

Sufficient decrease is all you need

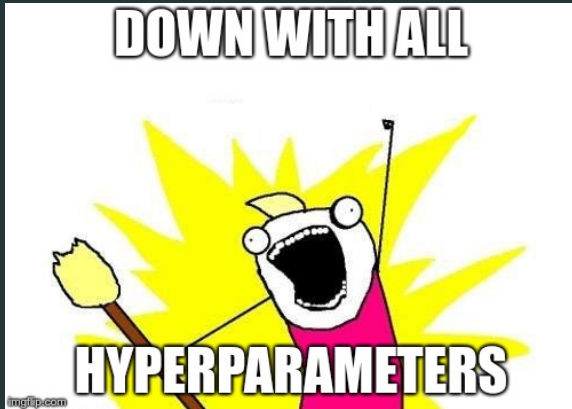
Fabian Pedregosa

DALI 2019



Berkeley

Motivation



Problem Setting

Most first-order optimization methods are of the form

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{p}_t \quad \text{where}$$

- \mathbf{p}_t is the update direction, determined by the algorithm.
- γ_t is a the step-size, free parameter.

How to choose the step-size?

- Chosen **in advance**, i.e., $\gamma_t = 1/L$ or $\gamma_t = L/\sqrt{t+1}$
+ Simple, - Suboptimal

Problem Setting

Most first-order optimization methods are of the form

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{p}_t \quad \text{where}$$

- \mathbf{p}_t is the update direction, determined by the algorithm.
- γ_t is a the step-size, free parameter.

How to choose the step-size?

- Chosen **in advance**, i.e., $\gamma_t = 1/L$ or $\gamma_t = L/\sqrt{t+1}$
+ Simple, - Suboptimal
- Exact line-search: $\gamma_t \in \mathbf{arg\,min}_{\gamma \geq 0} f(\mathbf{x}_t + \gamma \mathbf{p}_t)$
+ Optimal, - Expensive

Problem Setting

Most first-order optimization methods are of the form

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{p}_t \quad \text{where}$$

- \mathbf{p}_t is the update direction, determined by the algorithm.
- γ_t is a the step-size, free parameter.

How to choose the step-size?

- Chosen **in advance**, i.e., $\gamma_t = 1/L$ or $\gamma_t = L/\sqrt{t+1}$
+ Simple, - Suboptimal
- Exact line-search: $\gamma_t \in \mathbf{arg\,min}_{\gamma \geq 0} f(\mathbf{x}_t + \gamma \mathbf{p}_t)$
+ Optimal, - Expensive
- **Adaptive** step-size (aka backtracking, inexact LS): Choose γ_t based on local sufficient decrease condition
+Efficient, +Simple

Adaptive Step-size Selection

For gradient descent ($\mathbf{p}_t = -\nabla f(\mathbf{x}_t)$), Armijo backtracking (Armijo, 1966) simple and efficient condition:

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

What about other methods?

\implies **This talk**

Outline

1. Adaptive **Three Operator Splitting** and structured saddle-point problems.
2. Adaptive **Frank-Wolfe** (and linearly-convergent variants).
3. **Perspectives**: stochastic optimization.

Three Operator Splitting (TOS) (Davis and Yin, 2017)

Solves optimization problems of the form

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

where f is convex and L -smooth, g, h are convex with access to $\text{prox}_{\gamma g}(x) \stackrel{\text{def}}{=} \arg \min_z g(z) + \frac{1}{2\gamma} \|x - z\|^2$, $\text{prox}_{\gamma h}$.

Iterates on $y_t \in \mathbb{R}^d$ given by

$$z_t = \text{prox}_{\gamma h}(y_t)$$

$$x_t = \text{prox}_{\gamma g}(2y_t - z_t - \gamma \nabla f(z_t))$$

$$y_{t+1} = y_t - z_t + x_t$$

Three Operator Splitting (TOS) (Davis and Yin, 2017)

Solves optimization problems of the form

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

where f is convex and L -smooth, g, h are convex with access to $\text{prox}_{\gamma g}(x) \stackrel{\text{def}}{=} \arg \min_z g(z) + \frac{1}{2\gamma} \|x - z\|^2$, $\text{prox}_{\gamma h}$.

Iterates on $y_t \in \mathbb{R}^d$ given by

$$z_t = \text{prox}_{\gamma h}(y_t)$$

$$x_t = \text{prox}_{\gamma g}(2y_t - z_t - \gamma \nabla f(z_t))$$

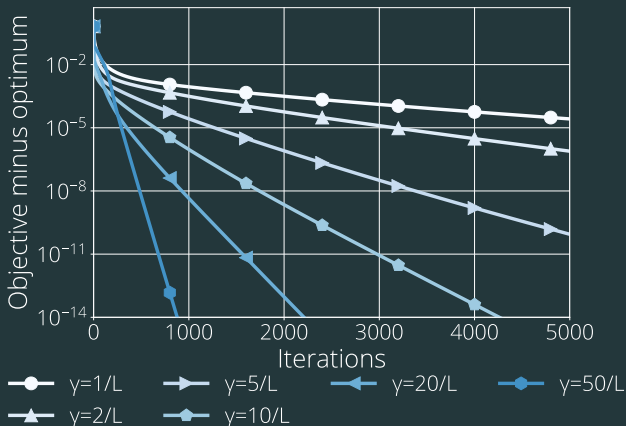
$$y_{t+1} = y_t - z_t + x_t$$

Generalizes proximal-gradient (FB, ISTA) and Douglas-Rachford.

Importance of step-size

Guaranteed convergence $\gamma < 2/L$, with $L = \text{Lipschitz const. of } \nabla f$.

In practice, best performance is when $\gamma \gg 2/L$



What can we do about it?

Revisiting the Three Operator Splitting

Saddle-point reformulation of original problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + g(\mathbf{x}) + h(\mathbf{x}) \\ = \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + g(\mathbf{x}) + \max_{\mathbf{u} \in \mathbb{R}^d} \{\langle \mathbf{x}, \mathbf{u} \rangle - h^*(\mathbf{u})\} \end{aligned}$$

Revisiting the Three Operator Splitting

Saddle-point reformulation of original problem

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

We can rewrite the three operator splitting as

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

$$\mathbf{u}_{t+1} = \mathbf{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma),$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$

Minimizing with respect to primal variable

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{u}_t) = \min_{\mathbf{x} \in \mathbb{R}^d} \underbrace{f(\mathbf{x}) + \langle \mathbf{x}, \mathbf{u}_t \rangle}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{proximal}}$$

Minimizing with respect to primal variable

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{u}_t) = \min_{\mathbf{x} \in \mathbb{R}^d} \underbrace{f(\mathbf{x}) + \langle \mathbf{x}, \mathbf{u}_t \rangle}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{proximal}}$$

- Proximal-gradient iteration, with $\mathbf{x} = \mathbf{z}_t$ as starting point:

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

Minimizing with respect to primal variable

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{u}_t) = \min_{\mathbf{x} \in \mathbb{R}^d} \underbrace{f(\mathbf{x}) + \langle \mathbf{x}, \mathbf{u}_t \rangle}_{\text{smooth}} + \underbrace{g(\mathbf{x})}_{\text{proximal}}$$

- Proximal-gradient iteration, with $\mathbf{x} = \mathbf{z}_t$ as starting point:

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

= first step of TOS

Minimizing with respect to the dual variable

$$\min_{\mathbf{u} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}_t, \mathbf{u}) = \min_{\mathbf{u} \in \mathbb{R}^d} \underbrace{h^*(\mathbf{u}) - \langle \mathbf{x}_t, \mathbf{u} \rangle}_{\text{proximal}}$$

Minimizing with respect to the dual variable

$$\min_{\mathbf{u} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}_t, \mathbf{u}) = \min_{\mathbf{u} \in \mathbb{R}^d} \underbrace{h^*(\mathbf{u}) - \langle \mathbf{x}_t, \mathbf{u} \rangle}_{\text{proximal}}$$

- Proximal-point iteration:

$$\mathbf{u}_{t+1} = \mathbf{prox}_{\sigma h^*}(\mathbf{u}_t + \sigma \mathbf{x}_{t+1})$$

= second update in TOS with $\sigma = 1/\gamma$

Extrapolation Step

Last update:

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t) \quad (\text{extrapolation step})$$

Extrapolation Step

Last update:

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t) \quad (\text{extrapolation step})$$

Verifies $\mathbf{z}_{t+1} \in \partial h^*(\mathbf{u}_{t+1})$.

Extrapolation Step

Last update:

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t) \quad (\text{extrapolation step})$$

Verifies $\mathbf{z}_{t+1} \in \partial h^*(\mathbf{u}_{t+1})$.

At optimum, we have $\mathbf{x}^* \in \partial h^*(\mathbf{u}^*)$.

\implies Solving the KKT conditions at \mathbf{u}_{t+1} .

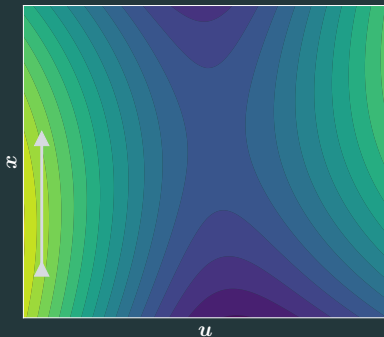
Revisiting the three operator splitting

Iteration 1: proximal-gradient step

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



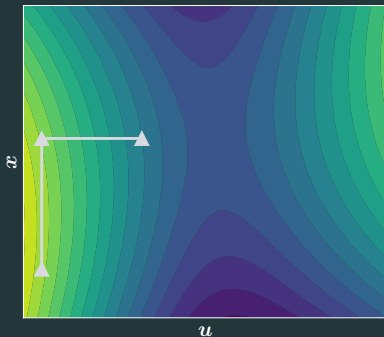
Revisiting the three operator splitting

Iteration 1: proximal-point step

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

$$\rightarrow \mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



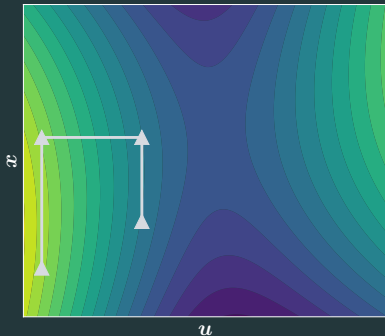
Revisiting the three operator splitting

Iteration 2: extrapolation

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\rightarrow \mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



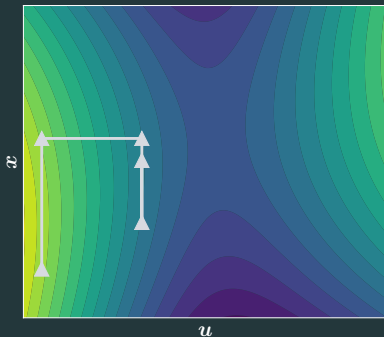
Revisiting the three operator splitting

Iteration 2: proximal-gradient step

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



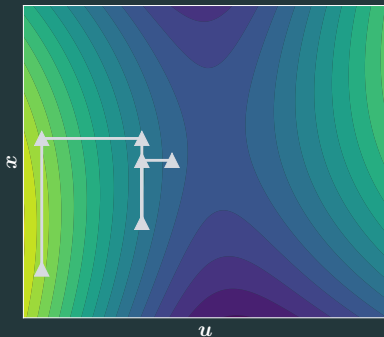
Revisiting the three operator splitting

Iteration 3: proximal-point step

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

$$\rightarrow \mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



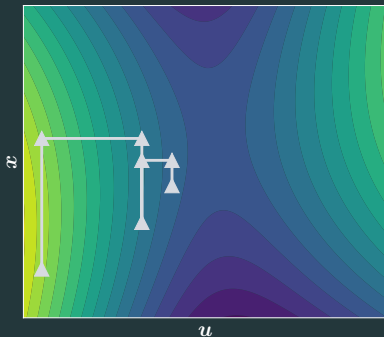
Revisiting the three operator splitting

Iteration 3: extrapolation

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\rightarrow \mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



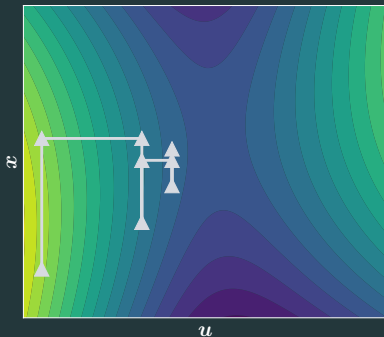
Revisiting the three operator splitting

Iteration 4: proximal-gradient step

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



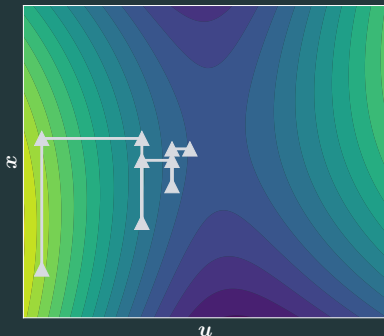
Revisiting the three operator splitting

Iteration 4: proximal-point step

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

$$\rightarrow \mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



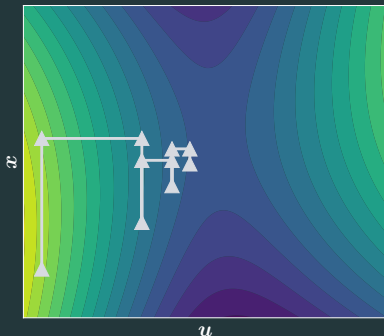
Revisiting the three operator splitting

Iteration 5: extrapolation

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\rightarrow \mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



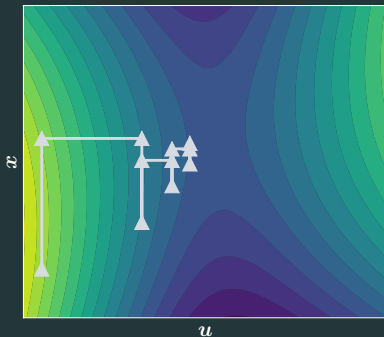
Revisiting the three operator splitting

Iteration 5: proximal-gradient step

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



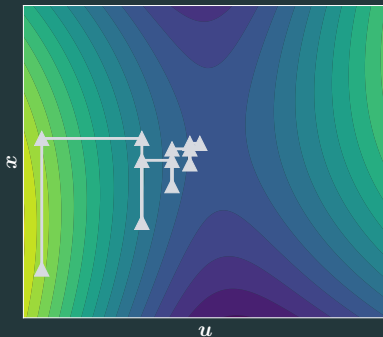
Revisiting the three operator splitting

Iteration 6: proximal-point step

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

$$\rightarrow \mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



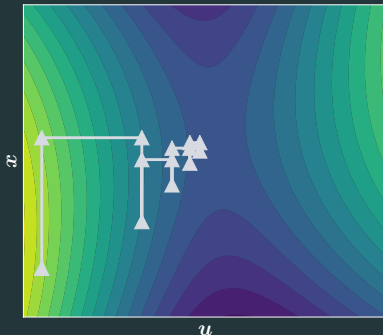
Revisiting the three operator splitting

Iteration 6: extrapolation

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\rightarrow \mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



Take-Home Message

TOS is (basically) alternated proximal-gradient and proximal-point

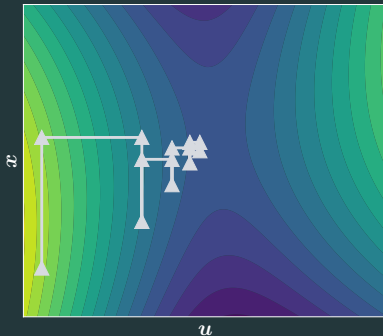
Revisiting the three operator splitting

Iteration 6: extrapolation

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

$$\mathbf{u}_{t+1} = \text{prox}_{h^*/\gamma}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma)$$

$$\rightarrow \mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma(\mathbf{u}_{t+1} - \mathbf{u}_t)$$



Take-Home Message

TOS is (basically) alternated proximal-gradient and proximal-point

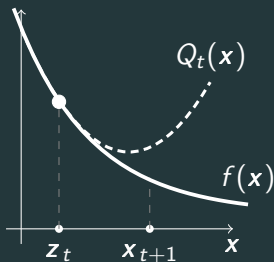
Can we use the adaptive step-size of proximal-gradient?

Adaptive TOS (Pedregosa and Gidel, 2018)

Let $Q_t(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{z}_t) + \langle \nabla f(\mathbf{z}_t), \mathbf{x} - \mathbf{z}_t \rangle + \frac{1}{2\gamma_t} \|\mathbf{x} - \mathbf{z}_t\|^2$.

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

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

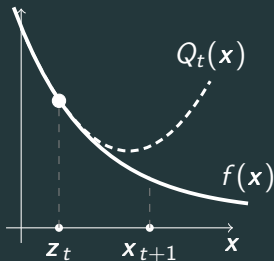


Adaptive TOS (Pedregosa and Gidel, 2018)

Let $Q_t(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{z}_t) + \langle \nabla f(\mathbf{z}_t), \mathbf{x} - \mathbf{z}_t \rangle + \frac{1}{2\gamma_t} \|\mathbf{x} - \mathbf{z}_t\|^2$.

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

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

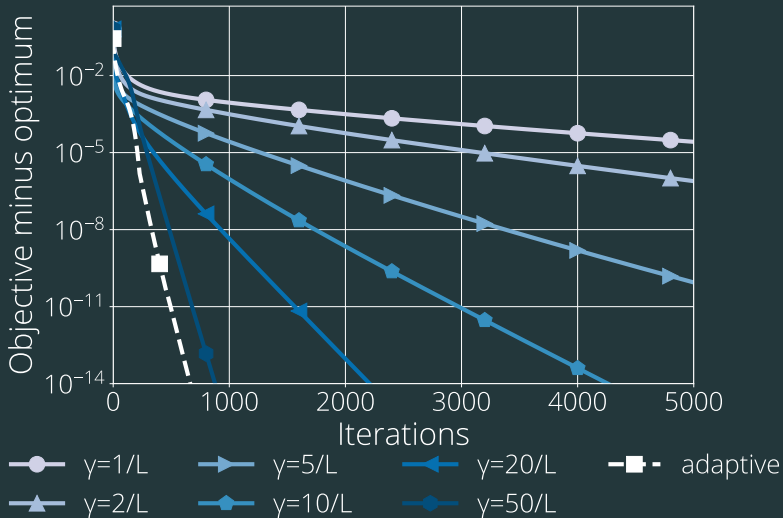


Run rest of algorithm with that step-size:

$$\mathbf{u}_{t+1} = \mathbf{prox}_{h^*/\gamma_t}(\mathbf{u}_t + \mathbf{x}_{t+1}/\gamma_t) \quad (1)$$

$$\mathbf{z}_{t+1} = \mathbf{x}_{t+1} - \gamma_t(\mathbf{u}_{t+1} - \mathbf{u}_t) \quad (2)$$

Performance of the adaptive step-size strategy



Performance is as good as best hand-tuned step-size

Convergence rates (informal)

As good as the original method with fixed step-size

Theorem (sublinear convergence rate)

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

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

Theorem

If f is L_f -smooth, μ -strongly convex and h is L_h -smooth then

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \left(1 - \min \left\{ \tau \frac{\mu}{L_f}, \frac{1}{1 + \gamma_0 L_h} \right\}\right)^{t+1} C_0 \quad (3)$$

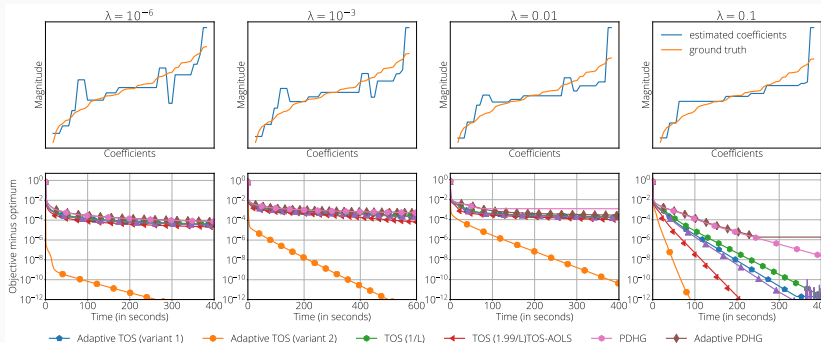
with $\tau =$ line search decrease factor, $C_0 =$ only depends on initial conditions.

Experiments

Logistic + Nearly-isotonic penalty

Problem

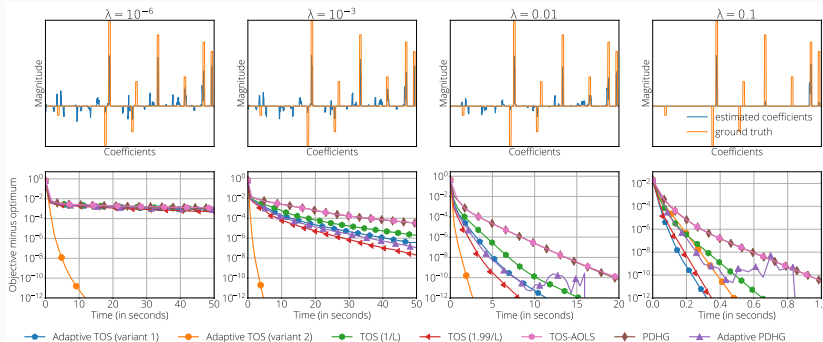
$$\arg \min_{\mathbf{x}} \text{logistic}(\mathbf{x}) + \lambda \sum_{i=1}^{p-1} \max\{\mathbf{x}_i - \mathbf{x}_{i+1}, 0\}$$



Logistic + Overlapping group lasso penalty

Problem

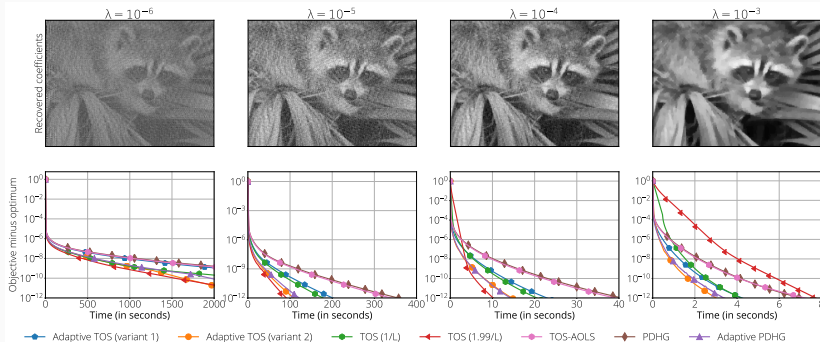
$$\arg \min_{\mathbf{x}} \text{logistic}(\mathbf{x}) + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{x}\|_2$$



Quadratic loss + total variation penalty

Problem

$$\arg \min_x \text{least_squares}(x) + \lambda \|x\|_{\text{TV}}$$



Adaptive Frank-Wolfe

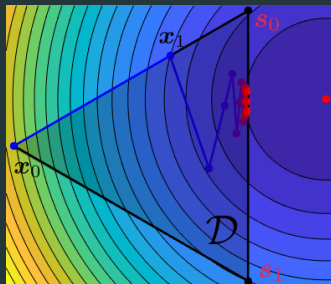
The Frank-Wolfe (FW) algorithm, aka conditional gradient

Problem: smooth f , compact \mathcal{D}

$$\arg \min_{x \in \mathcal{D}} f(x)$$

Algorithm 1: Frank-Wolfe (FW)

- 1 for $t = 0, 1 \dots$ do
 - 2 $\mathbf{s}_t \in \arg \min_{s \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$
 - 3 $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$
 - 4 Find γ_t , e.g., by line-search:
 $\gamma_t \in \arg \min_{\gamma \in [0,1]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$
 - 5 $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma \mathbf{d}_t$
-



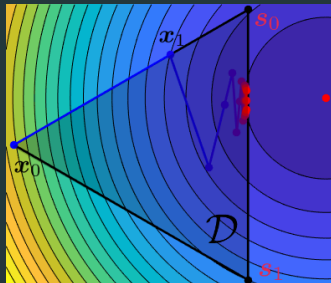
The Frank-Wolfe (FW) algorithm, aka conditional gradient

Problem: smooth f , compact \mathcal{D}

$$\arg \min_{x \in \mathcal{D}} f(x)$$

Algorithm 1: Frank-Wolfe (FW)

- 1 for $t = 0, 1 \dots$ do
 - 2 $\mathbf{s}_t \in \arg \min_{s \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$
 - 3 $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$
 - 4 Find γ_t , e.g., by line-search:
 $\gamma_t \in \arg \min_{\gamma \in [0,1]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$
 - 5 $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma \mathbf{d}_t$
-



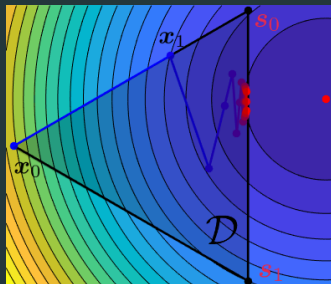
The Frank-Wolfe (FW) algorithm, aka conditional gradient

Problem: smooth f , compact \mathcal{D}

$$\arg \min_{x \in \mathcal{D}} f(x)$$

Algorithm 1: Frank-Wolfe (FW)

- 1 for $t = 0, 1 \dots$ do
 - 2 $\mathbf{s}_t \in \arg \min_{s \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$
 - 3 $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$
 - 4 Find γ_t , e.g., by line-search:
 $\gamma_t \in \arg \min_{\gamma \in [0,1]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$
 - 5 $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma \mathbf{d}_t$
-



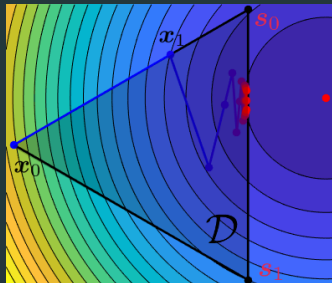
The Frank-Wolfe (FW) algorithm, aka conditional gradient

Problem: smooth f , compact \mathcal{D}

$$\arg \min_{x \in \mathcal{D}} f(x)$$

Algorithm 1: Frank-Wolfe (FW)

```
1 for  $t = 0, 1 \dots$  do
2    $\mathbf{s}_t \in \arg \min_{s \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$ 
3    $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$ 
4   Find  $\gamma_t$ , e.g., by line-search:
    $\gamma_t \in \arg \min_{\gamma \in [0,1]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$ 
5    $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma \mathbf{d}_t$ 
```



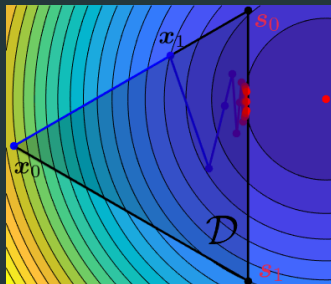
The Frank-Wolfe (FW) algorithm, aka conditional gradient

Problem: smooth f , compact \mathcal{D}

$$\arg \min_{x \in \mathcal{D}} f(x)$$

Algorithm 1: Frank-Wolfe (FW)

- 1 for $t = 0, 1 \dots$ do
 - 2 $\mathbf{s}_t \in \arg \min_{s \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$
 - 3 $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$
 - 4 Find γ_t , e.g., by line-search:
 $\gamma_t \in \arg \min_{\gamma \in [0,1]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$
 - 5 $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma \mathbf{d}_t$
-



Setting the step-size in Frank-Wolfe

- Exact line-search only feasible for quadratic objective.

Setting the step-size in Frank-Wolfe

- Exact line-search only feasible for quadratic objective.
- “Oblivious” step-size $\gamma_t = 2/(t + 2)$ is convergent, but slow.

Setting the step-size in Frank-Wolfe

- Exact line-search only feasible for quadratic objective.
- “Oblivious” step-size $\gamma_t = 2/(t + 2)$ is convergent, but slow.
- New linearly-convergent variants (Lacoste-Julien and Jaggi, 2015) assume access to exact line search.

Setting the step-size in Frank-Wolfe

- Exact line-search only feasible for quadratic objective.
- “Oblivious” step-size $\gamma_t = 2/(t + 2)$ is convergent, but slow.
- New linearly-convergent variants (Lacoste-Julien and Jaggi, 2015) assume access to exact line search.

We would like:

- Efficient.
- Adaptive to local geometry.
- That achieves best possible rates in every situation.

Setting the step-size in Frank-Wolfe

- Exact line-search only feasible for quadratic objective.
- “Oblivious” step-size $\gamma_t = 2/(t + 2)$ is convergent, but slow.
- New linearly-convergent variants (Lacoste-Julien and Jaggi, 2015) assume access to exact line search.

We would like:

- Efficient.
- Adaptive to local geometry.
- That achieves best possible rates in every situation.

Is it possible?

A practical issue

- In a polytope, FW moves in the direction of a vertex.

A practical issue

- In a polytope, FW moves in the direction of a vertex.
- Two vertices can be far apart \implies optimal step-size does not vary smoothly.

A practical issue

- In a polytope, FW moves in the direction of a vertex.
- Two vertices can be far apart \implies optimal step-size does not vary smoothly.
- Directly using Armijo/Sufficient decrease conditions (Dunn, 1980) is particularly difficult because of this.

The Adaptive FW algorithm (Pedregosa et al., 2018)

Key Idea

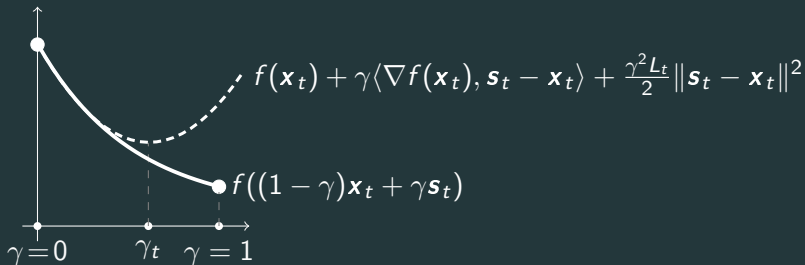
Estimate smoothness, not step-size!

Algorithm 2: The Adaptive Frank-Wolfe algorithm (AdaFW)

```
1 for  $t = 0, 1 \dots$  do
2    $\mathbf{s}_t \in \arg \min_{\mathbf{s} \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$ 
3    $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$ 
4   Find  $L_t$  that verifies sufficient decrease (4), with
5    $\gamma_t = \min \left\{ \frac{\langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle}{L_t \|\mathbf{d}_t\|^2}, 1 \right\}$ 
6    $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$ 
```

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \gamma_t \langle \nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle + \frac{\gamma_t^2 L_t}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2 \quad (4)$$

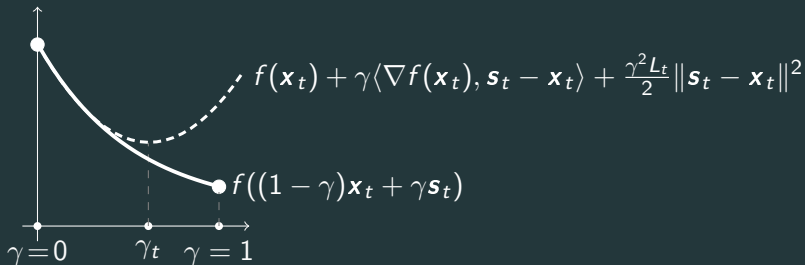
The Adaptive FW algorithm (Pedregosa et al., 2018)¹



- Worst-case, $L_t = L$. Often $L_t \ll L \implies$ larger step-size.

¹Fabian Pedregosa, Armin Askari, Geoffrey Negiar, and Martin Jaggi (2018). “Step-Size Adaptivity in Projection-Free Optimization”. In: *ArXiv*.

The Adaptive FW algorithm (Pedregosa et al., 2018)¹



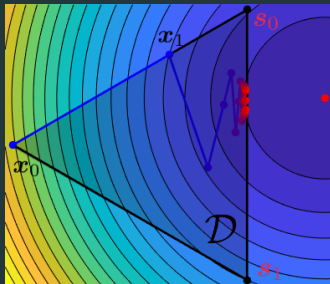
- Worst-case, $L_t = L$. Often $L_t \ll L \implies$ larger step-size.
- Two extra function evaluations per iteration. Often given as byproduct of gradient.

¹Fabian Pedregosa, Armin Askari, Geoffrey Negiar, and Martin Jaggi (2018). “Step-Size Adaptivity in Projection-Free Optimization”. In: *ArXiv*.

Extension to other FW variants

Zig-Zagging phenomena in FW

The Frank-Wolfe algorithm zig-zags when the solution lies in a face of the boundary.



Some FW variants have been developed to address this issue.

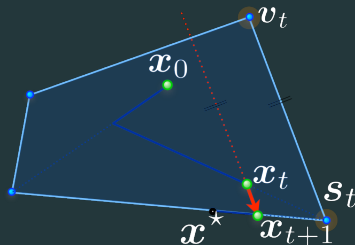
Pairwise FW

Key idea

1. Keep track of previously added vertices (activeset) \mathcal{S}_t .
2. Move weight mass between two vertices in each step.

Algorithm 3: Pairwise FW (Lacoste-Julien and Jaggi, 2015)

```
1 for  $t = 0, 1 \dots$  do
2    $\mathbf{s}_t \in \arg \min_{\mathbf{s} \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$ 
3    $\mathbf{v}_t \in \arg \max_{\mathbf{s} \in \mathcal{S}_t} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$ 
4    $\mathbf{d}_t = \mathbf{s}_t - \mathbf{v}_t$ 
5   Find  $\gamma_t$  by line-search:  $\gamma_t \in$ 
      $\arg \min_{\gamma \in [0, \gamma_t^{\max}]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$ 
6    $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$ 
```



Adaptive Away-steps and Pairwise FW (Pedregosa et al., 2018)

Convergence of Away-steps and Pairwise FW

- Linear convergence for strongly convex functions on polytopes (Lacoste-Julien and Jaggi, 2015).

Adaptive Away-steps and Pairwise FW (Pedregosa et al., 2018)

Convergence of Away-steps and Pairwise FW

- Linear convergence for strongly convex functions on polytopes (Lacoste-Julien and Jaggi, 2015).
- Can we design variants with adaptive step-size?

Adaptive Away-steps and Pairwise FW (Pedregosa et al., 2018)

Convergence of Away-steps and Pairwise FW

- Linear convergence for strongly convex functions on polytopes (Lacoste-Julien and Jaggi, 2015).
- Can we design variants with adaptive step-size?

Introducing Adaptive Away-steps and Adaptive Pairwise

Choose L_t such that it verifies

$$f(\mathbf{x}_t + \gamma_t \mathbf{d}_t) \leq f(\mathbf{x}_t) + \gamma_t \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{\gamma_t^2 L_t}{2} \|\mathbf{d}_t\|^2$$

$$\text{with } \gamma_t = \min \left\{ \frac{\langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle}{L_t \|\mathbf{d}_t\|^2}, \gamma_{\max} \right\}$$

Theory for Adaptive Step-size variants

Strongly convex f

Pairwise and Away-steps converge linearly on a polytope. For each “good step” we have:

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq (1 - \rho)(f(\mathbf{x}_t) - f(\mathbf{x}^*)) , \quad \rho > 0$$

Theory for Adaptive Step-size variants

Strongly convex f

Pairwise and Away-steps converge linearly on a polytope. For each “good step” we have:

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq (1 - \rho)(f(\mathbf{x}_t) - f(\mathbf{x}^*)) , \quad \rho > 0$$

Convex f

For all FW variants, $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathcal{O}(1/t)$

Theory for Adaptive Step-size variants

Strongly convex f

Pairwise and Away-steps converge linearly on a polytope. For each “good step” we have:

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq (1 - \rho)(f(\mathbf{x}_t) - f(\mathbf{x}^*)) , \quad \rho > 0$$

Convex f

For all FW variants, $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathcal{O}(1/t)$

Non-Convex f

For all FW variants, $\max_{\mathbf{s} \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s} \rangle \leq \mathcal{O}(1/\sqrt{t})$

Same rate as with exact line search

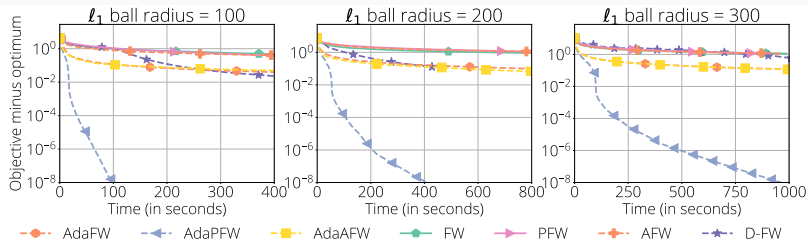
Experiments

Experiments RCV1

Problem: l_1 -constrained logistic regression

$$\arg \min_{\|x\|_1 \leq \alpha} \frac{1}{n} \sum_{i=1}^n \varphi(\mathbf{a}_i^T x, b_i) \quad \text{with } \varphi = \text{logistic loss.}$$

Dataset	dimension	density	\bar{L}_t/L
RCV1	47236	10^{-3}	1.3×10^{-2}

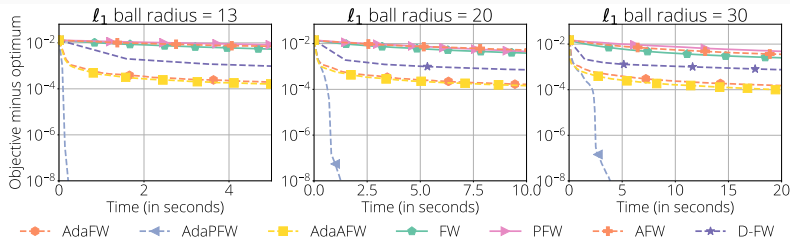


Experiments Madelon

Problem: ℓ_1 -constrained logistic regression

$$\arg \min_{\|x\|_1 \leq \alpha} \frac{1}{n} \sum_{i=1}^n \varphi(\mathbf{a}_i^T x, b_i) \quad \text{with } \varphi = \text{logistic loss.}$$

Dataset	dimension	density	\bar{L}_t/L
Madelon	500	1.	3.3×10^{-3}

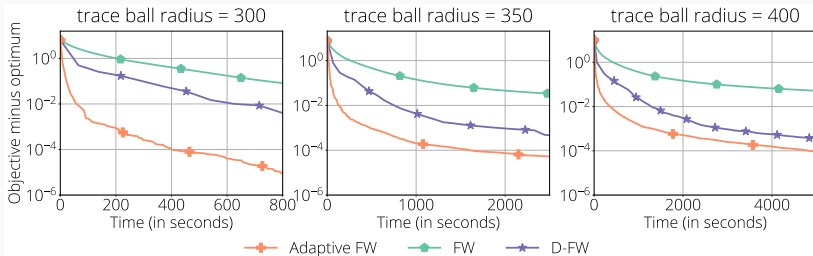


Experiments MovieLens 1M

Problem: trace-norm constrained robust matrix completion

$$\arg \min_{\|x\|_* \leq \alpha} \frac{1}{|B|} \sum_{(i,j) \in B}^n h(\mathbf{X}_{i,j}, \mathbf{A}_{i,j}) \quad \text{with } h = \text{Huber loss.}$$

Dataset	dimension	density	\bar{L}_t/L
MovieLens 1M	22,393,987	0.04	1.1×10^{-2}



Perspectives

Stochastic optimization

Problem

$$\arg \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

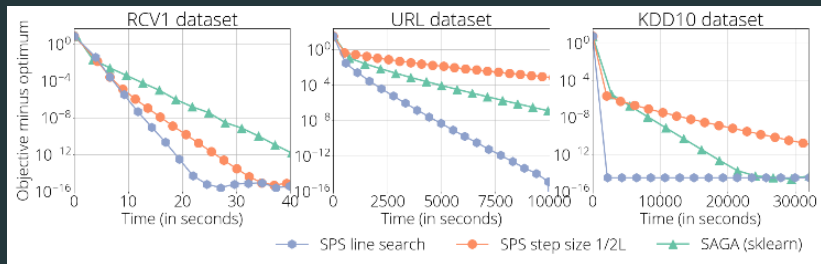
Main challenge: How to evaluate Armijo condition $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma \|\mathbf{p}_t\|^2$ without access to f ?

Experiments stochastic line search

Heuristic from (Schmidt et al. 2017)² to estimate L :

$$f_i(\mathbf{x}_t - \frac{1}{L} \nabla f_i(\mathbf{x}_t)) \leq f_i(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f_i(\mathbf{x}_t)\|^2$$

with i random index sampled at iter t .



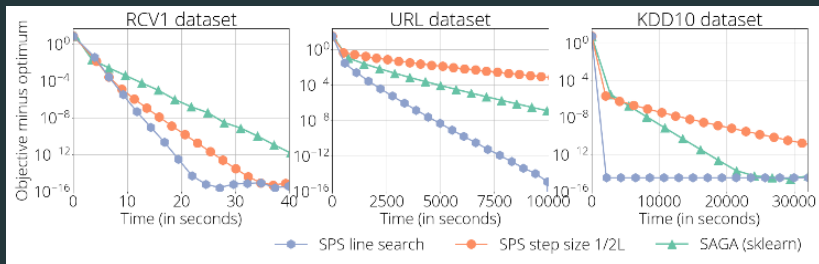
²Mark Schmidt, Nicolas Le Roux, and Francis Bach (2017). "Minimizing finite sums with the stochastic average gradient". In: *Mathematical Programming*.

Experiments stochastic line search

Heuristic from (Schmidt et al. 2017)² to estimate L :

$$f_i(\mathbf{x}_t - \frac{1}{L} \nabla f_i(\mathbf{x}_t)) \leq f_i(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f_i(\mathbf{x}_t)\|^2$$

with i random index sampled at iter t .



Can we prove convergence of this (or similar) method?

²Mark Schmidt, Nicolas Le Roux, and Francis Bach (2017). "Minimizing finite sums with the stochastic average gradient". In: *Mathematical Programming*.

Recent developments

- (Shang et al., 2018) Adaptive step-size for SVRG. Condition is costly to evaluate.

Recent developments

- (Shang et al., 2018) Adaptive step-size for SVRG. Condition is costly to evaluate.
- (Paquette and Scheinberg, 2018) Evaluates a stochastic version of the Armijo condition. Accepts step-size based on concept of reliable/unreliable estimate.

Recent developments

- (Shang et al., 2018) Adaptive step-size for SVRG. Condition is costly to evaluate.
- (Paquette and Scheinberg, 2018) Evaluates a stochastic version of the Armijo condition. Accepts step-size based on concept of reliable/unreliable estimate.
(step-size tends to be really small)

Conclusion

- Applicability of sufficient decrease beyond classical framework.

Conclusion

- Applicability of sufficient decrease beyond classical framework.
- Sufficient decrease condition to set step-size in TOS and FW and variants.

Conclusion

- Applicability of sufficient decrease beyond classical framework.
- Sufficient decrease condition to set step-size in TOS and FW and variants.
 - Faster
 - (Mostly) Hyperparameter-free.

Conclusion

- Applicability of sufficient decrease beyond classical framework.
- Sufficient decrease condition to set step-size in TOS and FW and variants.
 - Faster
 - (Mostly) Hyperparameter-free.
- Perspectives in stochastic optimization.

Thanks for your attention

References



Armijo, Larry (1966). "Minimization of functions having Lipschitz continuous first partial derivatives". In: *Pacific Journal of Mathematics*.



Davis, Damek and Wotao Yin (2017). "A three-operator splitting scheme and its optimization applications". In: *Set-valued and variational analysis*.



Dunn, Joseph C (1980). "Convergence rates for conditional gradient sequences generated by implicit step length rules". In: *SIAM Journal on Control and Optimization*.



Lacoste-Julien, Simon and Martin Jaggi (2015). "On the global linear convergence of Frank-Wolfe optimization variants". In: *Advances in Neural Information Processing Systems*.



Paquette, Courtney and Katya Scheinberg (2018). "A stochastic line search method with convergence rate analysis". In: *arXiv preprint arXiv:1807.07994*.



Pedregosa, Fabian et al. (2018). "Step-Size Adaptivity in Projection-Free Optimization". In: *ArXiv*.



Pedregosa, Fabian and Gauthier Gidel (2018). “Adaptive Three Operator Splitting” .
In: *Proceedings of the 35th International Conference on Machine Learning*.



Schmidt, Mark, Nicolas Le Roux, and Francis Bach (2017). “Minimizing finite sums
with the stochastic average gradient” . In: *Mathematical Programming*.



Shang, Fanhua et al. (2018). “Guaranteed Sufficient Decrease for Stochastic Variance
Reduced Gradient Optimization” . In: *arXiv preprint arXiv:1802.09933*.