2 \mathbf{h}_1 \dot{\mathbf{x}} \quad \nabla^2 \mathbf{h}_2 \dot{\mathbf{x}} \quad \dots \quad \nabla^2 \mathbf{h}_c \dot{\mathbf{x}}] \quad (37)$$

Therefore, the second order dynamics for  $\mathbf{x}$  are given by:

$$\ddot{\mathbf{x}} = \frac{d}{dt} [-\nabla f(\mathbf{x}) - \mathcal{J}\mathbf{h}(\mathbf{x})\boldsymbol{\mu}] = -\frac{d}{dt} [\nabla f] - \mathcal{J}\mathbf{h}\dot{\boldsymbol{\mu}} - \frac{d}{dt} [\mathcal{J}\mathbf{h}] \boldsymbol{\mu} \quad (38a)$$

$$= -\nabla^2 f \dot{\mathbf{x}} - \mathcal{J}\mathbf{h}\dot{\boldsymbol{\mu}} - \sum_{c'} \mu_{c'} \nabla^2 \mathbf{h}_{c'} \dot{\mathbf{x}} \quad (38b)$$

$$= - \underbrace{\left( \nabla^2 f + \sum_{c'} \mu_{c'} \nabla^2 \mathbf{h}_{c'} \right)}_{\Phi} \dot{\mathbf{x}} - \mathcal{J}\mathbf{h}\dot{\boldsymbol{\mu}} \quad (38c)$$

The second order dynamics for  $\boldsymbol{\mu}$  are given by:

$$\ddot{\boldsymbol{\mu}} = \frac{d}{dt} [\kappa_i \mathbf{h} + \kappa_p \mathcal{J}\mathbf{h}^\top \dot{\mathbf{x}}] = \alpha \dot{\mathbf{h}} + \kappa_p \Xi \dot{\mathbf{x}} + \kappa_p \mathcal{J}\mathbf{h}^\top \ddot{\mathbf{x}}, \quad (39)$$

where  $\Xi$  is defined as:

$$\Xi \triangleq \frac{d}{dt} [\mathcal{J}\mathbf{h}^\top] \dot{\mathbf{x}} = [\dot{\mathbf{x}}^\top \nabla^2 \mathbf{h}_1 \dot{\mathbf{x}} \quad \dot{\mathbf{x}}^\top \nabla^2 \mathbf{h}_2 \dot{\mathbf{x}} \quad \dots \quad \dot{\mathbf{x}}^\top \nabla^2 \mathbf{h}_c \dot{\mathbf{x}}]^\top. \quad (40)$$

□

## D.2. Dynamics of GD/PI flow for a constrained quadratic program

Let  $\mathbf{H} \in \mathbb{R}^{n \times n}$  be positive semi-definite and consider the convex quadratic program with  $c$  linear constraints:

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{c}^\top \mathbf{x} \quad \text{subject to} \quad \mathbf{A} \mathbf{x} - \mathbf{b} = 0. \quad (41)$$

The Lagrangian min-max game associated with the problem in Eq. (41) is given by:

$$\mathfrak{L}(\mathbf{x}, \boldsymbol{\mu}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{c}^\top \mathbf{x} + \boldsymbol{\mu}^\top (\mathbf{A} \mathbf{x} - \mathbf{b}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x} + \mathbf{c}^\top \mathbf{x} + \boldsymbol{\mu}^\top \mathbf{A} \mathbf{x} - \boldsymbol{\mu}^\top \mathbf{b}. \quad (42)$$

The linearity of the constraints in Eq. (41) implies  $\mathcal{J}\mathbf{h} = \mathbf{A}^\top$  and  $\nabla^2 g_{c'} = \mathbf{0}$  for  $c' \in [c]$ , thus  $\Phi = \mathbf{H}$  and  $\Xi = \mathbf{0}$ . Therefore, we obtain a homogeneous system of second-order differential equations with constant coefficients:

$$\begin{bmatrix} \ddot{\mathbf{x}} \\ \ddot{\boldsymbol{\mu}} \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{H} & \mathbf{A}^\top \\ \mathbf{A}(\kappa_p \mathbf{H} - \kappa_i \mathbf{I}) & \kappa_p \mathbf{A} \mathbf{A}^\top \end{bmatrix}}_{\mathbf{U}} \begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\boldsymbol{\mu}} \end{bmatrix} = \mathbf{0}. \quad (43)$$

A simple state transformation  $\mathbf{z} = [\mathbf{x}, \boldsymbol{\mu}, \dot{\mathbf{x}}, \dot{\boldsymbol{\mu}}]^\top$  and  $\dot{\mathbf{z}} = [\dot{\mathbf{x}}, \dot{\boldsymbol{\mu}}, \ddot{\mathbf{x}}, \ddot{\boldsymbol{\mu}}]^\top$  yields:

$$\dot{\mathbf{z}} = - \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \mathbf{0} & \mathbf{U} \end{bmatrix} \mathbf{z} = - \begin{bmatrix} \mathbf{0}_{(n+c) \times (n+c)} & \mathbf{I}_{(n+c) \times (n+c)} \\ \mathbf{0}_{(n+c) \times (n+c)} & \begin{bmatrix} \mathbf{H} & \mathbf{A}^\top \\ \mathbf{A}(\kappa_p \mathbf{H} - \kappa_i \mathbf{I}) & \kappa_p \mathbf{A} \mathbf{A}^\top \end{bmatrix} \end{bmatrix} \mathbf{z} \quad (44)$$

Therefore, this  $2(n+c)$ -dimensional linear system has zero as an eigenvalue with algebraic multiplicity  $n+c$ , and the remaining eigenvalues correspond to the spectrum of  $-\mathbf{U}$ .

When the matrix  $\mathbf{H}$  is zero, we recover the smooth bilinear games considered by Gidel et al. (2019a, Eq. 18) in their study of negative momentum. In this case, the matrix  $\mathbf{U}$  looks like:

$$-\mathbf{U} = - \begin{bmatrix} \mathbf{0} & \mathbf{A}^\top \\ -\kappa_i \mathbf{A} & \kappa_p \mathbf{A} \mathbf{A}^\top \end{bmatrix} \quad (45)$$

It is easy to see that large enough values of  $\kappa_p$  cause the eigenvalues of the matrix to have negative real parts, and thus make the system converge. However, if  $\kappa_p = 0$  (i.e. gradient descent-ascent), the eigenvalues of this matrix are either 0 or pure

imaginary. This fact is in line with existing results in the literature on the lack of convergence gradient descent-ascent on bilinear games (Gidel et al., 2019a).

**Case of one-variable and one constraint.** It is instructive to analyze the spectrum of  $\mathbf{U}$  in the case of a problem with a one-dimensional primal variable and a single constraint (and thus one multiplier). In this case,  $\mathbf{U}$  and its eigenvalues take the form:

$$-\mathbf{U} = - \begin{bmatrix} h & a \\ a(\kappa_p h - \kappa_i) & \kappa_p a^2 \end{bmatrix} \quad \lambda = \frac{-(h + \kappa_p a^2) \pm \sqrt{(h + \kappa_p a^2)^2 - 4a^2 \kappa_i}}{2} \quad (46)$$

As before, the eigenvalues of this matrix depend on the choice of  $\kappa_p$ . This is illustrated in Fig. 12.

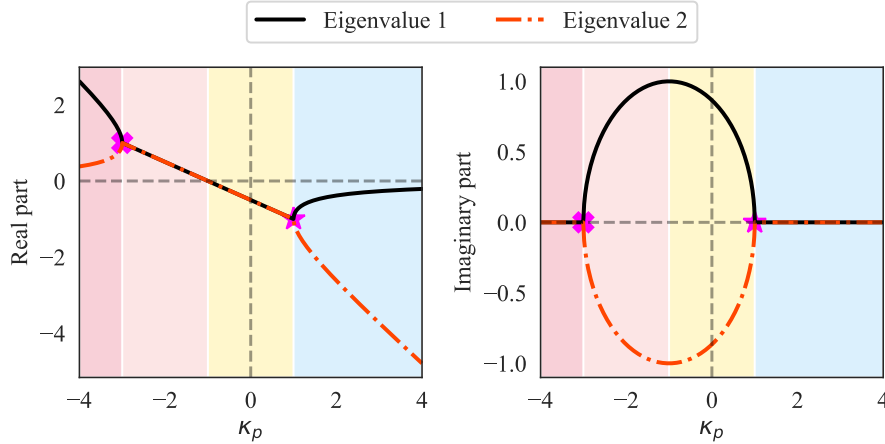


Figure 12: Eigenvalues for Eq. (41) as a function of  $\kappa_p$  in the one-dimensional case. **A positive value of  $\kappa_p$  (denoted by  $\star$ ) achieves critical damping (i.e. equal convergence rate for both dimensions).** This plot uses  $h = 1$ ,  $a = -1$  and  $\kappa_i = 1$ .

Note that when the discriminant of Eq. (46) is zero, both eigenvalues match (and must thus be real). When this occurs and both eigenvalues are negative, the system converges and does so at the same rate in both dimensions. This is akin to the notion of *critical damping* from the control theory literature.

The values of  $\kappa_p$  that make the discriminant zero are  $\kappa_p^* = \frac{-h \pm 2|a|\sqrt{\kappa_i}}{a^2}$ , leading to the eigenvalues  $\lambda(\kappa_p^*) = \frac{-(h + \kappa_p^* a^2)}{2} = \mp a\sqrt{\kappa_i}$ . These values of  $\kappa_p^*$  are marked with  $\star$  and  $\times$  in Fig. 12. Note that out of the two values of  $\kappa_p$  producing matching eigenvalues, only the choice  $\kappa_p^* > 0$  yields a convergent system.

More generally, depending on  $\kappa_p$ , the system exhibits different behaviors:

- **Divergence.** In the red regions, the system *diverges*; in light red region, this happens together with oscillations. Note how all the divergent configurations use a negative value of  $\kappa_p$ . The fuchsia cross ( $\times$ ) denotes the value of  $\kappa_p$  for which both dimensions diverge at the same rate.
- **Underdamping.** In the yellow region, the system is *underdamped* and *converges with oscillations*. Interestingly, this system admits some negative values of  $\kappa_p$  (of sufficiently small magnitude) while remaining convergent.
- **Critical damping.** The fuchsia star ( $\star$ ) shows the  $\kappa_p$  value that makes both dimensions of the system converge *at the same rate*. Note that this *critical damping* regime is achieved at a strictly positive value of  $\kappa_p$ , and thus is not achievable by gradient ascent.
- **Overdamping.** In the blue region, the system is *convergent without oscillations* but *overdamped* since the dimension corresponding to the black eigenvalue converges more slowly.

## E. Illustrative 2D nonconvex problem

We demonstrate the behavior of  $\nu$ PI on the two-dimensional, nonconvex, equality-constrained problem in Eq. (47). This problem was proposed by Boyd (2021). The setting is simple enough to allow for visualizing the optimization paths of each optimization variable and multiplier, while also being challenging due to nonconvexity.

$$\min_{\mathbf{x}=(x_1,x_2)} f(\mathbf{x}) \triangleq \left\| \begin{bmatrix} x_1 + e^{-x_2} \\ x_1^2 + 2x_2 + 1 \end{bmatrix} \right\|_2^2 \quad \text{subject to } h(\mathbf{x}) \triangleq x_1 + x_1^3 + x_2 + x_2^2 = 2. \quad (47)$$

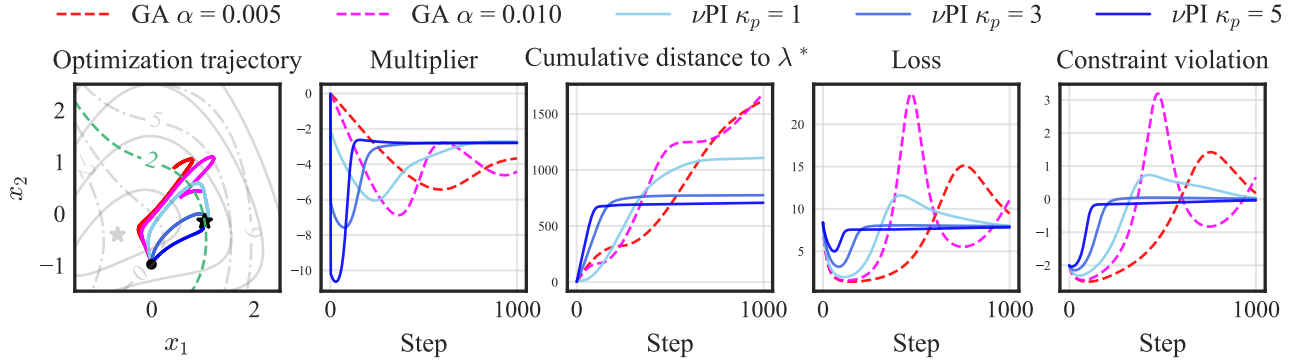


Figure 13: Optimization trajectories for different algorithms on a 2D nonconvex equality-constrained minimization problem.  $\nu$ PI runs use  $\nu = 0$  and  $\kappa_i = 0.01$ . The light gray  $\star$  marks the *unconstrained* optimum, while the black  $\star$  marks the *constrained* optimum. Level sets correspond to the objective function (solid) and constraint (dashed).

**GA trajectories.** In Fig. 13, GA trajectories are initially drawn toward the direction of the unconstrained optimum since multipliers grow slowly at first. As training progresses, the constraint plays a more significant role. With a step-size that is too small ( $\alpha = 0.005$ ), the trajectory does not converge to the global optimum. In contrast, the system reaches the global constrained optimum point when employing a larger step-size ( $\alpha = 0.01$ ). This is achieved while incurring in *oscillations*. The phase change from not converging with a small step-size, to converging with oscillations indicates that GA is not suitable for obtaining critical damping when solving the problem.

**$\nu$ PI trajectories.** The three blue trajectories in Fig. 13 show different behaviors of  $\nu$ PI: underdamping (light blue,  $\kappa_p = 1$ ), almost-critical damping ( $\kappa_p = 3$ ) and overdamping (dark blue,  $\kappa_p = 5$ ). Note the monotonic effect of  $\kappa_p$  on the damping of the system.  $\nu$ PI provides the flexibility to obtain different levels of constraint overshoot, and can achieve feasibility and convergence at different speeds. This added flexibility leads to enhanced control over the dynamics of the system relative to GA, thus enabling applications of  $\nu$ PI to safety-sensitive tasks.

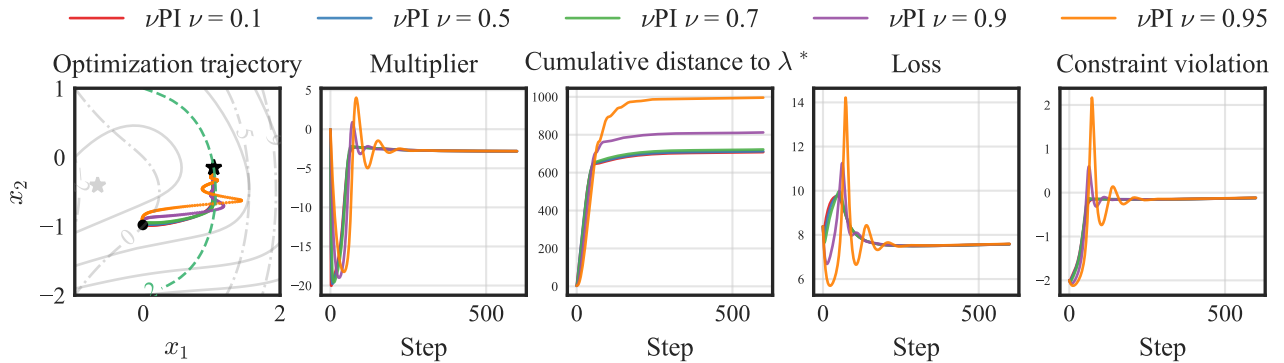


Figure 14: Optimization trajectories for the  $\nu$ PI algorithm under different choices of  $\nu$ .  $\nu$ PI runs use  $\kappa_i = 0.01$  and  $\kappa_p = 10$ .

**Ablation on  $\nu$ .** In Fig. 14, we zoom in on the effect of  $\nu$  for fixed choices of  $\kappa_p$  and  $\kappa_i$ . A  $\nu$  closer to 0 tends towards PI, whereas a  $\nu$  closer to one gives more importance to historical constraint violations. We observe that a larger  $\nu$  behaves

qualitatively similar to positive momentum: the multiplier tends to increase faster if the constraint is not satisfied for a period of time. In this example, this leads to oscillations as shown for  $\nu = 0.95$ . Since this problem is deterministic, using a non-zero  $\nu$  does not show any advantage. Our fairness experiments showcase an application where  $\nu > 0$  is beneficial.

## F. Experimental details

Our implementation use PyTorch 2.0.0 (Paszke et al., 2019) and the Cooper library for Lagrangian constrained optimization (Gallego-Posada & Ramirez, 2022). Our code is available at: <https://github.com/motahareh-sohrabi/nuPI>.

### F.1. Linear SVM experiments

In our experiment with linear SVMs, we focus on two linearly separable classes from the Iris dataset (Fisher, 1988). We select 100 instances from the Iris setosa and Iris versicolor species, which are two linearly separable classes. Each data point in this dataset has four features. We selected 70% of data for training and the rest for validation. This gives the algorithm 70 Lagrange multipliers to learn.

We know that a unique  $\lambda^*$  exists in our experiments. The linearly independent constraint qualification (LICQ) holds for the selected data, so the Karush-Kuhn-Tucker (KKT) conditions imply the existence and uniqueness of  $\lambda^*$  at the constrained optimum  $x^*$ . All of the methods that do not diverge achieve perfect training and validation accuracy in this task.

**Experiment configuration and hyperparameters.** Throughout all of the experiments, we fixed the primal optimizer and only changed the dual optimizer. The primal optimizer is gradient descent with momentum, with a step size of  $10^{-3}$  and momentum of 0.9.

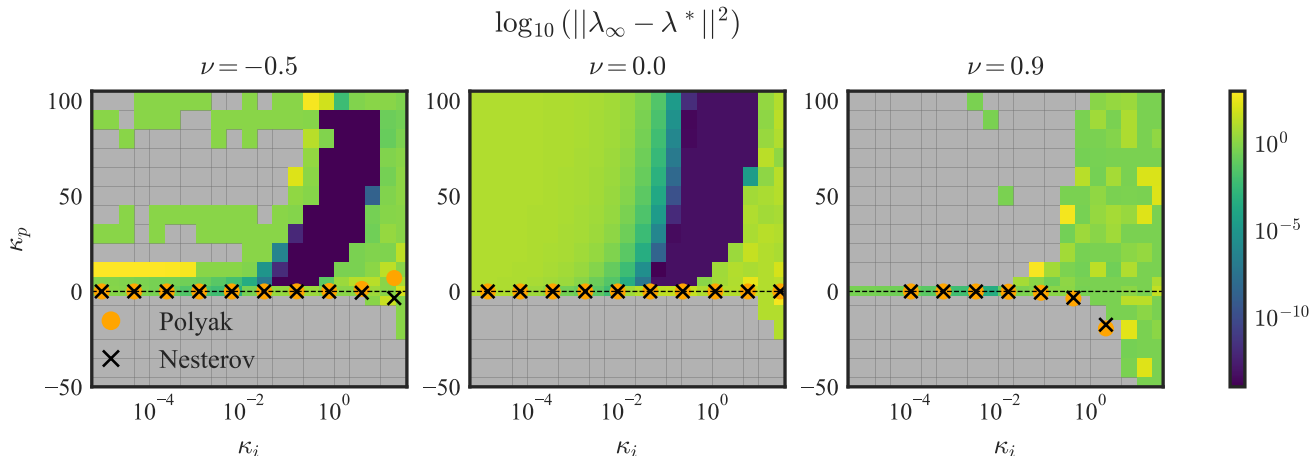


Figure 15: Distance to optimal Lagrange multipliers for different selections of parameter  $\nu$  in  $\nu$ PI algorithm. We also show where the equivalent  $\kappa_p$  and  $\kappa_i$  parameters for  $\text{NESTEROV}(\alpha, \beta = \nu)$  and  $\text{POLYAK}(\alpha, \beta = \nu)$  lie according to Thm. 1 for different values of  $\alpha$ . **The  $\nu$ PI algorithm with  $\nu = 0$  can give the highest number of converging step-sizes. While negative  $\nu = 0.5$  induces a range of converging step-sizes as well, there is no value of  $\kappa_i$  that the algorithm converges for  $\nu = 0.9$ .** The gray color shows the runs exceeding a distance of  $10^3$  to  $\lambda^*$ .

**Different values of the parameter  $\nu$  in  $\nu$ PI algorithm.** We examine how changing the parameter  $\nu$  in the Algo. 1 can affect the convergence of the SVM task with different choices of  $\kappa_i$  and  $\kappa_p$ . Fig. 15 shows how  $\nu$ PI behaves when choosing a negative, zero and positive value of  $\nu$ . While  $\nu = -0.5$  can lead to a converging algorithm for some step-sizes,  $\nu = 0.0$  offers a wider range of converging step-sizes. There is no choice of step-size for which the  $\nu$ PI algorithm with a positive value of  $\nu = 0.9$  converges to  $\lambda^*$ .

**Relationship between momentum and  $\nu$ PI algorithms.** Thm. 1 suggests that Polyak and Nesterov momentum algorithms can be instantiated using a specific choice of  $\kappa_i$  and  $\kappa_p$  in the Algo. 1. In Fig. 15 we show where  $\text{POLYAK}(\alpha, \beta)$  and  $\text{NESTEROV}(\alpha, \beta)$  lie for a choice of  $\alpha$ 's and with  $\beta = \nu$ . For each pair of  $(\alpha_i, \beta)$ , we calculate the value of  $\kappa_i$  and  $\kappa_p$  that recover  $\text{POLYAK}(\alpha_i, \beta)$  and  $\text{NESTEROV}(\alpha_i, \beta)$  according to Thm. 1.

- When  $\beta = \nu = 0$  there is no momentum and both POLYAK and NESTEROV reduce to gradient ascent. Therefore, all of the dots indicating momentum methods lie in the  $\kappa_p = 0$  line for the middle plot of Fig. 15.
- A positive  $\beta = \nu = 0.9$  (Fig. 15, right) corresponds to a common momentum choice in minimization problems. We can see how there is no step-size value for which POLYAK or NESTEROV converge in this task. This observation supports the claim of Gidel et al. (2019a) on the ineffectiveness of positive momentum for convergence in games.
- A negative  $\beta = \nu = -0.5$  leads to convergence for some choices of step-size. However, these do not correspond to what POLYAK and NESTEROV can achieve. This is consistent with our observation in Fig. 6 (left), indicating the necessity of adding a (positive)  $\kappa_p$  term to the optimizer in order to achieve convergence. This highlights the benefits of the increased generality of  $\nu$ PI.

Moreover, we observe that POLYAK with negative momentum achieves a positive  $\kappa_p$  while all momentum choices for NESTEROV lead to negative  $\kappa_p$  values. This further supports the experimental results of Gidel et al. (2019a), where POLYAK is used when they want to experiment with negative momentum. Our hypothesis is that negative momentum with Polyak is successful in games because it can induce a positive  $\kappa_p$ .

## F.2. Sparsity experiments

**Background.** Louizos et al. (2018) propose a re-parameterization of models that allows applying  $L_0$ -norm regularization on their weights. They propose the use of stochastic gates  $z$  that indicate whether each parameter is active or not, where  $z$  follows a hard-concrete distribution parameterized by  $\phi$ . Employing the re-parameterization trick allows the computation of gradients of the  $L_0$ -norm of the model with respect to  $\phi$ . Gallego-Posada et al. (2022) formulate a constrained optimization problem that prescribes the desired sparsity of the model as a constraint.

$$\min_{\mathbf{w}, \phi \in \mathbb{R}^d} \mathbb{E}_{z|\phi} [L(\mathbf{w} \odot z | \mathcal{D})] \quad \text{s.t.} \quad \frac{\mathbb{E}_{z|\phi} [\|z\|_0]}{\#(\mathbf{w})} \leq \epsilon, \quad (48)$$

where  $\mathbf{x}$  are the parameters of the model,  $L$  is an ERM objective, and  $\mathcal{D}$  is a dataset. The constraint is normalized with the total number of parameters of the model  $\#(\cdot)$ , so that the constraint level  $\epsilon \in [0, 1]$  corresponds to the maximum allowed *proportion* of active parameters. For details on the re-parameterization, and a closed form expression for  $\mathbb{E}_{z|\phi} [\|z\|_0]$ , see Louizos et al. (2018); Gallego-Posada et al. (2022).

**Hard-concrete distribution.** The  $L_0$ -norm re-parameterization proposed by Louizos et al. (2018) considers a hard-concrete distribution for the stochastic gates of the model. The hard-concrete distribution is based on a stretched and clamped concrete distribution (Maddison et al., 2017). Similar to Louizos et al. (2018); Gallego-Posada et al. (2022), we choose a temperature of  $2/3$  for the concrete distribution, and a stretching interval of  $[-0.1, 1.1]$ .

**Architecture.** We consider ResNet-18 (He et al., 2016) models with basic residual blocks for our sparsity experiments, which have a total of approximately 11.2 million parameters. Following Louizos et al. (2018) and Gallego-Posada et al. (2022), we employ output feature map sparsity on the first convolutional layer of each residual block, whereas the following convolutional layer and the residual connection are kept to be fully dense. The first convolutional layer of the model and the linear output layer are also kept fully dense. This model counts with 8 sparsifiable convolutional layers.

**Choice of sparsity levels.** Although Gallego-Posada et al. (2022) consider up to 80% structured sparsity (20% density) for ResNet-18 models, Gale et al. (2019) indicate that it is possible to train ResNet-50 models with structured sparsity of up to 95% (5% density or less), without incurring on a catastrophic loss on model accuracy. Therefore, we consider sparsity levels of between 30% and 90% (70% and 10% density, respectively).

**Primal optimization.** We consider an optimization pipeline for the model that incorporates standard techniques used to train  $L_0$ -sparse ResNet-18 models on CIFAR10. For the weights of the model, we use SGD with a momentum of 0.9, an initial learning rate of 0.01, and a cosine annealing learning rate scheduler (Loshchilov & Hutter, 2017).

We initialize the gates with a *dropout init* of 0.01, effectively yielding a fully dense model at initialization. Akin to Gallego-Posada et al. (2022), we use ADAM (Kingma & Ba, 2014) with a step-size of  $8 \cdot 10^{-4}$  to optimize the  $\phi$  parameters of the stochastic gates. When applying  $L_2$ -norm regularization on the parameters, we detach the contribution of the gates as recommended by Gallego-Posada et al. (2022).

**Dual optimization.** For sparsity experiments, we consider  $\nu = 0$ . Since the constraint is deterministic given the state of the model (it does not need to be estimated from mini-batches), we consider the use of an EMA to not be crucial for

this task. Unless otherwise stated, we use a dual step-size of  $8 \cdot 10^{-4}$  for all dual optimizer choices (as was provided by Gallego-Posada et al. (2022)). We decide against tuning the dual step-size separately for each optimizer to highlight the flexibility of  $\nu$ PI: given a step-size that was tuned to yield good results for GA,  $\nu$ PI may produce better-behaved dynamics. All of our sparsity experiments use a batch size of 128 and are over 200 epochs.

### F.3. Fairness experiments

**Dataset.** In this experiment we consider the adult dataset (Becker & Kohavi, 1996), pre-processed following Zafar et al. (2019). The raw data comprises eight categorical and six continuous attributes. After processing, the data is comprised of 50-dimensional sparse feature vectors. The train and test sets consist of 30,162 and 15,060 samples, respectively.

**Background.** We consider a fairness task under the disparate impact constraint (Zafar et al., 2019) shown in Eq. (13). This constraint is also known as statistical parity and demographic parity (Corbett-Davies et al., 2017; Dwork et al., 2012). We consider two sensitive attributes in the adult dataset: sex, denoted as  $A_1 = \{male, female\}$ , and race, denoted as  $A_2 = \{White, Black, Asian-Pac-Islander, Amer-Indian-Eskimo, Other\}$ . Equation (13) entails the intersection of both attributes, leading to  $|A_1| \times |A_2| = 10$  constraints.

**Architecture and primal optimization.** We train a 2-hidden-layer neural network with hidden sizes of (100, 100) similar to the experimental setup of Cotter et al. (2019). In order to choose the primal optimizer hyperparameters, we trained the unconstrained problem and chose the parameters of the run with the highest training accuracy. We fixed this primal optimizer across our constrained experiments to be ADAM ( $\alpha = 10^{-2}$ ).

**Dual optimization.** We chose the best step-size for GA aimed at minimizing training accuracy, while ensuring that the maximum violation achieves the lowest possible value. This led to a dual step-size of  $\alpha = 0.03$ . We then fixed this value as the  $\kappa_i$  parameter of  $\nu$ PI and ran a grid search to find the best  $\kappa_p$ . The grid search for  $\kappa_p$  considered (logarithmically spaced) values in  $[0.01, 100]$ . The best results were found with  $\kappa_p = 5$ .

Due to the noise in the constraints, we also experimented with the effect of  $\nu$  on the optimization dynamics. We tried  $\nu$  values of 0.0, 0.5, 0.9, 0.95, and 0.99. We noticed that higher values of  $\nu$  can improve the learning dynamics, with the best results achieved at 0.99. Setting  $\nu = 0$  results in noisy Lagrange multipliers, which lead to unstable optimization. This is illustrated in Fig. 7.

## G. Comprehensive results on the sparsity task

In this section we provide extensive experiment results for our sparsity experiments, complementing §5.3. We conducted experiments with global and layer-wise sparsity targets, at  $\epsilon = 70\%$ ,  $50\%$ ,  $30\%$ ,  $10\%$  density levels. The shaded region of our plots corresponds to the feasible set. “Relative violations” are computed by dividing the absolute constraint violations by the target density.

### G.1. Global

For global sparsity experiments (Figs. 16 to 19 and Tables 3 to 6), we observe a general trend for models that overshoot into becoming excessively sparse to achieve a lower training performance. This insight is also generally true for validation accuracy. In particular, gradient ascent and momentum methods consistently exhibit overshoot, whereas  $\nu$ PI and gradient ascent with dual restarts do not overshoot and achieve good final performance. Dual restarts generally produce (slightly) infeasible solutions. Note that at  $\epsilon = 10\%$ , negative momentum runs do not overshoot, but positive momentum runs do.

### G.2. Layer-wise

We perform layer-wise sparsity experiments with  $\epsilon = 10\%$ ,  $30\%$ ,  $50\%$ ,  $70\%$  density targets (Tables 7 to 10). We observe a similar trend to global sparsity experiments: GA and momentum methods overshoot, while  $\nu$ PI and GA with dual restarts reliably achieve feasible solutions, with small levels of overshoot. Moreover, we observe that the violations of  $\nu$ PI and GA with dual restarts span a small range of values relative to other methods. This highlights the robustness of  $\nu$ PI since the  $\kappa_p$  coefficient did not need to be tuned independently for each constraint.



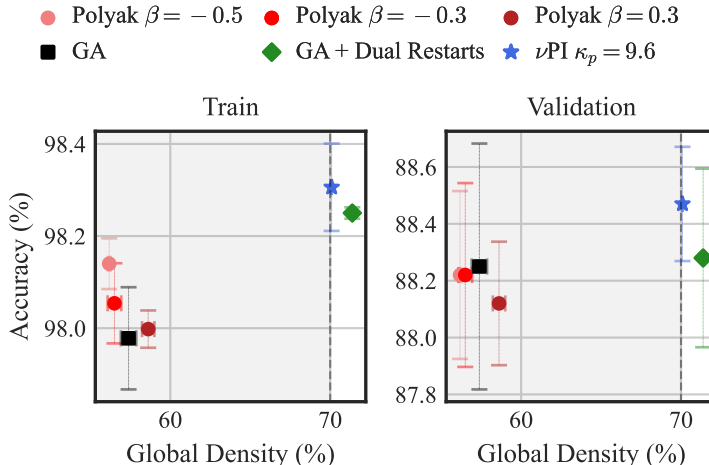


Figure 16: CIFAR10 trade-off plot for *global* sparsity under a 70% density target.  **$\nu$ PI successfully achieves the desired sparsity while achieving the highest train accuracy.** The shaded region is the feasible set. As higher density correlates to higher train accuracy, overshooting to a lower density is undesirable. All optimizers use the same step-size. *This figure is the same as Fig. 8. We repeat it here for the reader's convenience.*

Table 3: CIFAR10 results for *global* sparsity under a 70% density target.  **$\nu$ PI successfully achieves the desired sparsity while achieving the highest train accuracy.** The results in this table correspond to those in Fig. 16. As higher density correlates to higher train accuracy, overshooting to a lower density is undesirable. All optimizers use the same step-size.

Method	Train Acc.	Test Acc.	Violation	Relative Violation
Polyak $\beta = -0.5$	$98.1 \pm 0.06$	$88.2 \pm 0.30$	$-13.8 \pm 0.09$	$-19.7 \pm 0.13$
Polyak $\beta = -0.3$	$98.1 \pm 0.09$	$88.2 \pm 0.32$	$-13.5 \pm 0.43$	$-19.3 \pm 0.61$
Polyak $\beta = 0.3$	$98.0 \pm 0.04$	$88.1 \pm 0.22$	$-11.4 \pm 0.39$	$-16.3 \pm 0.55$
GA	$98.0 \pm 0.11$	$88.2 \pm 0.43$	$-12.6 \pm 0.48$	$-18.0 \pm 0.69$
GA + Dual Restarts	$98.2 \pm 0.01$	$88.3 \pm 0.31$	$1.4 \pm 0.07$	$2.0 \pm 0.10$
<i>Ours</i> - $\nu$ PI $\kappa_p = 9.6$	$98.3 \pm 0.09$	$88.5 \pm 0.20$	$0.1 \pm 0.01$	$0.1 \pm 0.02$

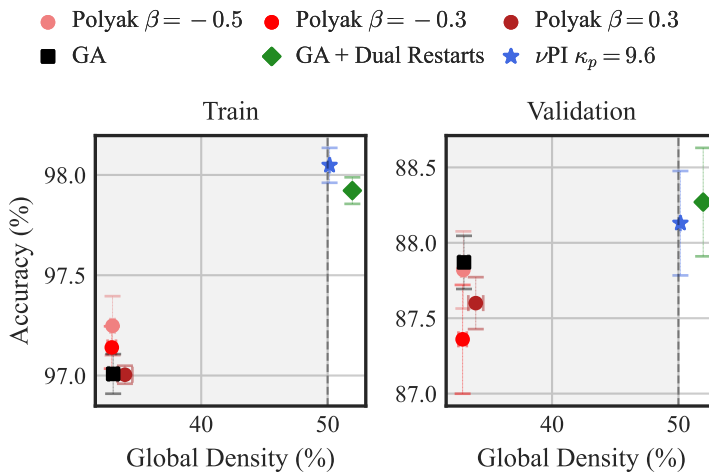


Figure 17: CIFAR10 trade-off plot for *global* sparsity under a 50% density target.

Table 4: CIFAR10 results for *global* sparsity under a 50% density target. The results in this table correspond to those in Fig. 17.

Method	Train Acc.	Test Acc.	Violation	Relative Violation
Polyak $\beta = -0.5$	$97.2 \pm 0.15$	$87.8 \pm 0.26$	$-17.0 \pm 0.17$	$-34.0 \pm 0.34$
Polyak $\beta = -0.3$	$97.1 \pm 0.11$	$87.4 \pm 0.36$	$-17.1 \pm 0.33$	$-34.2 \pm 0.66$
Polyak $\beta = 0.3$	$97.0 \pm 0.04$	$87.6 \pm 0.17$	$-16.0 \pm 0.59$	$-32.1 \pm 1.18$
GA	$97.0 \pm 0.10$	$87.9 \pm 0.18$	$-17.0 \pm 0.36$	$-33.9 \pm 0.72$
GA + Dual Restarts	$97.9 \pm 0.07$	$88.3 \pm 0.36$	$2.0 \pm 0.09$	$3.9 \pm 0.18$
<i>Ours</i> - $\nu$ PI $\kappa_p = 9.6$	$98.0 \pm 0.09$	$88.1 \pm 0.35$	$0.2 \pm 0.03$	$0.3 \pm 0.05$

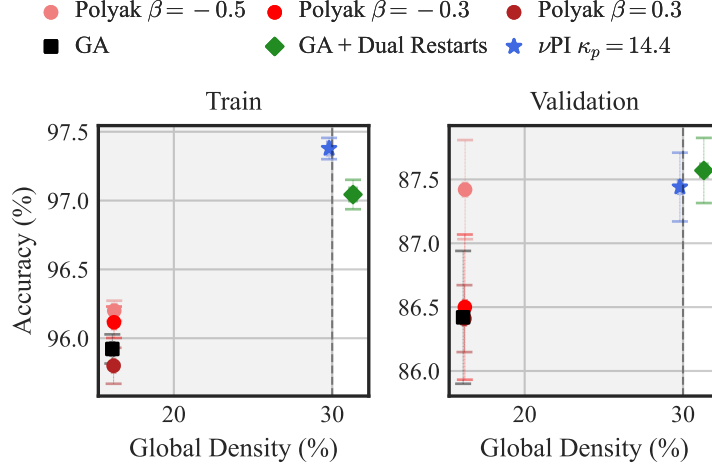


Figure 18: CIFAR10 trade-off plot for *global* sparsity under a 30% density target.

Table 5: CIFAR10 results for *global* sparsity under a 30% density target. The results in this table correspond to those in Fig. 18.

Method	Train Acc.	Test Acc.	Violation	Relative Violation
Polyak $\beta = -0.5$	$96.2 \pm 0.07$	$87.4 \pm 0.39$	$-13.8 \pm 0.13$	$-45.9 \pm 0.43$
Polyak $\beta = -0.3$	$96.1 \pm 0.11$	$86.5 \pm 0.57$	$-13.8 \pm 0.11$	$-45.9 \pm 0.36$
Polyak $\beta = 0.3$	$95.8 \pm 0.13$	$86.4 \pm 0.26$	$-13.8 \pm 0.09$	$-46.0 \pm 0.31$
GA	$95.9 \pm 0.11$	$86.4 \pm 0.52$	$-13.9 \pm 0.11$	$-46.3 \pm 0.36$
GA + Dual Restarts	$97.0 \pm 0.11$	$87.6 \pm 0.26$	$1.3 \pm 0.22$	$4.4 \pm 0.73$
<i>Ours</i> - $\nu$ PI $\kappa_p = 14.4$	$97.4 \pm 0.08$	$87.4 \pm 0.27$	$-0.2 \pm 0.11$	$-0.7 \pm 0.38$

Table 6: CIFAR10 results for *global* sparsity under a 10% density target. The results in this table correspond to those in Fig. 19.

Method	Train Acc.	Test Acc.	Violation	Relative Violation
Polyak $\beta = -0.5$	$94.6 \pm 0.17$	$85.2 \pm 0.93$	$0.7 \pm 0.13$	$7.0 \pm 1.31$
Polyak $\beta = -0.3$	$94.3 \pm 0.06$	$84.7 \pm 0.71$	$0.2 \pm 0.14$	$2.0 \pm 1.39$
Polyak $\beta = 0.3$	$92.6 \pm 0.10$	$81.9 \pm 0.77$	$-2.1 \pm 0.08$	$-21.4 \pm 0.82$
GA	$93.8 \pm 0.20$	$81.9 \pm 2.68$	$-0.9 \pm 0.09$	$-9.0 \pm 0.93$
GA + Dual Restarts	$94.3 \pm 0.08$	$85.5 \pm 0.59$	$0.4 \pm 0.03$	$3.5 \pm 0.35$
<i>Ours</i> - $\nu$ PI $\kappa_p = 1.6$	$94.1 \pm 0.08$	$83.4 \pm 1.83$	$-0.2 \pm 0.14$	$-1.6 \pm 1.44$

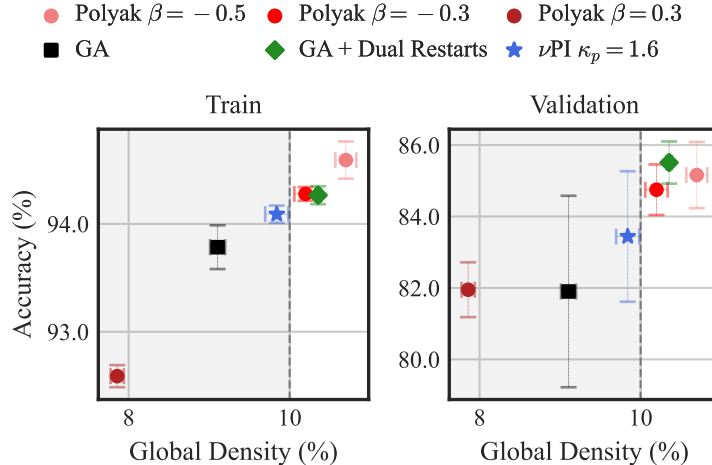

 Figure 19: CIFAR10 trade-off plot for *global* sparsity under a 10% density target.

 Table 7: CIFAR10 results for *layer-wise* sparsity under a 70% density target. GA and momentum methods overshoot to *different* values for each constraint.  **$\nu$ PI achieves the desired sparsity on all layers while achieving the highest train accuracy.** As higher density correlates to higher train accuracy, overshooting to a lower density is undesirable. All optimizers use the same step-size. *This table is the same as Table 1. We repeat it here for the reader’s convenience.*

Method	Accuracy		Violation			Relative Violation		
	Train	Test	Min	Max	Range	Min	Max	Range
Polyak $\beta = -0.5$	$91.9 \pm 0.18$	$83.6 \pm 1.40$	$-26.5 \pm 0.81$	$-7.9 \pm 0.86$	$18.9 \pm 1.31$	$-37.8 \pm 1.15$	$-11.4 \pm 1.22$	$27.0 \pm 1.87$
Polyak $\beta = -0.3$	$92.1 \pm 0.07$	$83.4 \pm 1.44$	$-27.1 \pm 0.73$	$-6.7 \pm 0.38$	$20.6 \pm 0.92$	$-38.8 \pm 1.05$	$-9.6 \pm 0.55$	$29.4 \pm 1.31$
Polyak $\beta = 0.3$	$91.9 \pm 0.20$	$82.5 \pm 1.50$	$-26.3 \pm 0.82$	$-2.3 \pm 0.69$	$24.0 \pm 0.88$	$-37.5 \pm 1.17$	$-3.2 \pm 0.99$	$34.3 \pm 1.26$
GA	$92.0 \pm 0.08$	$84.1 \pm 1.97$	$-27.8 \pm 0.49$	$-5.2 \pm 0.39$	$22.0 \pm 0.56$	$-39.6 \pm 0.70$	$-7.4 \pm 0.55$	$31.4 \pm 0.80$
GA + Dual Restarts	$95.0 \pm 0.22$	$85.3 \pm 0.61$	$-0.0 \pm 0.00$	$1.2 \pm 0.28$	$1.2 \pm 0.28$	$-0.0 \pm 0.00$	$1.8 \pm 0.40$	$1.8 \pm 0.40$
Ours - $\nu$ PI $\kappa_p = 8.0$	$95.1 \pm 0.06$	$86.2 \pm 0.46$	$-1.7 \pm 0.27$	$0.1 \pm 0.04$	$1.8 \pm 0.29$	$-2.4 \pm 0.38$	$0.2 \pm 0.06$	$2.5 \pm 0.42$

 Table 8: CIFAR10 results for *layer-wise* sparsity under a 50% density target. As higher density correlates to higher train accuracy, overshooting to a lower density is undesirable. All optimizers use the same step-size.

Method	Accuracy		Violation			Relative Violation		
	Train	Test	Min	Max	Range	Min	Max	Range
Polyak $\beta = -0.5$	$87.5 \pm 0.17$	$80.2 \pm 2.65$	$-32.6 \pm 0.95$	$-15.9 \pm 0.80$	$16.4 \pm 1.31$	$-65.1 \pm 1.89$	$-31.7 \pm 1.59$	$32.9 \pm 2.61$
Polyak $\beta = -0.3$	$87.7 \pm 0.21$	$80.3 \pm 2.13$	$-29.5 \pm 0.69$	$-15.4 \pm 0.81$	$13.6 \pm 1.23$	$-59.1 \pm 1.38$	$-30.8 \pm 1.62$	$27.2 \pm 2.46$
Polyak $\beta = 0.3$	$87.4 \pm 0.21$	$79.7 \pm 3.18$	$-30.9 \pm 0.66$	$-14.1 \pm 0.25$	$16.9 \pm 0.70$	$-61.8 \pm 1.33$	$-28.3 \pm 0.49$	$33.9 \pm 1.40$
GA	$87.6 \pm 0.12$	$77.5 \pm 3.14$	$-29.7 \pm 1.02$	$-14.2 \pm 0.60$	$14.7 \pm 1.15$	$-59.4 \pm 2.04$	$-28.3 \pm 1.20$	$29.4 \pm 2.30$
GA + Dual Restarts	$92.8 \pm 0.07$	$83.5 \pm 0.81$	$-0.0 \pm 0.01$	$1.0 \pm 0.46$	$1.0 \pm 0.46$	$-0.0 \pm 0.02$	$1.9 \pm 0.92$	$1.9 \pm 0.93$
Ours - $\nu$ PI $\kappa_p = 8.0$	$93.2 \pm 0.06$	$83.6 \pm 0.87$	$-1.5 \pm 0.13$	$0.1 \pm 0.08$	$1.6 \pm 0.19$	$-2.9 \pm 0.26$	$0.2 \pm 0.16$	$3.2 \pm 0.37$

## H. Additional Experiments

In this section, we include additional experimental results on the sparsity-constrained task. We analyze the dynamics of the multiplier throughout training in Appx. H.1, and conduct ablation studies on  $\kappa_p$  for  $\nu$ PI, the momentum coefficients of POLYAK and NESTEROV, and the step-size of ADAM.

### H.1. Dynamics

The dynamics shown in Fig. 20 illustrate the change of the constraint violation and multipliers throughout optimization. We observe that GA, POLYAK, and ADAM quickly lead to overshoot into the feasible region, leading to overly sparse models.

Table 9: CIFAR10 results for *layer-wise* sparsity under a 30% density target. As higher density correlates to higher train accuracy, overshooting to a lower density is undesirable. All optimizers use the same step-size.

Method	Accuracy		Min	Violation		Range	Relative Violation		
	Train	Test		Max	Min		Max	Range	
Polyak $\beta = -0.5$	81.8 $\pm$ 0.19	63.5 $\pm$ 18.58	-25.2 $\pm$ 1.54	-17.0 $\pm$ 0.37	8.4 $\pm$ 1.73	-84.1 $\pm$ 5.12	-56.8 $\pm$ 1.23	28.0 $\pm$ 5.77	
Polyak $\beta = -0.3$	82.1 $\pm$ 0.54	63.3 $\pm$ 8.65	-25.1 $\pm$ 1.15	-16.4 $\pm$ 0.38	8.7 $\pm$ 0.91	-83.5 $\pm$ 3.84	-54.7 $\pm$ 1.27	29.0 $\pm$ 3.04	
Polyak $\beta = 0.3$	81.8 $\pm$ 0.32	72.7 $\pm$ 3.36	-25.1 $\pm$ 2.12	-17.5 $\pm$ 0.24	7.4 $\pm$ 2.12	-83.6 $\pm$ 7.07	-58.5 $\pm$ 0.79	24.8 $\pm$ 7.07	
GA	81.8 $\pm$ 0.44	72.7 $\pm$ 4.40	-24.8 $\pm$ 1.11	-17.0 $\pm$ 0.60	8.5 $\pm$ 1.22	-82.5 $\pm$ 3.69	-56.7 $\pm$ 1.99	28.2 $\pm$ 4.07	
GA + Dual Restarts	89.7 $\pm$ 0.23	82.9 $\pm$ 2.59	-0.0 $\pm$ 0.00	0.9 $\pm$ 0.33	0.9 $\pm$ 0.33	-0.0 $\pm$ 0.01	3.0 $\pm$ 1.10	3.0 $\pm$ 1.10	
Ours - $\nu$ PI $\kappa_p = 12.0$	89.8 $\pm$ 0.11	82.0 $\pm$ 2.45	-0.3 $\pm$ 0.13	0.3 $\pm$ 0.03	0.6 $\pm$ 0.12	-0.8 $\pm$ 0.42	1.0 $\pm$ 0.11	2.1 $\pm$ 0.39	

Table 10: CIFAR10 results for *layer-wise* sparsity under a 10% density target. As higher density correlates to higher train accuracy, overshooting to a lower density is undesirable. All optimizers use the same step-size.

Method	Accuracy		Min	Violation		Range	Relative Violation		
	Train	Test		Max	Min		Max	Range	
Polyak $\beta = -0.5$	71.3 $\pm$ 0.61	61.0 $\pm$ 9.50	-10.0 $\pm$ 0.14	-5.9 $\pm$ 0.47	4.0 $\pm$ 0.48	-100.0 $\pm$ 1.36	-58.7 $\pm$ 4.74	40.5 $\pm$ 4.83	
Polyak $\beta = -0.3$	70.9 $\pm$ 0.60	49.5 $\pm$ 16.33	-10.0 $\pm$ 0.01	-5.9 $\pm$ 0.60	4.1 $\pm$ 0.60	-100.0 $\pm$ 0.11	-58.9 $\pm$ 5.97	41.1 $\pm$ 5.95	
Polyak $\beta = 0.3$	69.2 $\pm$ 0.71	56.3 $\pm$ 15.05	-10.0 $\pm$ 0.02	-6.7 $\pm$ 0.06	3.3 $\pm$ 0.08	-100.0 $\pm$ 0.15	-67.3 $\pm$ 0.65	32.7 $\pm$ 0.79	
GA	71.0 $\pm$ 0.32	49.6 $\pm$ 11.1	-10.0 $\pm$ 0.19	-6.1 $\pm$ 0.25	3.9 $\pm$ 0.42	-100.0 $\pm$ 1.91	-61.2 $\pm$ 2.54	38.8 $\pm$ 4.24	
GA + Dual Restarts	83.1 $\pm$ 0.27	73.1 $\pm$ 4.87	-0.0 $\pm$ 0.00	1.6 $\pm$ 0.14	1.6 $\pm$ 0.14	-0.0 $\pm$ 0.02	16.1 $\pm$ 1.39	16.1 $\pm$ 1.40	
Ours - $\nu$ PI $\kappa_p = 12.0$	81.4 $\pm$ 0.39	42.8 $\pm$ 14.54	-1.9 $\pm$ 0.34	0.9 $\pm$ 0.46	3.2 $\pm$ 0.72	-19.1 $\pm$ 3.41	9.3 $\pm$ 4.62	31.7 $\pm$ 7.17	

As training progresses however, these methods move closer to the boundary of the feasible region, reversing the initial overshoot. This recovery is most notorious for ADAM, whose multiplier decreases quickly after feasibility. GA with dual restarts sets the value of the multiplier to zero as soon as feasibility is achieved, thus preventing an incursion into the feasible set.  $\nu$ PI produces well-behaved multipliers and successfully avoids constraint overshoot.

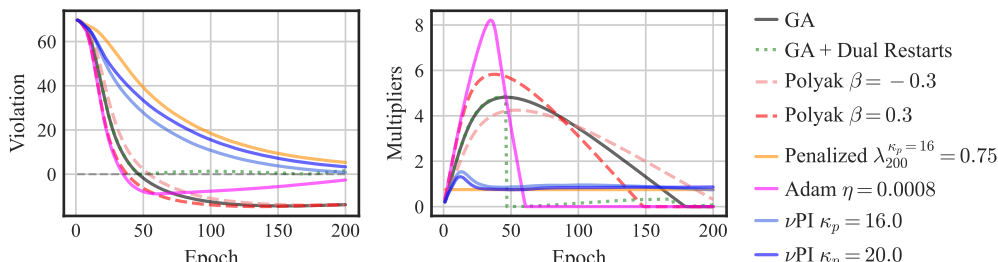


Figure 20: Dynamics plot for global sparsity under a 30% density target.

Note that for this sparsity task, it is reasonable to expect that the constraint is active at the constrained optimum since more model capacity correlates with better performance. However, note that  $\nu$ PI is the only method that provides a non-zero estimate of the Lagrange multiplier. The usefulness of Lagrange multiplier estimates is highlighted in Fig. 21.

Figure 21 considers unconstrained  $L_0$ -regularization experiments. We use the final value of the multipliers corresponding to  $\nu$ PI ( $\kappa_p = 16$ ) and  $\nu$ PI ( $\kappa_p = 20$ ) runs as the (fixed) penalty coefficient for 200 epochs in the penalized formulation of the problem, akin to Louizos et al. (2018). We also include an experiment using the multiplier estimate from GA (equal to zero).

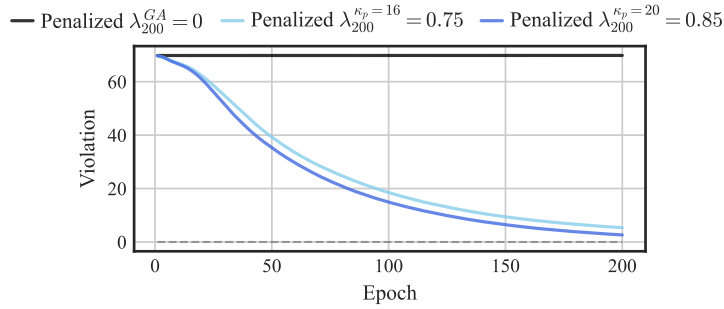


Figure 21: Dynamics plot for global sparsity under a 30% density target.

Unsurprisingly, the run with the multiplier of zero leads to 100% density, since the sparsity penalty does not exert any influence during training. In contrast, the runs with the  $\nu$ PI multiplier estimates not only lead to sparse models but are also very close to the desired model density by the end of training. This is remarkable since the problem we are solving is nonconvex, and optimal Lagrange multiplier values may not even exist.

### H.2. Ablation on the value of $\kappa_p$

In this section, we fix  $\kappa_i$  for  $\nu$ PI and ablate on the hyperparameter  $\kappa_p$  for two sparsity levels. The results are presented in Fig. 22 and Tables 11 and 12.

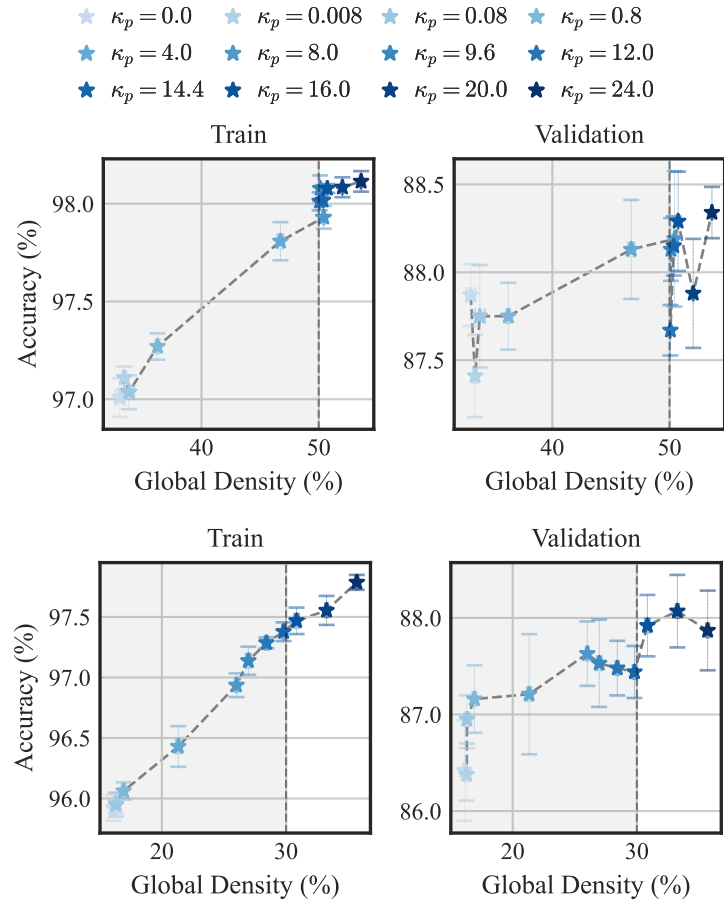


Figure 22: Ablation on trade-offs achievable by  $\nu$ PI under global density targets of 50% (top) and 30% (bottom).

We see that a larger  $\kappa_p$  leads to more damping and less overshoot. Note that there is a strong correlation between training accuracy and model density. Hence, it is important to be able to control overshoot in sparsity constraints and take advantage of the maximum allowed density for the sake of accuracy. There is a range of  $\kappa_p$  that can achieve such desired sparsity. The same trend roughly extends to validation accuracy (with some caveats due to generalization errors).

Table 11: Ablation on the  $\kappa_p$  hyperparameter for a CIFAR10 task with a global density target of  $\epsilon = 50\%$ .  $\kappa_p$  **monotonically controls the degree of damping and constraint overshoot.**

$\nu$ PI $\kappa_p$	Train Acc.	Test Acc.	Violation	Relative Violation
0	97.0 ± 0.10	87.9 ± 0.18	-17.0 ± 0.36	-33.9 ± 0.72
0.008	97.1 ± 0.06	87.4 ± 0.23	-16.6 ± 0.14	-33.2 ± 0.28
0.08	97.0 ± 0.09	87.7 ± 0.29	-16.2 ± 0.42	-32.4 ± 0.84
0.8	97.3 ± 0.07	87.7 ± 0.19	-13.8 ± 0.13	-27.5 ± 0.27
4	97.8 ± 0.10	88.1 ± 0.28	-3.3 ± 0.23	-6.5 ± 0.46
8	97.9 ± 0.06	88.2 ± 0.39	0.4 ± 0.03	0.9 ± 0.06
9.6	98.1 ± 0.07	88.1 ± 0.18	0.1 ± 0.02	0.2 ± 0.04
12	98.0 ± 0.05	87.7 ± 0.14	0.1 ± 0.01	0.2 ± 0.03
14.4	98.0 ± 0.08	88.2 ± 0.17	0.4 ± 0.02	0.7 ± 0.05
16	98.1 ± 0.02	88.3 ± 0.28	0.7 ± 0.02	1.5 ± 0.04
20	98.1 ± 0.05	87.9 ± 0.31	2.0 ± 0.03	4.0 ± 0.07
24	98.1 ± 0.05	88.3 ± 0.15	3.6 ± 0.03	7.2 ± 0.06

Table 12: Ablation on the  $\kappa_p$  hyperparameter for a CIFAR10 task with a global density target of  $\epsilon = 30\%$ .

$\nu$ PI $\kappa_p$	Train Acc.	Test Acc.	Violation	Relative Violation
0	95.9 ± 0.11	86.4 ± 0.52	-13.9 ± 0.11	-46.3 ± 0.36
0.008	96.0 ± 0.08	86.4 ± 0.27	-13.7 ± 0.13	-45.7 ± 0.43
0.08	95.9 ± 0.10	86.9 ± 0.25	-13.7 ± 0.18	-45.6 ± 0.59
0.8	96.1 ± 0.07	87.2 ± 0.35	-13.1 ± 0.13	-43.6 ± 0.45
4	96.4 ± 0.17	87.2 ± 0.62	-8.7 ± 0.16	-28.9 ± 0.54
8	96.9 ± 0.10	87.6 ± 0.33	-4.0 ± 0.10	-13.3 ± 0.32
9.6	97.1 ± 0.12	87.5 ± 0.45	-3.0 ± 0.18	-10.1 ± 0.60
12	97.3 ± 0.04	87.5 ± 0.28	-1.6 ± 0.11	-5.3 ± 0.37
14.4	97.4 ± 0.08	87.4 ± 0.27	-0.2 ± 0.11	-0.7 ± 0.38
16	97.5 ± 0.11	87.9 ± 0.32	0.8 ± 0.17	2.8 ± 0.57
20	97.6 ± 0.12	88.1 ± 0.38	3.3 ± 0.11	10.9 ± 0.36
24	97.8 ± 0.06	87.9 ± 0.41	5.7 ± 0.11	19.0 ± 0.36

### H.3. ADAM

We also experimented with a range of learning choices for ADAM to explore their effect on constraint satisfaction and overshoot. The results are shown in Fig. 23, and Tables 13 and 14.

We observe that the influence of ADAM’s learning on the constraint overshoot is not monotonic. When the step-size is small, ADAM runs do not satisfy the constraint at the end of training. As the step-size increases, satisfaction is achieved together with varying degrees of overshoot into the feasible region. A range of larger step-sizes that lie at the sweet spot of almost exact constraint satisfaction.

The sensitivity and non-monotonicity of the step-size make the tuning of the step-size hyperparameter for ADAM challenging. Note that we restricted our experiments to the default EMA coefficients for ADAM following PyTorch:  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ .

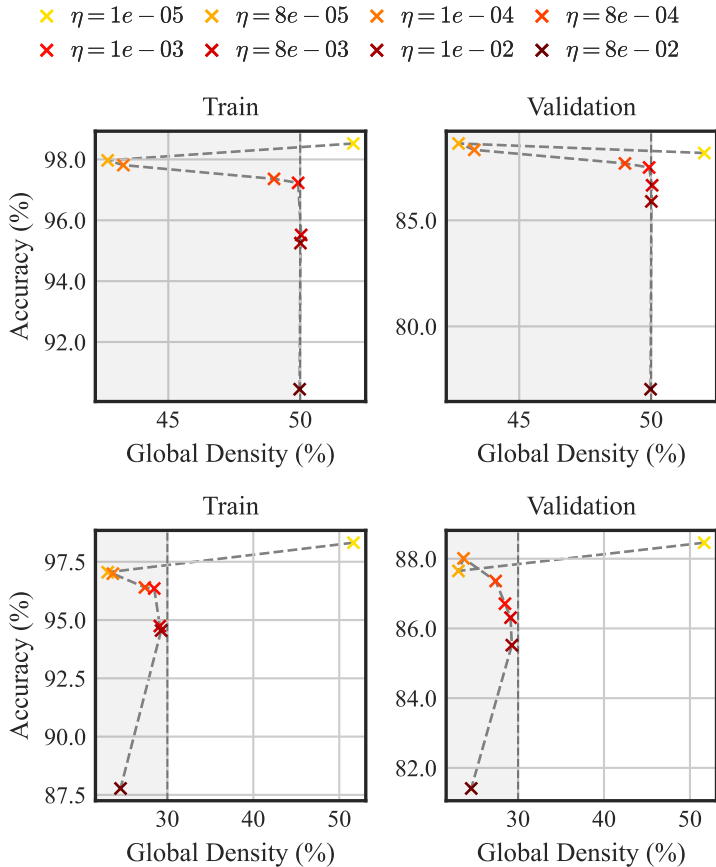


Figure 23: Ablation on the density-accuracy trade-offs achievable by ADAM under global density targets of 50% (top) and 30% (bottom).

Table 13: Ablation on the step-size hyperparameter for ADAM on a CIFAR10 task with a global density target of  $\epsilon = 50\%$ .

Adam $\eta$	Train Acc.	Test Acc.	Violation
$1 \cdot 10^{-5}$	98.52	88.17	2.02
$8 \cdot 10^{-5}$	97.97	88.62	-7.30
$1 \cdot 10^{-4}$	97.81	88.32	-6.70
$8 \cdot 10^{-4}$	97.36	87.68	-0.99
$1 \cdot 10^{-3}$	97.23	87.49	-0.08
$8 \cdot 10^{-3}$	95.52	86.65	0.04
$1 \cdot 10^{-2}$	95.25	85.89	0.01
$8 \cdot 10^{-2}$	90.45	77.04	-0.02

#### H.4. Momentum

We carried out similar ablations on the momentum coefficient of POLYAK and NESTEROV using both positive and negative values. The results are shown in Fig. 24, and Tables 15 and 16. We observe significant overshoot into the feasible region for all attempted values, compared to the desired target density of 30%.

Table 14: Ablation on the step-size hyperparameter for ADAM on a CIFAR10 task with a global density target of  $\epsilon = 30\%$ .

Adam $\eta$	Train Acc.	Test Acc.	Violation
$1 \cdot 10^{-5}$	98.32	88.46	21.66
$8 \cdot 10^{-5}$	97.05	87.65	-6.94
$1 \cdot 10^{-4}$	96.99	88.01	-6.34
$8 \cdot 10^{-4}$	96.40	87.36	-2.61
$1 \cdot 10^{-3}$	96.35	86.71	-1.53
$8 \cdot 10^{-3}$	94.74	86.31	-0.88
$1 \cdot 10^{-2}$	94.55	85.52	-0.73
$8 \cdot 10^{-2}$	87.78	81.41	-5.45

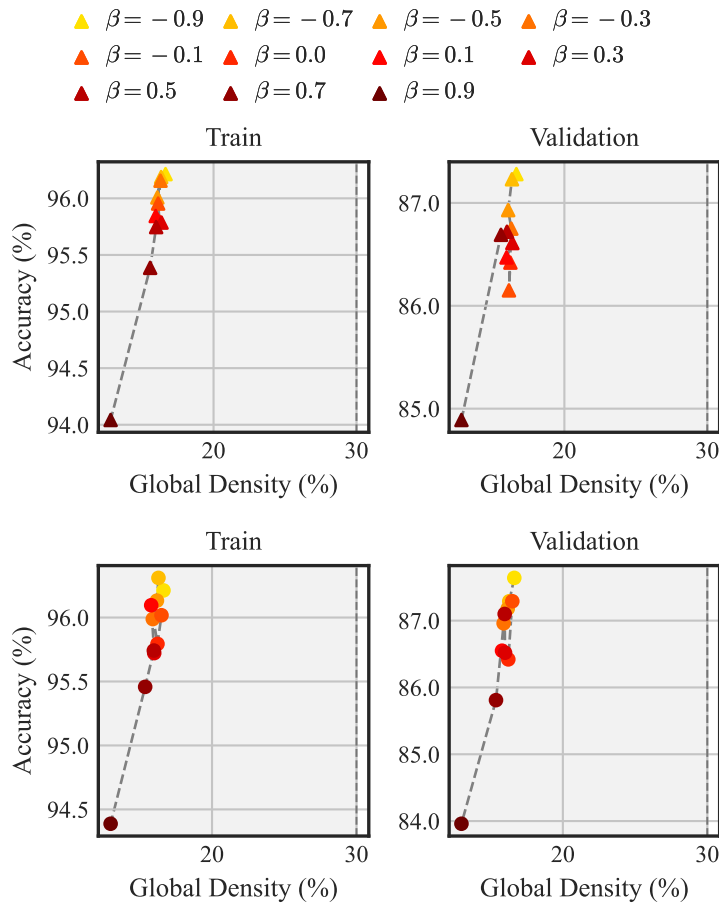


Figure 24: Trade-off plot under a 30% global density target for NESTEROV (top) and POLYAK (bottom) momentum.



Table 15: Ablation on the momentum hyperparameter for NESTEROV on a CIFAR10 task with a 30% global density target.

Nesterov $\beta$	Train Acc.	Test Acc.	Violation
-0.9	96.21	87.28	-13.36
-0.7	96.19	87.23	-13.67
-0.5	96.01	86.93	-13.93
-0.3	96.16	86.75	-13.71
-0.1	95.95	86.15	-13.88
0.0	95.79	86.42	-13.78
0.1	95.84	86.47	-14.05
0.3	95.79	86.61	-13.64
0.5	95.75	86.72	-14.02
0.7	95.39	86.69	-14.44
0.9	94.04	84.89	-17.20

Table 16: Ablation on the momentum hyperparameter for POLYAK on a CIFAR10 task with a 30% global density target.

Polyak $\beta$	Train Acc.	Test Acc.	Violation
-0.9	96.21	87.64	-13.36
-0.7	96.31	87.29	-13.71
-0.5	96.13	87.18	-13.82
-0.3	95.99	86.96	-14.10
-0.1	96.02	87.29	-13.50
0.0	95.79	86.42	-13.78
0.1	96.10	86.55	-14.21
0.3	95.72	86.52	-14.00
0.5	95.74	87.10	-14.03
0.7	95.46	85.81	-14.63
0.9	94.39	83.96	-17.02