# Step size selection in Frank Wolf

# Overview

# Problem statement
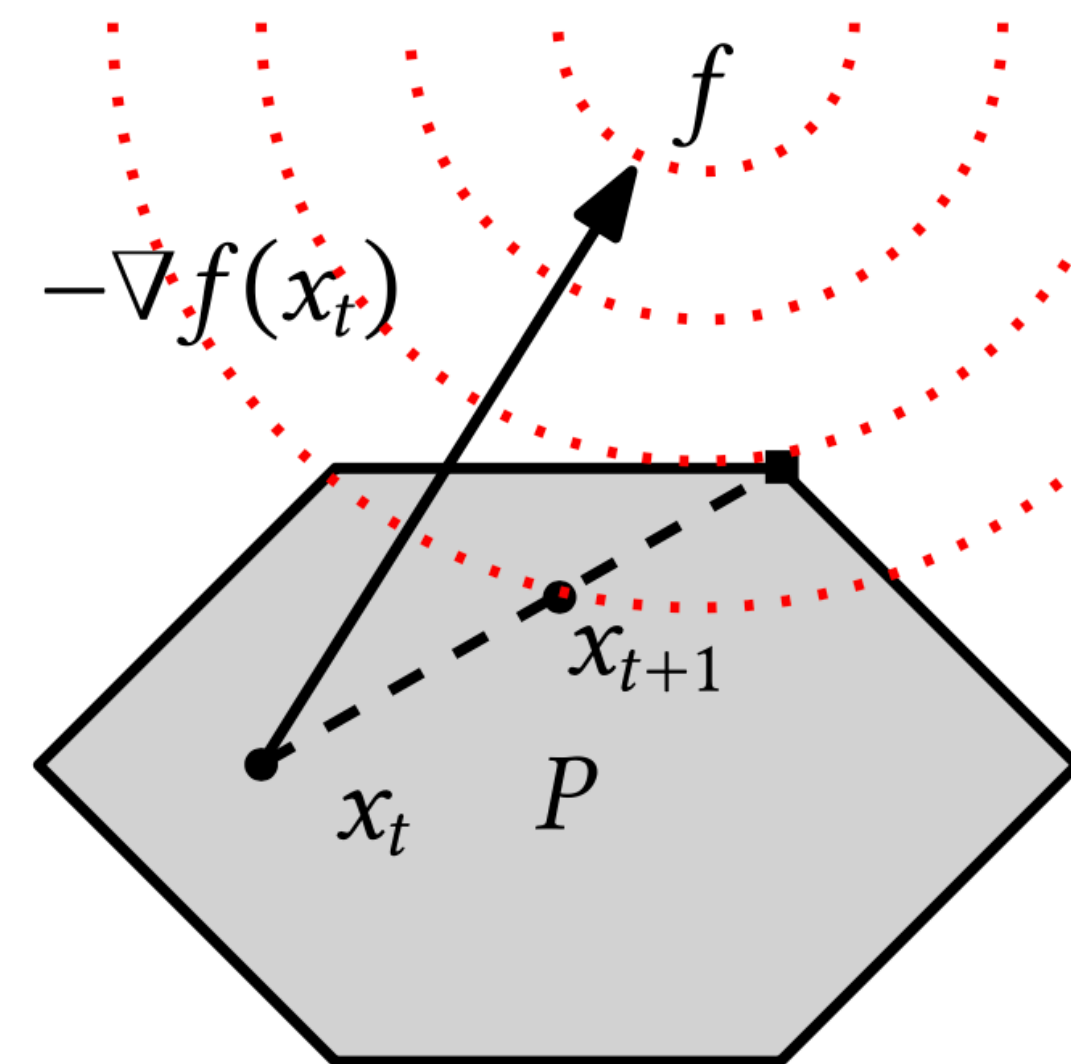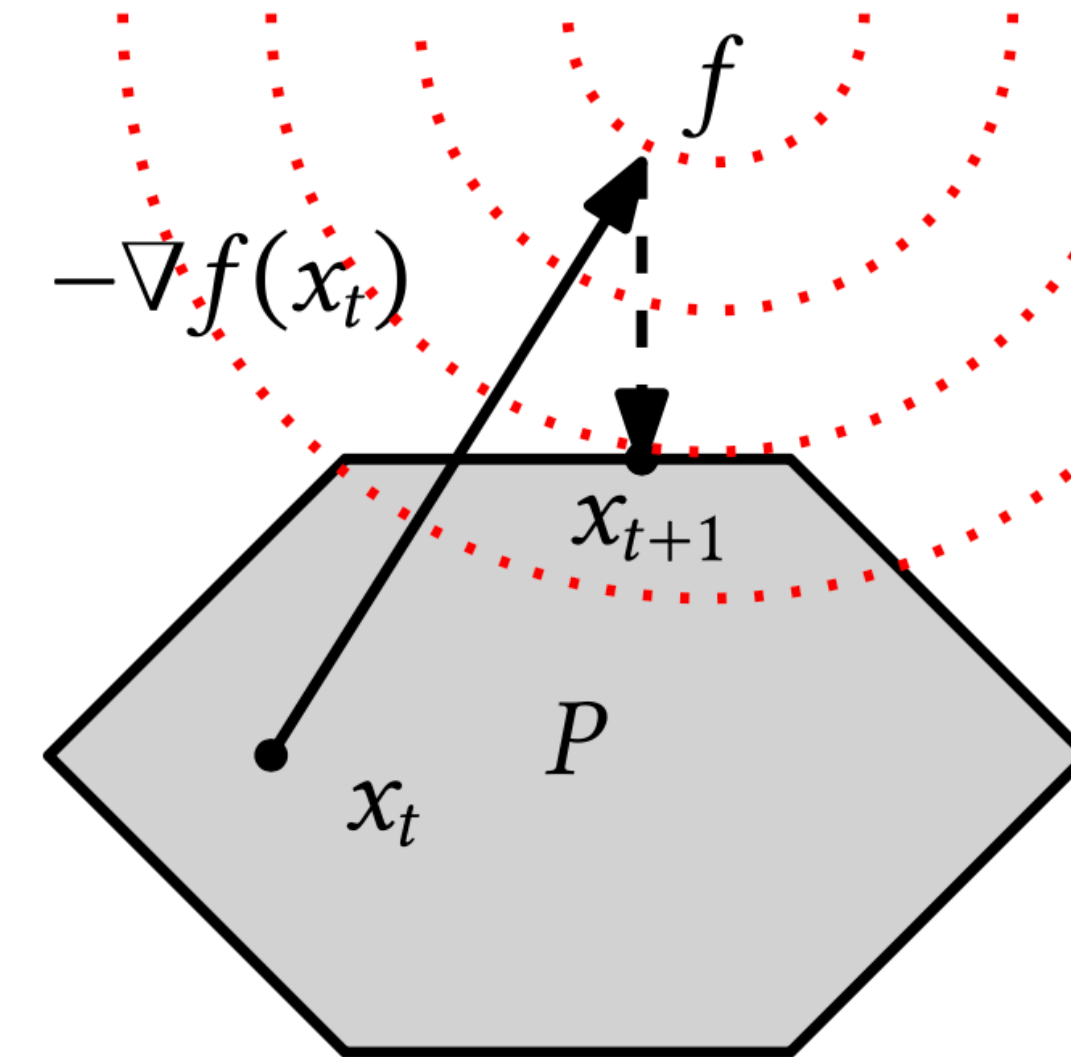
$$\text{Minimize } f(\mathbf{x})$$

$$\text{Subject to } \mathbf{x} \in \mathscr{D}$$

- $\mathscr{D}$ - **compact**, **convex** set in a vector space

- $f : \mathscr{D} \to \mathbf{R}$ is a **convex**, $L$**-smooth**, function

# Linear Minimization Oracle

$$s = \arg\min_{s \in \mathscr{D}} \quad s^T \nabla f(x_k)$$

- Finds a vector $s$ in feasible set $\mathscr{D}$, which aligns most with $\nabla f(x_k)$.

- Vector $s$ has the largest projection on $-\nabla f(x_k)$. Usually a vertex of the domain.

# Algorithm

1. $s_k = \arg \min\limits_{s \in \mathscr{D}} \ s^T \nabla f(x_k)$

2. $x_{k+1} = (1 - \gamma_k)x_k + \gamma_k s_k$

, where $\gamma_k \in [0,1]$ is a step-size.

- Both $x_k, s_k \in \mathscr{D}$. Convex combination of them is going to remain in the set. $x_{k+1} \in \mathscr{D}$
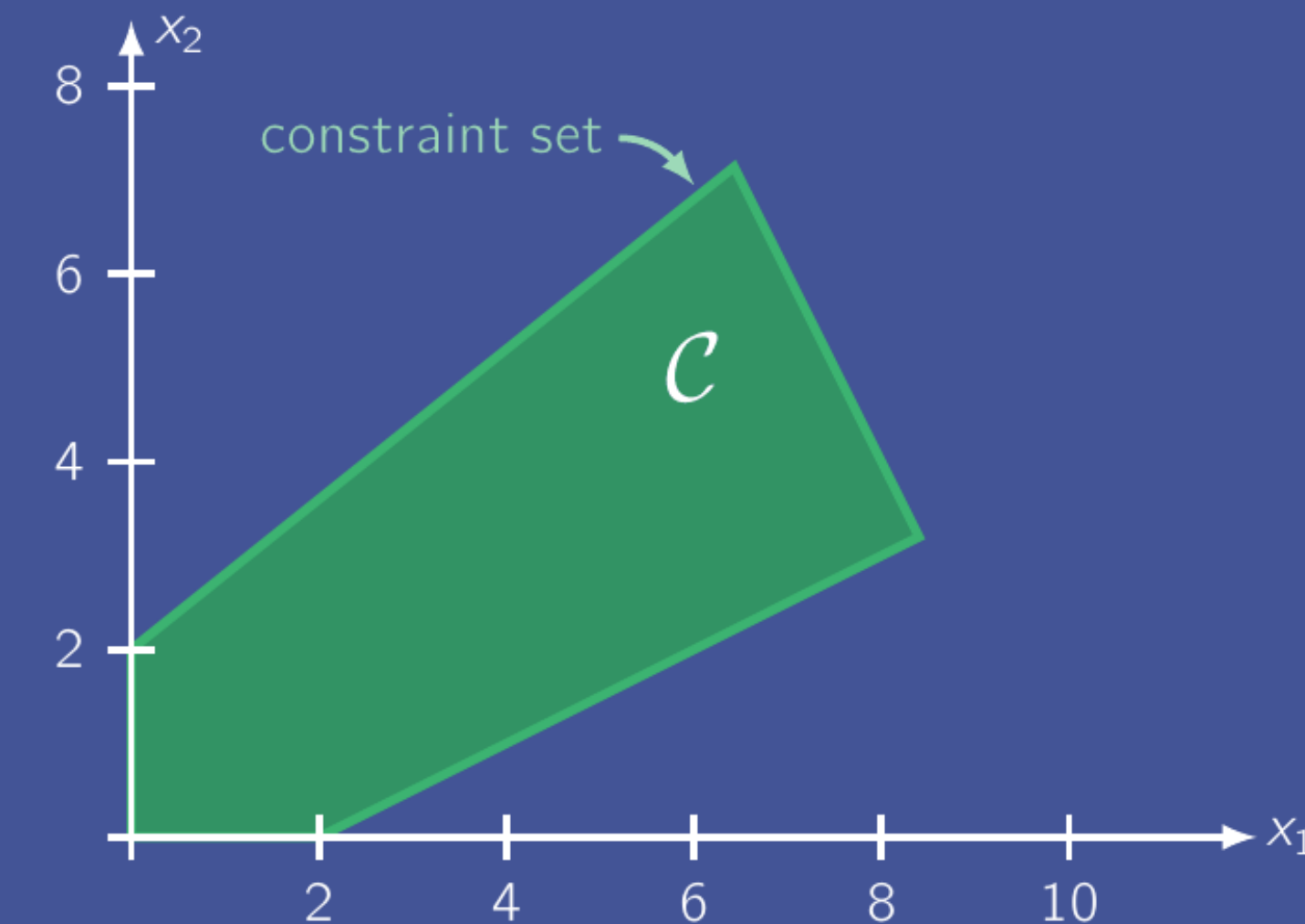
# Properties

- Convergence rate $O(1/k)$

- No need to do projection step (linear optimization vs. quadratic)

- Solve high dimensional problems

- Sparse solutions

- Designed for smooth, convex $f$

- Poor performance near optimum

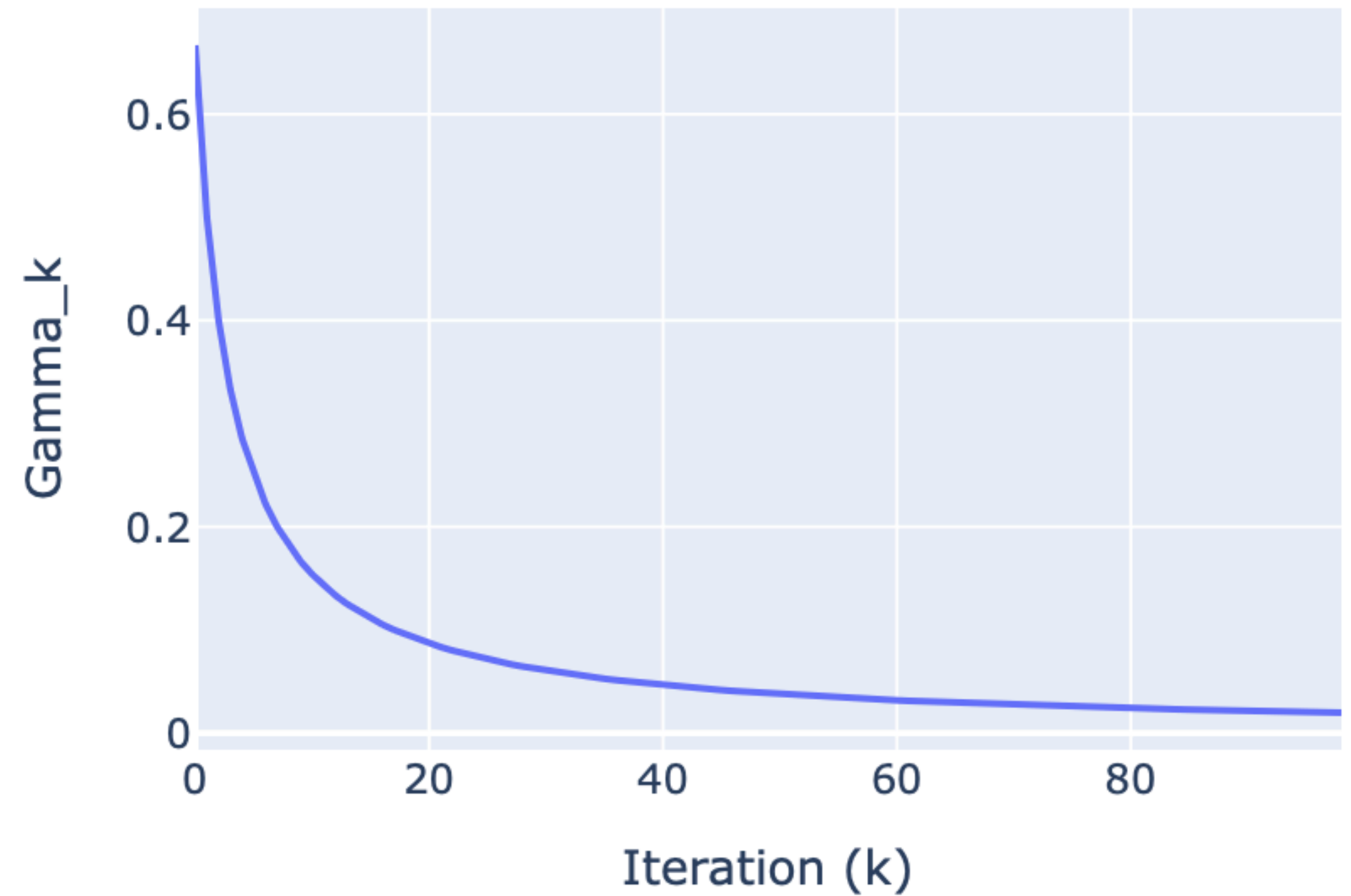- Complex, non-linear boundaries increase computation costs.

# Step size

# Line step size

$$\gamma_k = \frac{2}{k+2}$$

- Straightforward and cheap to compute.

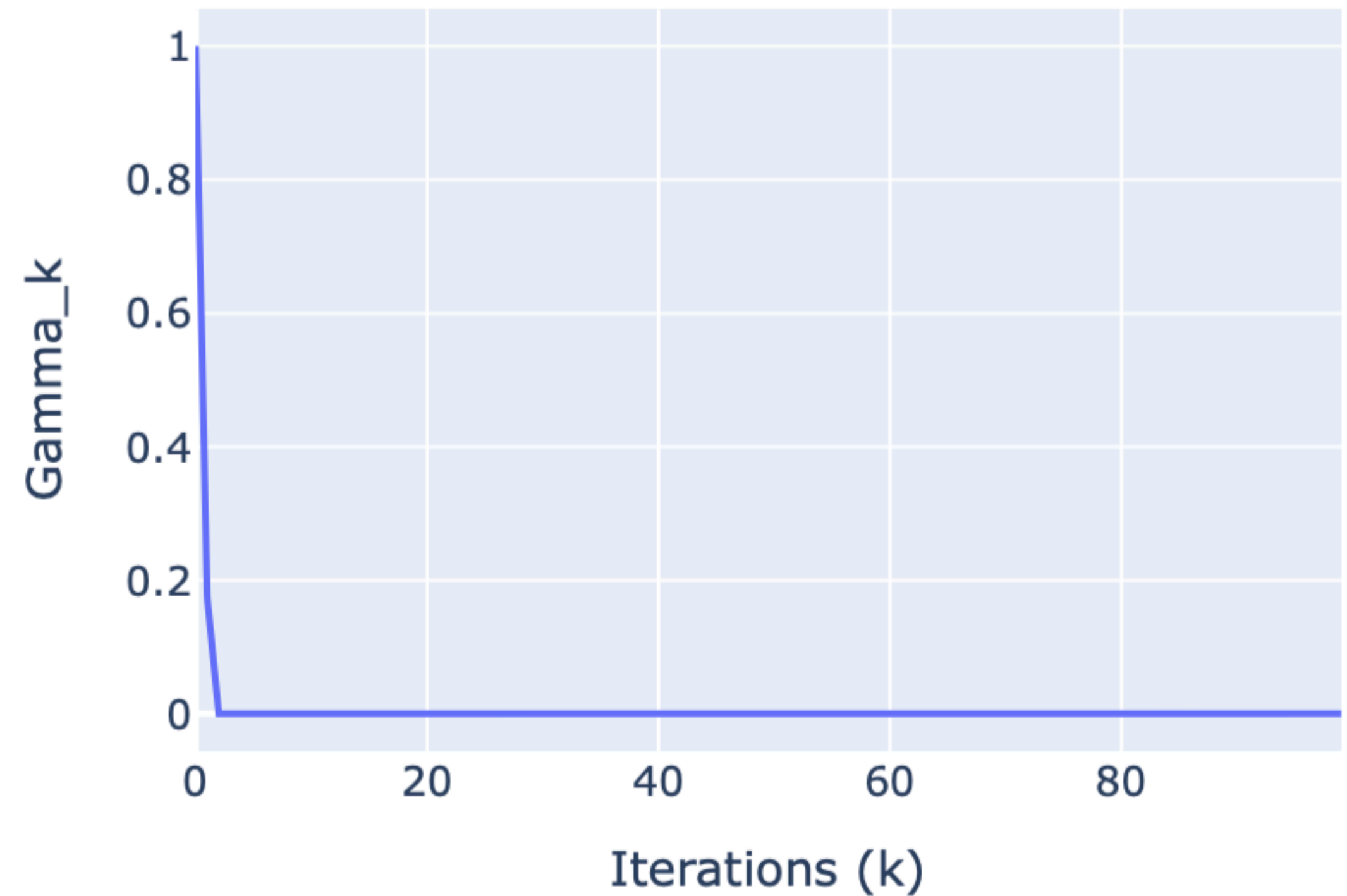- Slow convergence near optimum

- No function adaptation

# Exact line-search

$$\gamma_k = \arg\min f((x_k + \gamma_k(s_k - x_k))$$

- Ensures highest decrease per iteration

- Costly optimization problem
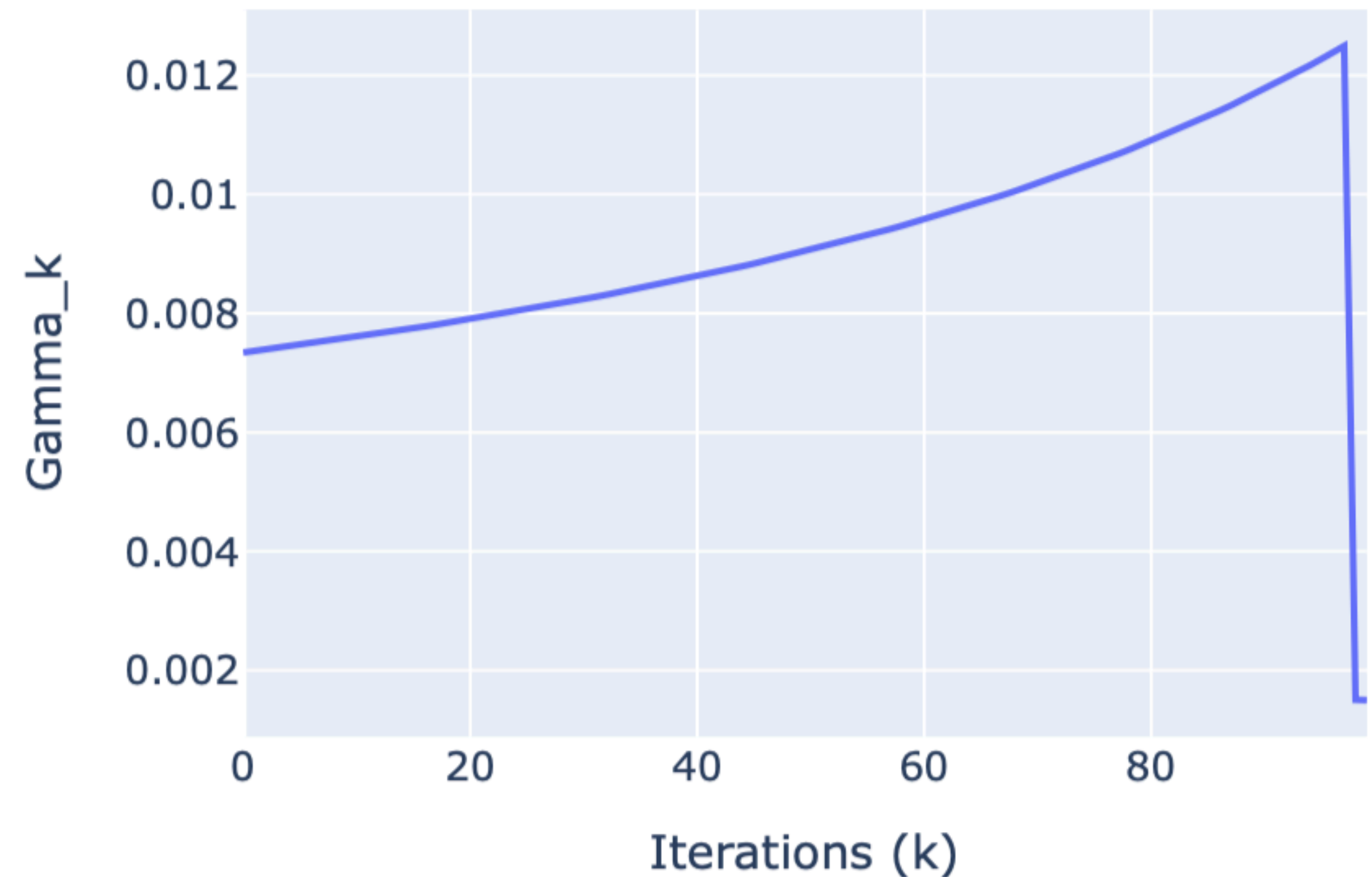
Exact line search

# Demyanov-Rubinov

$$\gamma_k = \min\{\frac{-\nabla f(x)^T(s_k - x_k)}{L\|s_k - x_k\|^2}, 1\}$$

- Goes to zero as we approach the optimum

- Responsive for geometry of $f$

- Require access to $L$

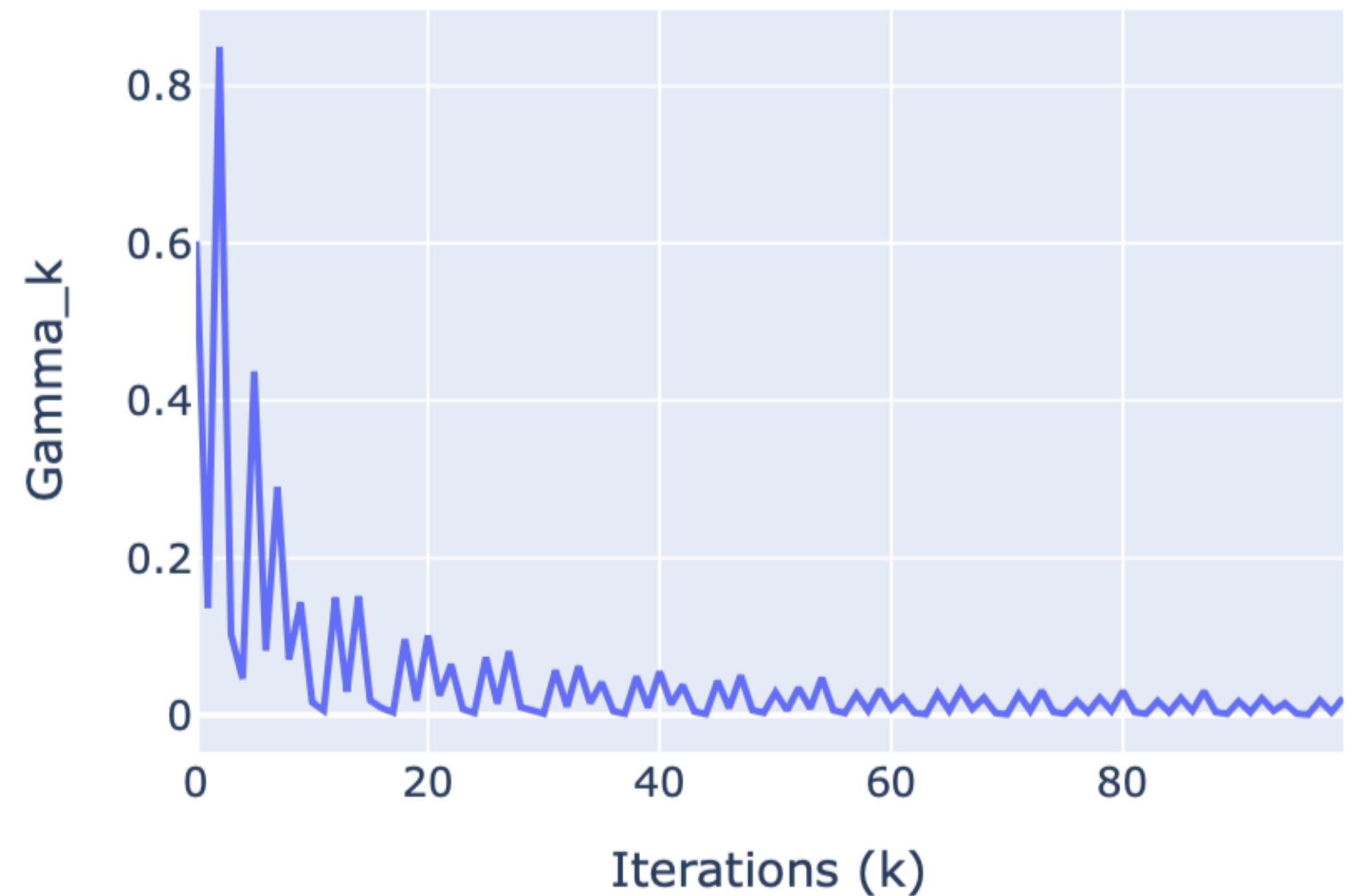- Unstable for small denominator (near optimum)



Demyanov Rubinov

# Backtracking line search

$$\gamma_k = \min\{\frac{-\nabla f(x)^T(s_k - x_k)}{M\|s_k - x_k\|^2}, 1\}$$

$M_t$ is approximation of $L$

- Doesn't require $L$

- Adaptive and stable progress

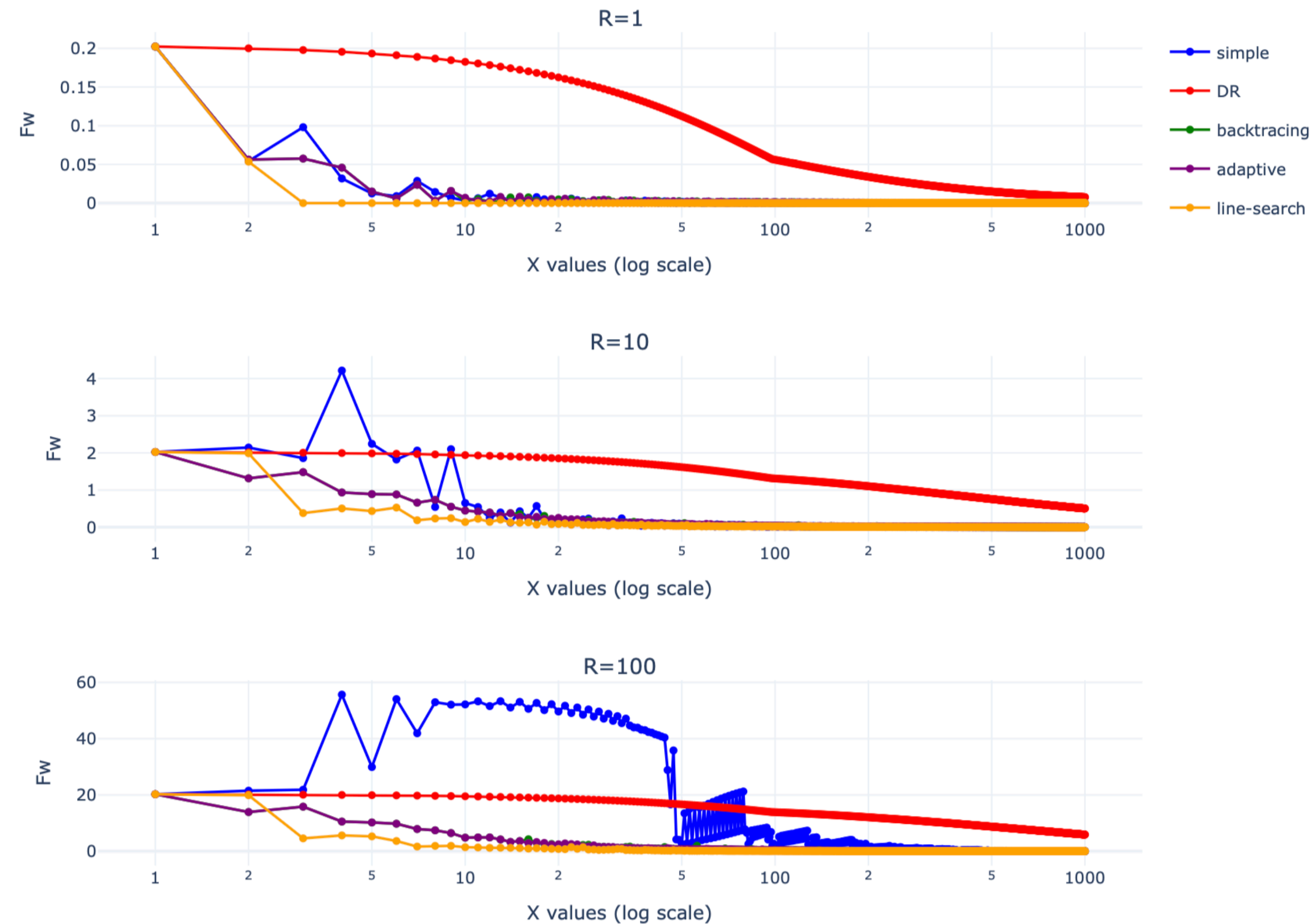- Require multiple evaluation of $f$



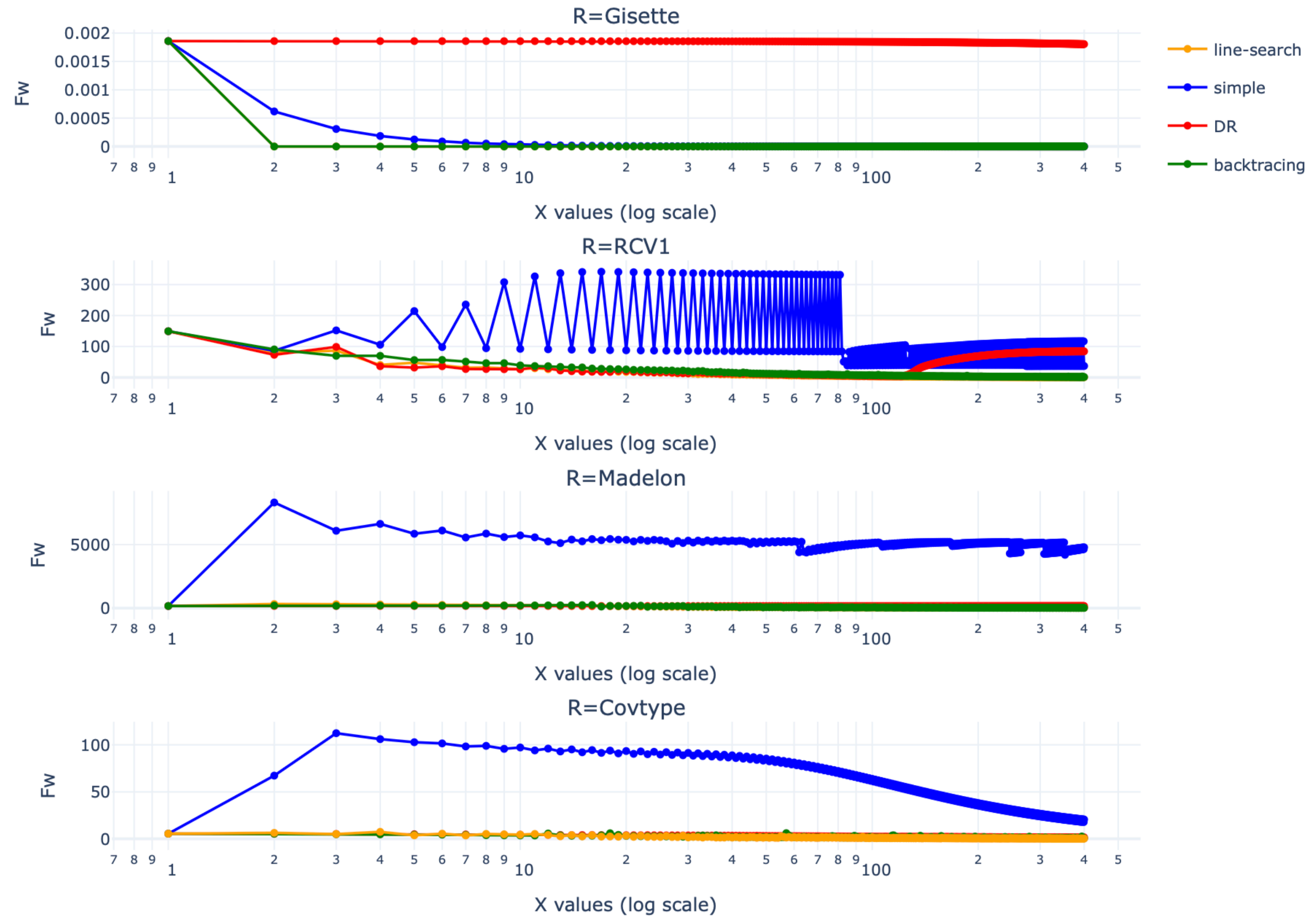Backtracking line search

# Experiments

# HW (Mushrooms)


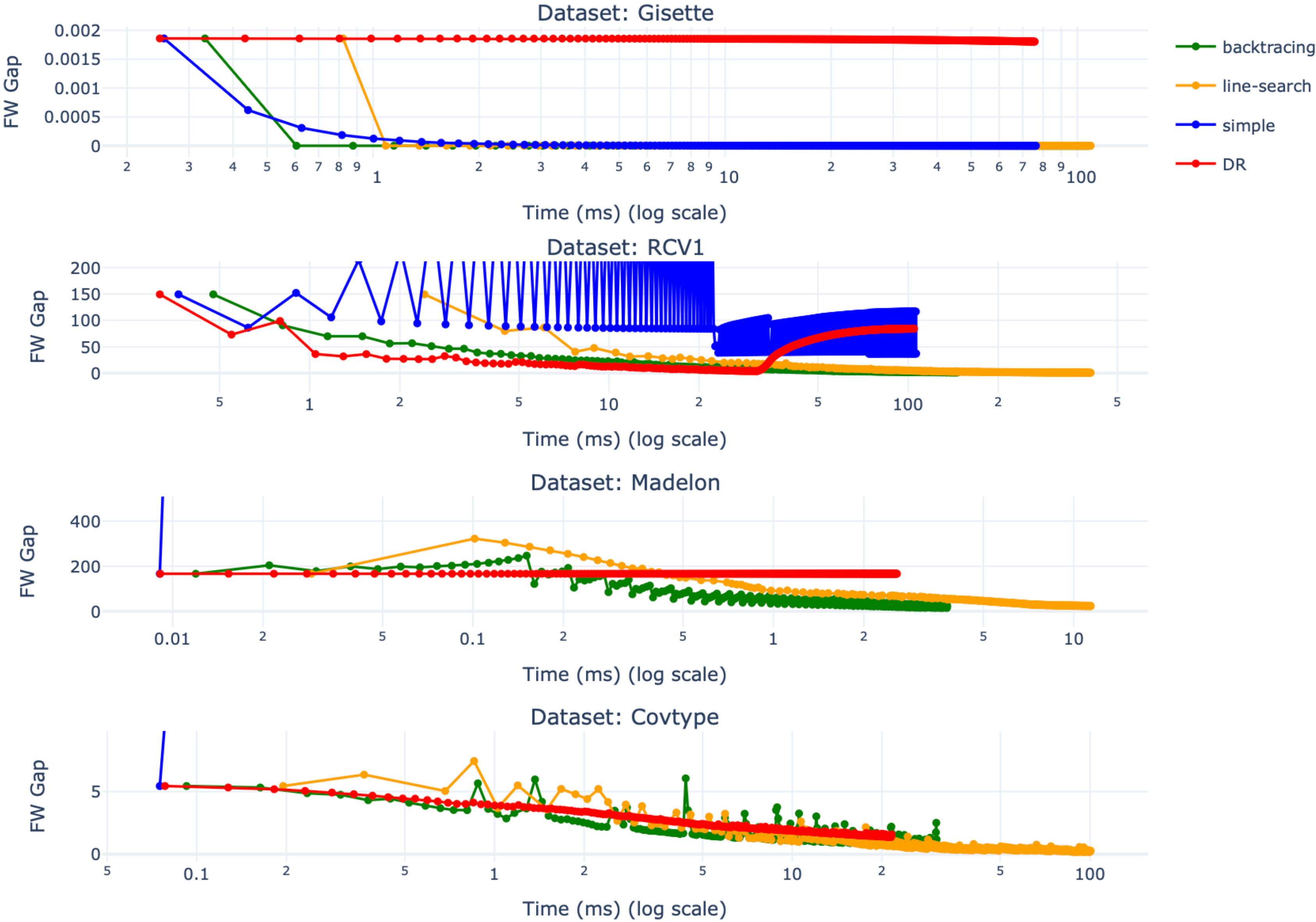
Line Plots for Metric: fw (Grouped Legend)

# Benchmark

- Gisette: Binary classification on 5000 features $R = 6e^{-3}$

- RCV1: Binary classification on 47236 features $R = 2e^4$

- Madelon: Binary classification on 500 features $R = 20$

- Covtype: Binary classification on 54 features $R = 200$

# Performance

|          | Simple           | DR              | Backtracking    | Line-Search     |
|----------|------------------|-----------------|-----------------|-----------------|
| **Gisette** | 1:16<br>5.24 it/s | 1:15<br>5.27 it/s | 1:49<br>3.66 it/s | 1:49<br>3.67 it/s |
| **RCV1**    | 1:43<br>3.87 it/s | 1:41<br>3.94 it/s | 2:25<br>2.75 it/s | 6:41<br>1 it/s |
| **Madelon** | 0:02<br>157.22 it/s | 0:02<br>159 it/s | 0:03<br>107 it/s | 0:11<br>35.84 it/s |
| **Covtype** | 0:21<br>18.84 it/s | 0:21<br>18.67 it/s | 0:30<br>13.02 it/s | 1:39<br>4.02 it/s |

Benchmark. Performance experiment

# Appendix

Line Plots for Metric: fw (Grouped Legend)

Mushrooms dataset. Performance experiment

# Cancer

- Breast Cancer:



Cancer dataset

Constraint: Trace

Constraint: L1