

HSE University

23 December

Step-size strategies in the conditional gradient algorithm

Ilya Levin

Pavel Zakharov

Valeriia Shcherbakova

Table of content

- Frank-Wolfe algorithm
- Predefined decreasing sequence
- Demyanov-Rubinov step-size
- Exact line-search
- Backtracking line-search
- Results

Frank-Wolfe algorithm

is a method for constrained optimization that solves problems of the form $\underset{\mathbf{x} \in \mathcal{D}}{\text{minimize}} f(\mathbf{x})$

where f is a smooth function for which we have access to its gradient and \mathcal{D} is a compact set. We also assume to have access to a *linear minimization oracle* over \mathcal{D} , that is, a routine that solves problems of the form

$$\mathbf{s}_t \in \arg \max_{\mathbf{s} \in \mathcal{D}} \langle -\nabla f(\mathbf{x}_t), \mathbf{s} \rangle$$

Frank-Wolfe algorithm pseudo code

Frank-Wolfe algorithm

Input: initial guess \mathbf{x}_0 , tolerance $\delta > 0$

For $t = 0, 1, \dots$ **do** $\mathbf{s}_t \in \arg \max_{\mathbf{s} \in \mathcal{D}} \langle -\nabla f(\mathbf{x}_t), \mathbf{s} \rangle$

 Set $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$ and $g_t = \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle$

 Choose step-size γ_t (to be discussed later)

$\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$.

end For loop

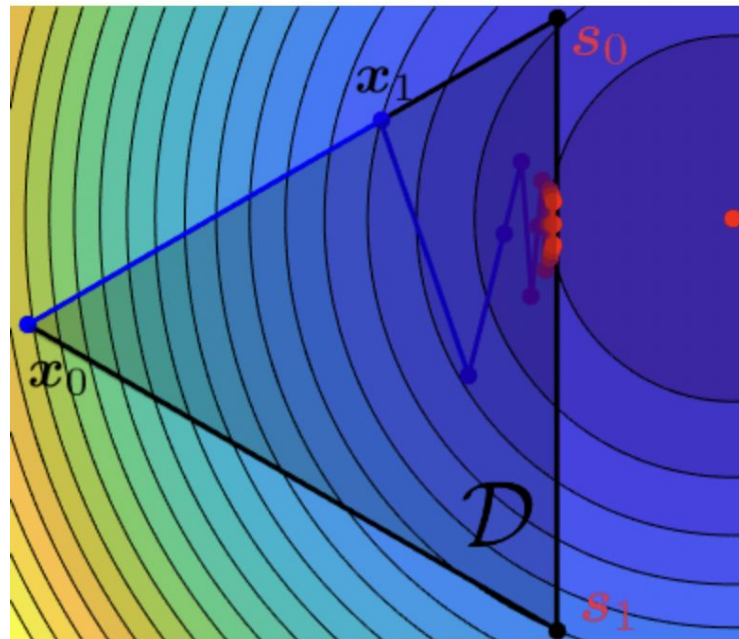
return \mathbf{x}_t

Predefined decreasing sequence

$$\gamma_t = \frac{2}{t+2}$$

is to choose the step-size according to the pre-defined decreasing sequence

- ✚ 1) is straightforward;
2) cheap to compute.
- 1) in practice it performs worse than the alternatives;
2) Oscillates a lot.



Demyanov-Rubinov

$$\gamma = \min \left\{ \frac{g_t}{\mathbf{L} \|\mathbf{d}_t\|^2}, 1 \right\}$$

is to choose the step-size according to the pre-defined decreasing sequence



- 1) Not that expensive to compute;
- 2) Good convergence in practice.



- 1) Can be problematic to estimate \mathbf{L} ;
- 2) Convergence in the neighbourhood of solution is not optimal.

Exact line-search

$$\gamma_{\star} \in \arg \min_{\gamma \in [0,1]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$$

takes the step-size that maximizes the decrease in objective along the update direction



- 1) Great convergence in practice;
- 2) Gives the highest decrease per iteration.



- 1) Pretty heavy computation if we don't have an access to the minimizer of f .

Backtracking line-search

$$\gamma_t = \min \left\{ \frac{g_t}{\mathbf{M}_t \|\mathbf{d}_t\|^2}, 1 \right\}$$

is to choose the step-size according to the pre-defined decreasing sequence

$$Q_t(\gamma, \mathbf{M}_t) \stackrel{\text{def}}{=} f(\mathbf{x}_t) - \gamma g_t + \frac{\gamma^2 \mathbf{M}_t}{2} \|\mathbf{d}_t\|^2$$

- + wildly successful and are a core part of any state-of-the-art implementation of (proximal) gradient descent and Quasi-Newton methods.

$$\mathbf{M}_t = \eta \mathbf{M}_{t-1}$$

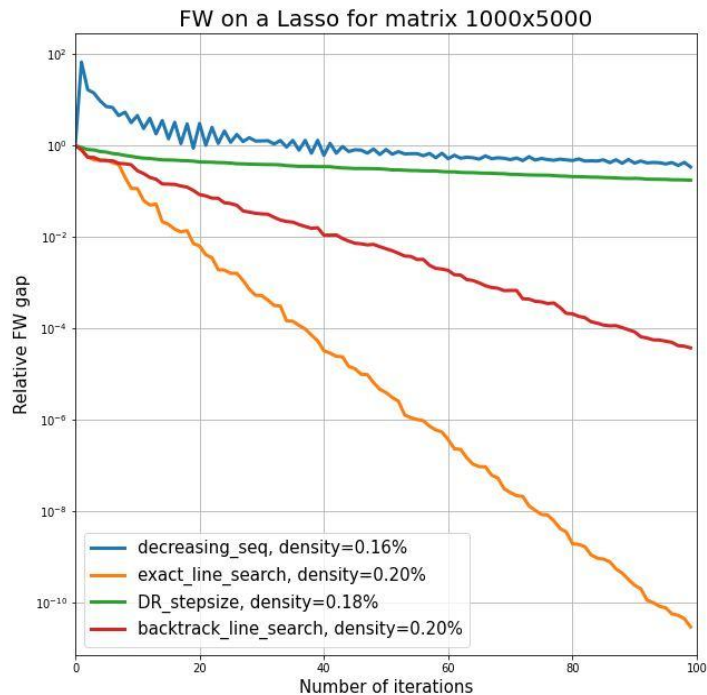
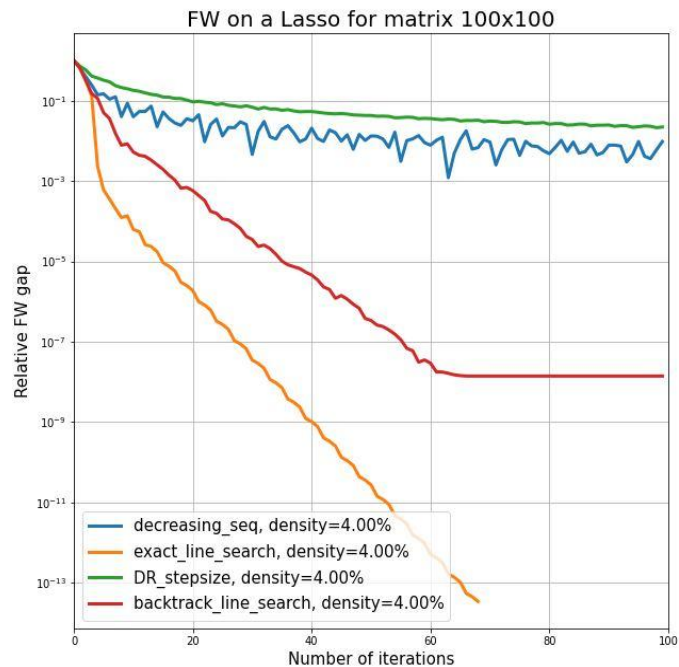
$$\gamma_t = \min \{ g_t / (\mathbf{M}_t \|\mathbf{d}_t\|^2), 1 \}$$

While $f(\mathbf{x}_t + \gamma_t \mathbf{d}_t) > Q_t(\gamma_t, \mathbf{M}_t)$ **do**

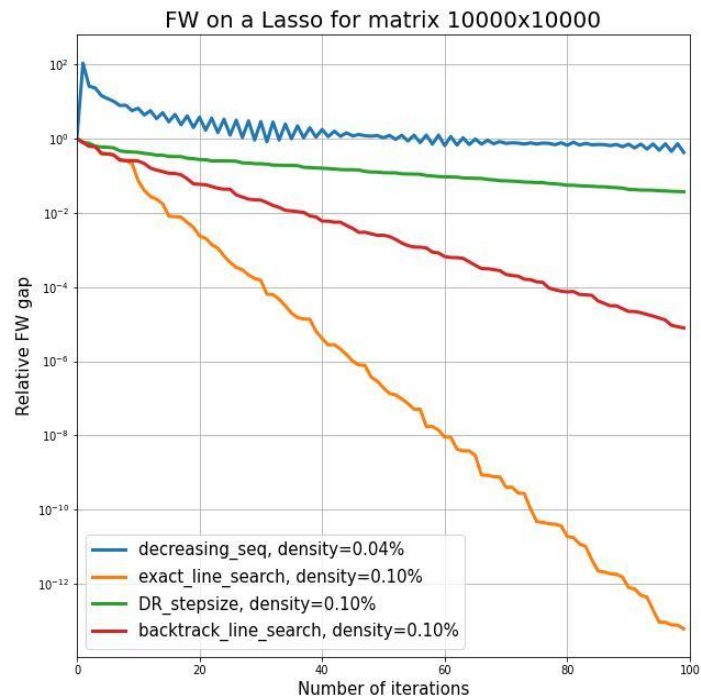
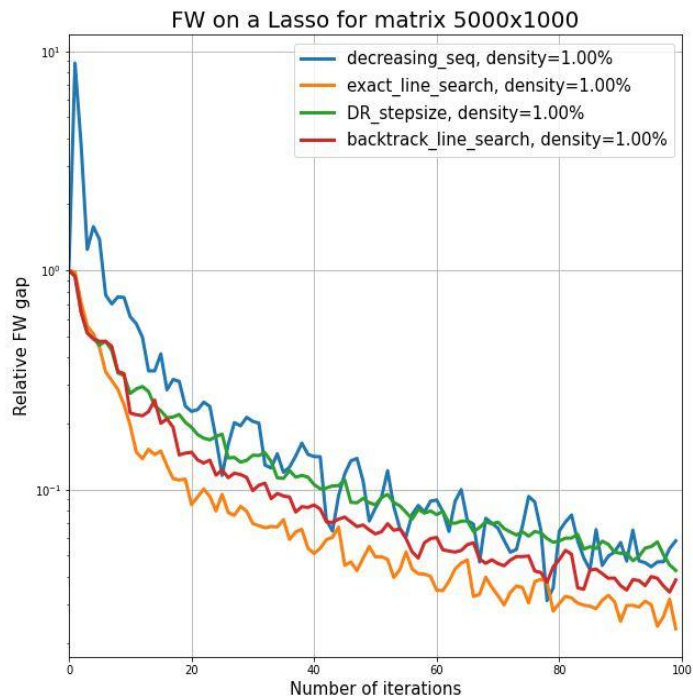
$$\mathbf{M}_t = \tau \mathbf{M}_t$$

- Highly depends on the input parameters

Our results: linear regression

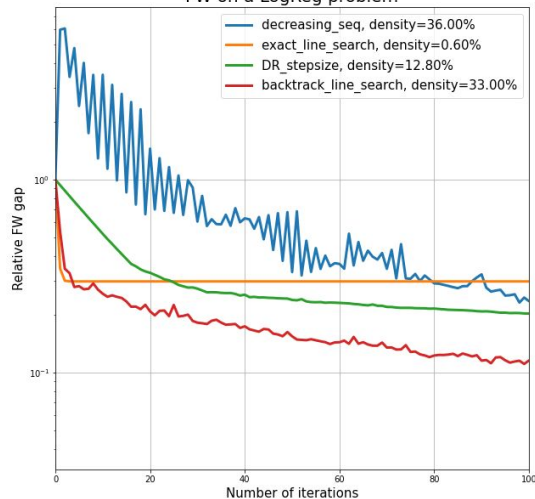


Synthetic dataset MSE:

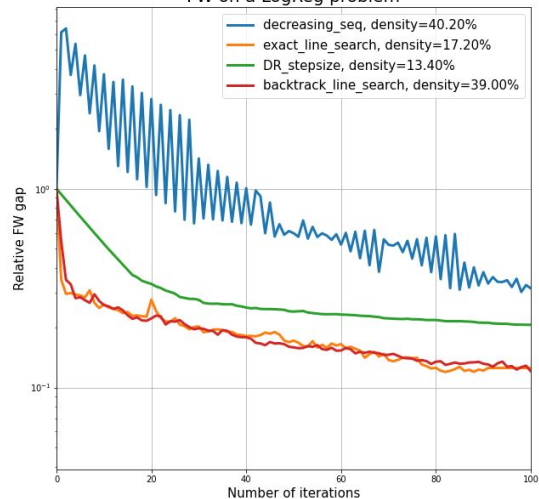


Modelon dataset logreg:

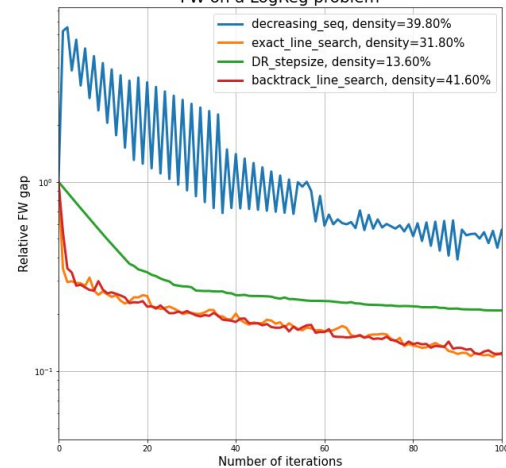
FW on a LogReg problem



FW on a LogReg problem



FW on a LogReg problem



Thank you for your attention!