

Adaptive Three Operator Splitting

Fabian Pedregosa^{†‡}, Gauthier Gidel^{*}

[†]UC Berkeley [‡]ETH Zurich ^{*}Mila and DIRO, Université de Montréal



Problem Setting

Goal: solve optimization problems of the form

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{x}) + g(\mathbf{x}) + h(\mathbf{x}), \quad (\text{OPT})$$

with access to ∇f , $\text{prox}_{\gamma g} \stackrel{\text{def}}{=} \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}\|^2$, $\text{prox}_{\gamma h}$.

- Optimization problems with smooth + non-smooth objectives are ubiquitous in machine learning.
- Many complex penalties can be written as sum of proximal: overlapping group lasso, ℓ_1 trend filtering, isotonic constraints, total variation, intersection of constraints, etc.

Three Operator Splitting

The **Three Operator Splitting** (TOS) is a recently proposed method to solve (OPT) (Davis and Yin 2017). Iterates on d -dim vector \mathbf{y} :

$$\mathbf{z} = \text{prox}_{\gamma h}(\mathbf{y}), \quad \mathbf{x} = \text{prox}_{\gamma g}(2\mathbf{y} - \mathbf{z} - \gamma \nabla f(\mathbf{z})), \quad \mathbf{y}^+ = \mathbf{y} - \mathbf{z} + \mathbf{x}$$

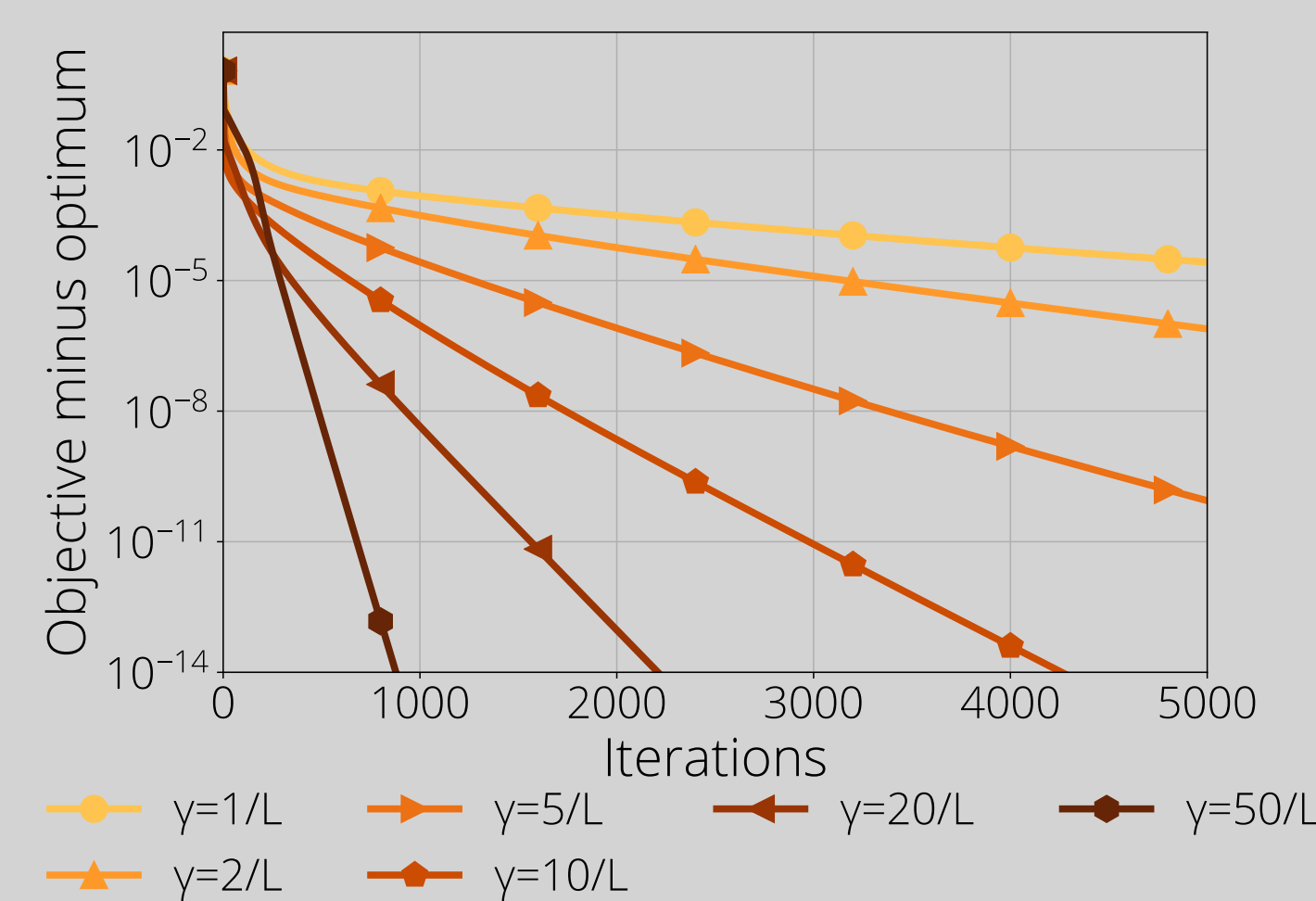
- Can be used to solve problems with an arbitrary number of proximal terms

$$\min_{\mathbf{x}} \varphi(\mathbf{x}) + \sum_{j=1}^k h_j(\mathbf{x}) \equiv \min_{\mathbf{X}} \underbrace{\varphi(\bar{\mathbf{X}})}_{=f(\mathbf{X})} + \underbrace{\sum_{j=1}^k h_j(\mathbf{X}_j)}_{=h(\mathbf{X})} + \underbrace{\iota\{\mathbf{X}_1 = \dots = \mathbf{X}_k\}}_{=g(\mathbf{X})}.$$

- Convergence has been proven for step-sizes $\gamma < 2/L$, with L = Lipschitz constant of ∇f .

- Often best step-size is orders of magnitude larger than theoretical one.

- Can we design a variant that computes step-size based exclusively on local information of the objective?



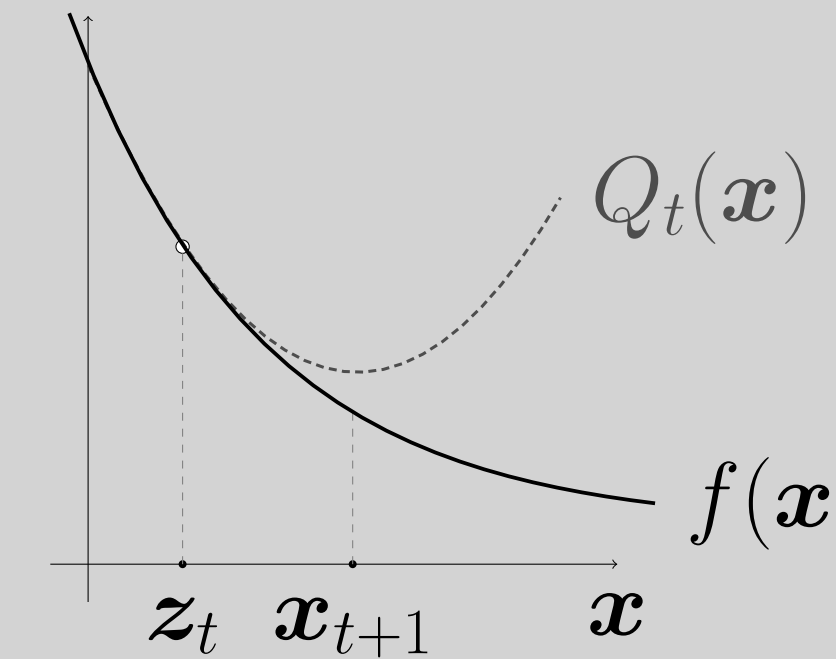
Contributions

- Three operator splitting with step-size adaptive to local information.
- Step size grows to largest admissible step-size, no hyperparameter to tune.
- Same theoretical guarantees than fixed step-size variant.
- New analysis of TOS based on saddle-point suboptimality: tighter rates and simpler proof.

Adaptive Three Operator Splitting

Key idea for adaptive step-size strategy:

- Construct surrogate quadratic of f at \mathbf{z}_t :
 $Q_t(\mathbf{x}, \gamma) \stackrel{\text{def}}{=} f(\mathbf{z}_t) + \langle \nabla f(\mathbf{z}_t), \mathbf{x} - \mathbf{z}_t \rangle + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{z}_t\|^2$.
- Choose γ_t such that $Q_t(\mathbf{x}_{t+1}, \gamma_t)$ is upper bound of $f(\mathbf{x}_{t+1})$.



Algorithm

1. Start with optimistic step-size γ_t and decrease it until:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{z}_t) + \langle \nabla f(\mathbf{z}_t), \mathbf{x}_{t+1} - \mathbf{z}_t \rangle + \frac{1}{2\gamma_t} \|\mathbf{x}_{t+1} - \mathbf{z}_t\|^2$$

with $\mathbf{x}_{t+1} = \text{prox}_{\gamma_t g}(\mathbf{z}_t - \gamma_t \nabla f(\mathbf{z}_t) + \mathbf{u}_t)$

2. Run the following updates selected step-size:

$$\mathbf{z}_{t+1} = \text{prox}_{\gamma h}(\mathbf{x}_{t+1} + \gamma_t \mathbf{u}_t), \quad \mathbf{u}_{t+1} = \mathbf{u}_t + (\mathbf{x}_{t+1} - \mathbf{z}_{t+1})$$

Analysis Framework

⚠ Iterates are not always feasible, objective can be $+\infty$. ⚠

A better gap function emerges from the saddle-point formulation

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + g(\mathbf{x}) + h(\mathbf{x}) = \min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{u} \in \mathbb{R}^d} \underbrace{f(\mathbf{x}) + g(\mathbf{x}) + \langle \mathbf{x}, \mathbf{u} \rangle - h^*(\mathbf{u})}_{\stackrel{\text{def}}{=} \mathcal{L}(\mathbf{x}, \mathbf{u})}$$

- $\mathcal{L}(\mathbf{x}_t, \mathbf{u}_t) < \infty$
- $\min_{\mathbf{u}} \mathcal{L}(\mathbf{x}_t, \mathbf{u}) = f(\mathbf{x}_t) + g(\mathbf{x}_t) + h(\mathbf{x}_t)$
- $\mathcal{L}(\mathbf{x}^*, \mathbf{u}) - \mathcal{L}(\mathbf{x}, \mathbf{u}^*) \leq 0$ for all $(\mathbf{x}, \mathbf{u}) \iff (\mathbf{x}^*, \mathbf{u}^*)$ is saddle point of $\mathcal{L} \iff \mathbf{x}^*$ is solution (OPT) and \mathbf{u}^* minimizes $(f+g)^*(-\mathbf{u}) + h^*(\mathbf{u})$.

Convergence Analysis

Let f, g, h be convex, and f also differentiable with Lipschitz gradient.

We define $s_t \stackrel{\text{def}}{=} \sum_{i=0}^{t-1} \gamma_i$, $\bar{\mathbf{x}}_t \stackrel{\text{def}}{=} \left(\sum_{i=0}^{t-1} \gamma_i \mathbf{x}_{i+1} \right) / s_t$, $\bar{\mathbf{u}}_t \stackrel{\text{def}}{=} \left(\sum_{i=0}^{t-1} \gamma_i \mathbf{u}_{i+1} \right) / s_t$.

Theorem (sublinear convergence). For any $(\mathbf{x}, \mathbf{u}) \in \text{dom} \mathcal{L}$:

$$\mathcal{L}(\bar{\mathbf{x}}_{t+1}, \mathbf{u}) - \mathcal{L}(\mathbf{x}, \bar{\mathbf{u}}_{t+1}) \leq \frac{\|\mathbf{z}_0 - \mathbf{x}\|^2 + \gamma_0^2 \|\mathbf{u}_0 - \mathbf{u}\|^2}{2s_t}.$$

Corollary (sublinear on objective). If h is β_h -Lipschitz,

$$(f+g+h)(\bar{\mathbf{x}}_{t+1}) - (f+g+h)(\mathbf{x}^*) \leq \frac{\|\mathbf{z}_0 - \mathbf{x}^*\|^2 + 2\gamma_0^2 (\|\mathbf{u}_0\|^2 + \beta_h^2)}{2s_t} = \mathcal{O}(1/t)$$

Convergence Analysis

Theorem (linear convergence). If f is μ_f -strongly cvx and h is L_h -smooth,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \min \left\{ \tau \frac{\mu_f}{L_f}, \frac{1}{1 + \gamma_0 L_h} \right\} \right)^{t+1} C_0,$$

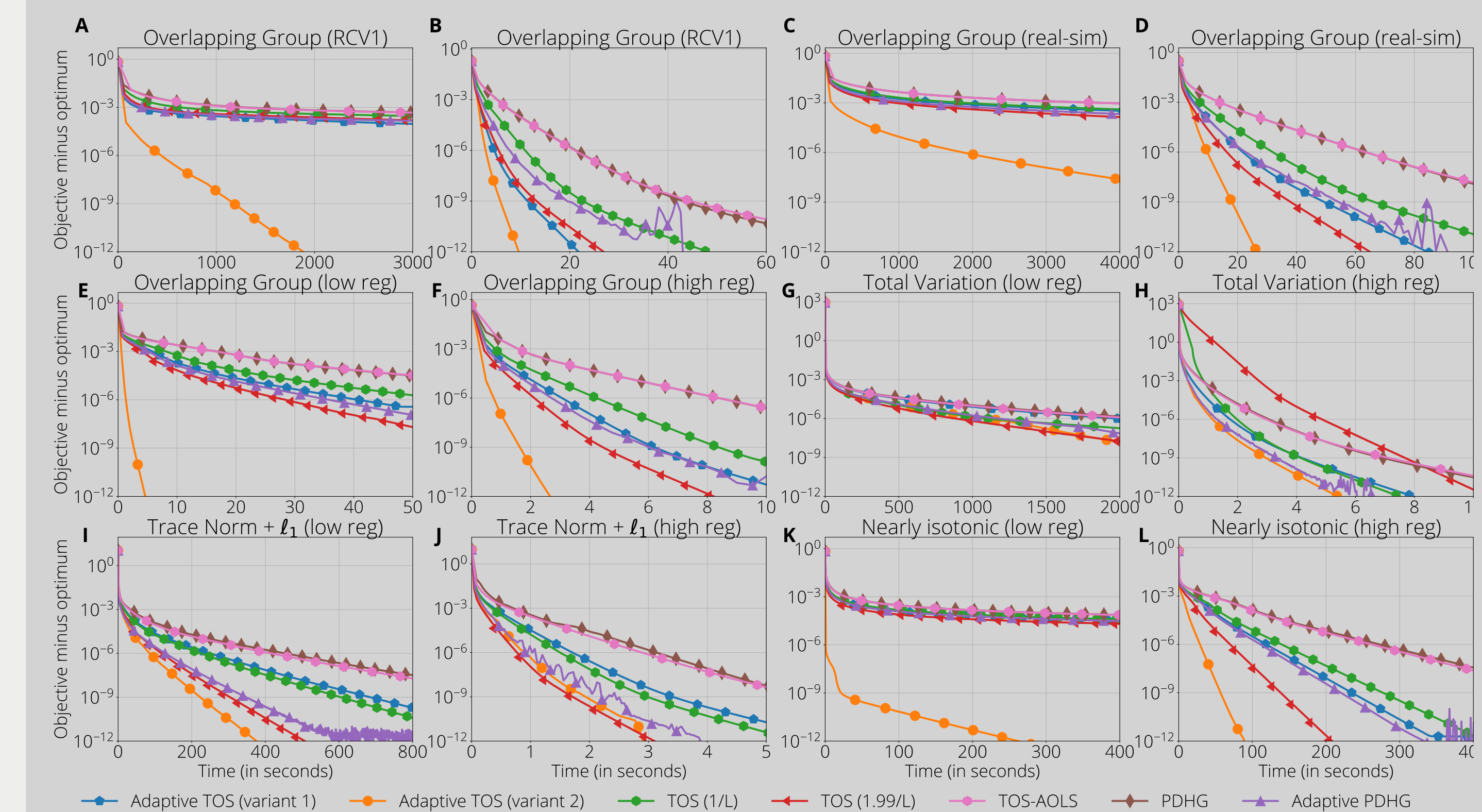
with $C_0 \stackrel{\text{def}}{=} 6\|\mathbf{z}_0 - \mathbf{x}^*\|^2 + \frac{6}{1-\sigma} \|\gamma_0(\mathbf{u}_0 - \mathbf{u}^*)\|^2$.

- Rate factor $\min \left\{ \frac{\mu_f}{L_f}, \frac{1}{1 + \gamma_0 L_h} \right\}$ larger than previous $\frac{\mu_f}{L_f} \times \frac{1}{(1 + \gamma_0 L_h)^2}$ (Davis and Yin 2015).

Experiments

Comparison with related splitting methods (Condat 2013, PDHG), (Giselsson et al. 2016, TOS-AOLS), (Malitsky and Pock 2018, Adaptive-PDHG).

Considered problems: logistic regression + nearly-isotonic, overlapping group lasso and least squares + total variation, large and small regularization regime.



Large computational gains with

- non-quadratic los functions
- low regularization regime.

References

- Condat, Laurent (2013). "A primal-dual splitting method for convex optimization involving lipschitzian, proximal and linear composite terms". In: *Journal of Optimization Theory and Applications*.
- Davis, Damek and Wotao Yin (2015). "A three-operator splitting scheme and its optimization applications". In: *arXiv preprint arXiv:1504.01032*.
- (2017). "A three-operator splitting scheme and its optimization applications". In: *Set-Valued and Variational Analysis*.
- Giselsson, Pontus et al. (2016). "Line search for averaged operator iteration". In: *Conference on Decision and Control (CDC)*.
- Malitsky, Yura and Thomas Pock (2018). "A first-order primal-dual algorithm with linesearch". In: *SIAM Journal on Optimization*.