
OpReg-Boost: Learning to Accelerate Online Algorithms with Operator Regression

Anonymous Authors¹

Abstract

In this paper we propose a learning-based acceleration scheme for online optimization. novel operator regression ...

1. Introduction

.... problem formulation, literature review, contribution statement

2. Operator Regression

2.1. Problem formulation

Consider the following *convex, time-varying* optimization problem

$$\mathbf{x}^*(t) \in \arg \min f(\mathbf{x}; t) + g(\mathbf{x}; t) \quad (1)$$

where $f : \mathbb{R}^n \times \mathbb{R}_+ \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \times \mathbb{R}_+ \rightarrow \mathbb{R} \cup \{+\infty\}$ are closed, convex and proper functions, optionally with $g \equiv 0$. In particular, we are interested in solving the following sequence of static problems derived from sampling (1) at the times $\{t_k\}_{k \in \mathbb{N}}$, $t_{k+1} - t_k = T_s$:

$$\mathbf{x}^*(t) \in \arg \min f(\mathbf{x}; t_k) + g(\mathbf{x}; t_k). \quad (2)$$

Problem (2) is convex but not strongly so, which means that the performance of online algorithms applied to it is worse than the performance attainable for strongly convex problems. As an example, consider the case $g \equiv 0$, and apply an online gradient descent: if the problem is convex the (static) regret is $O(\sqrt{T})$, with T the number of sampling times, while if the problem is strongly convex the regret is $O(\log(T))$, section 3.1 in (Hazan, 2016). [Instead of citing Hazan, we should derive a result in an appendix, for example for the fixed point residual](#)

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

As we can see from the results just mentioned, in general in online optimization it is not possible to achieve zero regret, due to the dynamic nature of the problem. However, for strongly convex problems smaller regrets can be achieved.

The goal then is to design learning techniques that allow to achieve strongly convex-like performance while tracking a solution trajectory of the problem (2) with good accuracy.

2.2. Operator regression

First of all, we define the observations

$$\mathbf{y}_i = \mathcal{T}_f \mathbf{x}_i + \mathbf{e}_i, \quad i \in [D]$$

of the operator we want to regularize. Now, using Fact 2.2 in (Ryu et al., 2020) we know that an operator \mathcal{T} is ζ -contractive interpolable if and only if it satisfies

$$\|\mathcal{T} \mathbf{x}_i - \mathcal{T} \mathbf{x}_j\|^2 \leq \zeta^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2, \quad \forall i, j \in [D], i \neq j.$$

Therefore, we can define the following constrained regression problem:

$$\begin{aligned} \hat{\mathbf{t}}_i &= \arg \min_{\mathbf{t}_i \in \mathbb{R}^n} \frac{1}{2} \sum_{i \in [D]} \|\mathbf{t}_i - \mathbf{y}_i\|^2 \\ \text{s.t. } \|\mathbf{t}_i - \mathbf{t}_j\|^2 &\leq \zeta^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2 \quad \forall i, j \in [D], i \neq j \end{aligned} \quad (3)$$

where again the cost function is a least square terms on the observations and the constraints enforce contractiveness.

Auto-tuning contraction We can modify (3) by including $w = \zeta^2$ as an unknown of the regression problem, which becomes

$$\begin{aligned} \hat{\mathbf{t}}_i &= \arg \min_{\mathbf{t}_i \in \mathbb{R}^n} \frac{1}{2} \sum_{i \in [D]} \|\mathbf{t}_i - \mathbf{y}_i\|^2 + \frac{c}{2} w^2 \\ \text{s.t. } \|\mathbf{t}_i - \mathbf{t}_j\|^2 - w \|\mathbf{x}_i - \mathbf{x}_j\|^2 &\leq 0 \quad \forall i, j \in [D], i \neq j \end{aligned} \quad (4)$$

this way the contraction constant does not need to be specified, and is *auto-tuned* by the regression.

PRS-based solver brief description of basic characteristics, reference to the appendix

3. OpReg-Boost

Let now $\mathcal{T}_k : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a solver for the convex problem observed at time t_k , and assume that the operator has a term $\mathcal{T}_{f,k}$ depending exclusively on f (e.g. we have $\mathcal{T}_k = \mathcal{T}'_k \circ \mathcal{T}_{f,k}$). Since the sampled problem is convex, the solver \mathcal{T}_k is non-expansive (actually, averaged, to guarantee convergence).

The idea then is to try and learn an approximation of \mathcal{T}_k that is *contractive* rather than non-expansive.

Again, w.l.o.g. we are only interested in learning the term $\mathcal{T}_{f,k}$ that depends on f . The algorithm can be described as follows: at each time t_k do:

1. sample a new problem, that is, observe $f(\mathbf{x}; t_k)$ and $g(\mathbf{x}; t_k)$;
2. approximate the term $\mathcal{T}_{f,k}$ of the solver with a contractive one, and
3. apply the resulting solver. For example if $\mathcal{T}_k = \mathcal{T}'_k \circ \mathcal{T}_{f,k}$, then we apply $\mathcal{T}_k = \mathcal{T}'_k \circ \hat{\mathcal{T}}_{f,k}$ where $\hat{\mathcal{T}}_{f,k}$ is the learned operator.

The learning step is performed using the (novel) constrained operator regression OpReg described in the following.

In particular we have the following algorithm [make this pseudo code](#) For $\ell \in 0, 1, \dots, M$ (for some $M \in \mathbb{N}$):

- let \mathbf{x}_k^ℓ be the current approximate solution, we choose $D - 1$ points around it, e.g.

$$\mathbf{x}_i = \mathbf{x}_k^\ell + \mathbf{p}_i, \quad i = 2, \dots, D, \quad \mathbf{p}_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$$

and set $\mathbf{x}_1 = \mathbf{x}_k^\ell$;

- we compute the data $\mathbf{t}_i = \mathcal{T}_{f,k} \mathbf{x}_i$, $i = 1, \dots, D$, for example $\mathbf{t}_i = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i; t_k)$;
- we solve the OpReg, and use the approximate operator at \mathbf{x}_k^ℓ , that is \mathbf{t}_1 , in the chosen solver \mathcal{T}_k :

$$\mathbf{x}_k^{\ell+1} = \text{prox}_{\alpha g}(\mathbf{t}_1)$$

where α is a step-size.

3.1. Interpolated version

4. Numerical Results

4.1. Simulations set-up

We consider the following time-varying problem:

$$\mathbf{x}^*(t_k) = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}(t_k)\|^2 + w \|\mathbf{x}\|_1 \quad (5)$$

with $n = 10$, \mathbf{A} matrix with maximum and minimum (non-zero) eigenvalues $\sqrt{L} = 10^8$, $\sqrt{\mu} = 1$, and with rank 5; $\mathbf{y}(t_k)$ has sinusoidal components with 3 zero components. Due to \mathbf{A} being rank deficient, the cost f is convex but not strongly so.

4.2. Results

References

- Hazan, E. Introduction to Online Convex Optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- Ryu, E. K., Taylor, A. B., Bergeling, C., and Giselsson, P. Operator Splitting Performance Estimation: Tight Contraction Factors and Optimal Parameter Selection. *SIAM Journal on Optimization*, 30(3):2251–2271, 2020.
- Simonetto, A. Smooth Strongly Convex Regression. In *2020 28th European Signal Processing Conference (EU-SIPCO)*, pp. 2130–2134, Amsterdam, January 2021. IEEE.

A. PRS-Based QCQP solver

In this section we present a solver for OpReg which can be efficiently parallelized, inspired by the approach in (Simonetto, 2021).

The idea is as follows: each pair of data points $i, j \in [D]$, $i \neq j$, gives rise to one constraint, for a total of $D(D-1)/2$ constraints. We define the following set of pairs

$$\mathcal{V} = \{e = (i, j) \mid i, j \in [D], i < j\}$$

which are ordered (that is, for example we take (1, 2) and not (2, 1), to avoid counting constraints twice). Clearly to each pair $e = (i, j)$ corresponds the constraint $\|\mathbf{t}_i - \mathbf{t}_j\|^2 \leq \zeta^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2$.

Let now $\mathbf{t}_{i,e}$ and $\mathbf{t}_{j,e}$ be copies of \mathbf{t}_i and \mathbf{t}_j associated to the e -th constraint; then we can equivalently reformulate the OpReg (3) as

$$\min_{\mathbf{t}_{i,e}, \mathbf{t}_{j,e}} \frac{1}{2(D-1)} \sum_{e \in \mathcal{V}} \left\| \begin{bmatrix} \mathbf{t}_{i,e} \\ \mathbf{t}_{j,e} \end{bmatrix} - \begin{bmatrix} \mathbf{y}_i \\ \mathbf{y}_j \end{bmatrix} \right\|^2 \quad (6a)$$

$$\text{s.t. } \|\mathbf{t}_{i,e} - \mathbf{t}_{j,e}\|^2 \leq \zeta^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2 \quad (6b)$$

$$\mathbf{t}_{i,e} = \mathbf{t}_{i,e'} \quad \forall e, e' \mid i \sim e, e'. \quad (6c)$$

Clearly (6) is a strongly convex problem with convex constraints defined in the variables $\mathbf{t}_{i,e}$.

A.1. PRS solver

Let ξ be the vector stacking all the $t_{i,e}$, then the problem is equivalent to

$$\min_{\xi} f(\xi) + g(\xi)$$

where

$$f(\xi) = \frac{1}{2(D-1)} \|\xi - \mathbf{y}\|^2 + f_1(\xi)$$

with f_1 the indicator function imposing (6b) and g the indicator function imposing the “consensus” constraints (6c). The problem can then be solved using the Peaceman-Rachford splitting (PRS) characterized by the following updates $\ell \in \mathbb{N}$:

$$\xi^\ell = \text{prox}_{\rho f}(\mathbf{z}^\ell) \quad (7a)$$

$$\mathbf{v}^\ell = \text{prox}_{\rho g}(2\xi^\ell - \mathbf{z}^\ell) \quad (7b)$$

$$\mathbf{z}^{\ell+1} = \mathbf{z}^\ell + \mathbf{v}^\ell - \xi^\ell. \quad (7c)$$

The proximal of g corresponds to the projection onto the consensus space, and thus can be characterized simply by

$$\mathbf{v}_{i,e}^\ell = \frac{1}{D-1} \sum_{e' | i \sim e'} (2t_{i,e'}^\ell - \mathbf{z}_{e'}^\ell).$$

Regarding the proximal of f , it is clear that f is separable, in the sense that it can be written as

$$f(\xi) = \sum_{e \in \mathcal{V}} \left[\frac{1}{2(D-1)} \left\| \begin{bmatrix} t_{i,e} \\ t_{j,e} \end{bmatrix} - \begin{bmatrix} \mathbf{y}_i \\ \mathbf{y}_j \end{bmatrix} \right\|^2 + \iota_e(t_{i,e}, t_{j,e}) \right]$$

where ι_e denotes the indicator function of (6b). Therefore, the update (7a) can be solved by solving (possibly in parallel) the problems

$$\begin{aligned} (t_{i,e}, t_{j,e}) = \arg \min & \left\{ \frac{1}{2(D-1)} \left\| \begin{bmatrix} t_{i,e} \\ t_{j,e} \end{bmatrix} - \begin{bmatrix} \mathbf{y}_i \\ \mathbf{y}_j \end{bmatrix} \right\|^2 + \frac{1}{2\rho} \left\| \begin{bmatrix} t_{i,e} \\ t_{j,e} \end{bmatrix} - \mathbf{z}_e^\ell \right\|^2 \right\} \\ \text{s.t.} \quad & \|t_{i,e} - t_{j,e}\|^2 \leq \zeta^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2. \end{aligned} \quad (8)$$

A.2. Local updates

The problems (8) are quadratic programs with quadratic constraints, that is, they can be written in the form

$$\min_{\xi} \frac{1}{2} \xi^\top P_0 \xi + \langle q_0, \xi \rangle \quad (9a)$$

$$\text{s.t.} \quad \frac{1}{2} \xi^\top P_1 \xi + \langle q_1, \xi \rangle + r_1 \leq 0. \quad (9b)$$

In particular, for the cost function we have

$$P_0 = \left(\frac{1}{D-1} + \frac{1}{\rho} \right) I_{2n}, \quad q_0 = - \left(\frac{1}{D-1} \begin{bmatrix} \mathbf{y}_i \\ \mathbf{y}_j \end{bmatrix} + \frac{1}{\rho} \mathbf{z}_e^\ell \right)$$

and for the constraint

$$P_1 = 2 \begin{bmatrix} I_n & -I_n \\ -I_n & I_n \end{bmatrix}, \quad q_1 = 0_{2n}, \quad r_1 = -\zeta^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2.$$

B. Interpolation Using MAP