



Smooth Constraint Convex Minimization via Conditional Gradients

Sebastian Pokutta

*Zuse Institute Berlin
Technische Universität Berlin*

Nice, 09/2019

Joint Work with... (in random order)



*Alexandre
D'Aspremont*



*Thomas
Kerdreux*



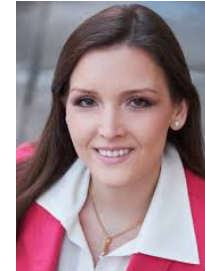
*Gábor
Braun*



*Cyrille
Combettes*



*Swati
Gupta*



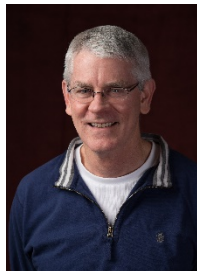
*Jelena
Diakonikolas*



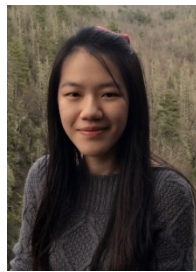
*Alejandro
Carderera*



*Robert
Hildebrand*



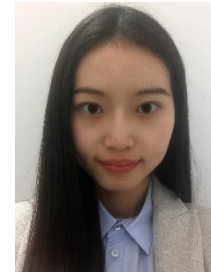
*Stephen
Wright*



*Yi
Zhou*



*George
Lan*



*Dan
Tu*



*Daniel
Zink*

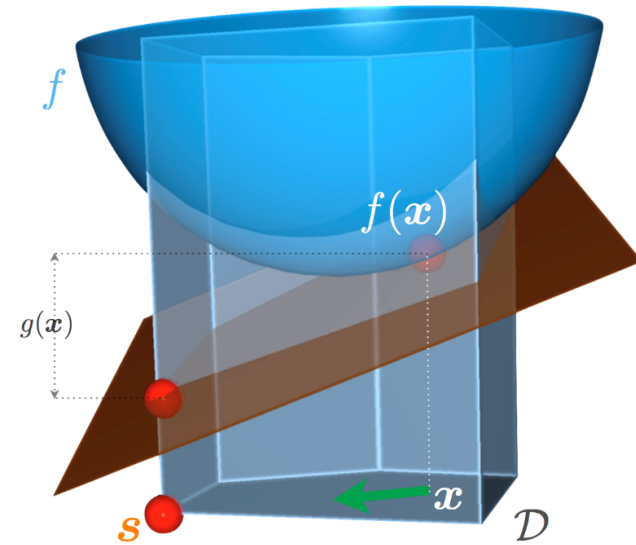
(Constraint) Convex Optimization

Convex Optimization:

Given a feasible region P solve the optimization problem:

$$\text{Min} \downarrow_{x \in P} f(x),$$

where f is a convex function (+ extra properties).



Source: [Jaggi 2013]

Our setup.

1. Access to P . **Linear Optimization (LO) oracle**: Given linear objective c
 $x \leftarrow \text{argmin} \downarrow_{v \in P} c^T v$
2. Access to f . **First-Order (FO) oracle**: Given x return
 $\nabla f(x)$ and $f(x)$

Why would you care for constraint convex optimization?

Setup captures various problems in Machine Learning, e.g.:

1. OCR (Structured SVM Training)

1. Marginal polytope over chain graph of letters of word and quadratic loss

2. Video Co-Localization

1. Flow polytope and quadratic loss

3. Lasso

1. Scaled ℓ_1 -ball and quadratic loss (regression)

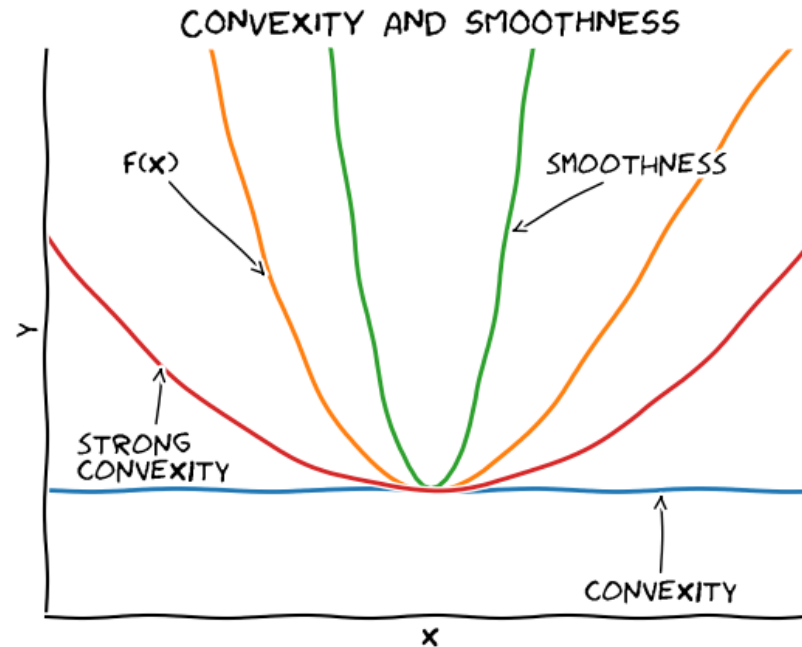
4. Regression over structured objects

1. Regression over convex hull of combinatorial atoms

5. Approximation of distributions

1. Bayesian inference, sequential kernel herding, ...

Smooth Convex Optimization 101



Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a function. We will use the following basic concepts:

Smoothness. $f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|x-y\|^2$

Convexity. $f(y) \geq f(x) + \nabla f(x)^T (y-x)$

Strong Convexity. $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|x-y\|^2$

=> Use unclear. Next step: Operationalize notions!

Measures of Progress: Smoothness and Idealized Gradient Descent

Consider an iterative algorithm of the form:

$$x_{t+1} \leftarrow x_t - \eta_t d_t$$

By definition of smoothness. $f(x_t) - f(x_{t+1}) \geq \eta_t \nabla f(x_t)^T d_t - \frac{L}{2} \eta_t^2 \|d_t\|^2$

Smoothness induces primal progress. Optimizing right-hand side:

$$f(x_t) - f(x_{t+1}) \geq (\nabla f(x_t)^T d_t)^2 / 2L \|d_t\|^2 \quad \text{for}$$

$$\eta_t \hat{=} (\nabla f(x_t)^T d_t) / L \|d_t\|$$

Idealized Gradient Descent (IGD). Choose $d_t \leftarrow x_t - x^*$ (non-det!)

$$f(x_t) - f(x_{t+1}) \geq (\nabla f(x_t)^T (x_t - x^*))^2 / 2L \|x_t - x^*\|^2$$

for $\eta_t \hat{=} (\nabla f(x_t)^T (x_t - x^*)) / L \|x_t - x^*\|$

Recall convexity: $f(y) \geq f(x) + \nabla f(x)^T (y - x)$

Primal bound from Convexity. $x \leftarrow x \downarrow t$ and $y \leftarrow x \uparrow^* \in \text{argmin}_{x \in P} f(x)$:

$$h \downarrow t := f(x \downarrow t) - f(x \uparrow^*) \leq \nabla f(x \downarrow t)^T (x \downarrow t - x \uparrow^*)$$

Plugging this into the progress from IGD and $\|x \downarrow t - x \uparrow^*\| \leq \|x \downarrow 0 - x \uparrow^*\|$.

$$f(x \downarrow t) - f(x \downarrow t+1) \geq (\nabla f(x \downarrow t)^T (x \downarrow t - x \uparrow^*))^2 / 2L \|x \downarrow t - x \uparrow^*\|^2 \geq h \downarrow t^2 / 2L \|x \downarrow 0 - x \uparrow^*\|^2$$

Rearranging provides contraction and convergence rate.

$$h \downarrow t+1 \leq h \downarrow t \cdot (1 - h \downarrow t / 2L \|x \downarrow 0 - x \uparrow^*\|^2) \Rightarrow h \downarrow T \leq 2L \|x \downarrow 0 -$$

Measures of Optimality: Strong Convexity

Recall strong convexity: $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|x-y\|^2$

Primal bound from Strong Convexity. $x \leftarrow x_{\downarrow t}$ and $y \leftarrow x_{\downarrow t} - \gamma(x_{\downarrow t} - x^*)$

$$h_{\downarrow t} := f(x_{\downarrow t}) - f(x^*) \leq (\nabla f(x_{\downarrow t})^T (x_{\downarrow t} - x^*))^2 / 2\mu \|x_{\downarrow t} - x^*\|^2$$

Plugging this into the progress from IGD.

$$f(x_{\downarrow t}) - f(x_{\downarrow t+1}) \geq (\nabla f(x_{\downarrow t})^T (x_{\downarrow t} - x^*))^2 / 2L \|x_{\downarrow t} - x^*\|^2 \geq \mu/L h_{\downarrow t}$$

Rearranging provides contraction and convergence rate.

$$h_{\downarrow t+1} \leq h_{\downarrow t} \cdot (1 - \mu/L) \Rightarrow h_{\downarrow T} \leq e^{-\mu/L T} \cdot h_{\downarrow 0}$$

From IGD to actual algorithms

Consider an algorithm of the form:

$$x_{t+1} \leftarrow x_t - \eta_t d_t$$

Scaling condition (Scaling). Show there exist α_t with

$$\nabla f(x_t)^T d_t / \|d_t\| \geq \alpha_t \nabla f(x_t)^T (x_t - x^*) / \|x_t - x^*\|$$

=> Lose an α_t^2 factor in iteration t . Bounds and rates follow.

Example. (Vanilla) Gradient Descent with $d_t \leftarrow -\nabla f(x_t)$

$$\nabla f(x_t)^T d_t / \|d_t\| = \|\nabla f(x_t)\|^2 \geq 1 \cdot \nabla f(x_t)^T (x_t - x^*) / \|x_t - x^*\|$$

Conditional Gradients (a.k.a. Frank-Wolfe Algorithm)

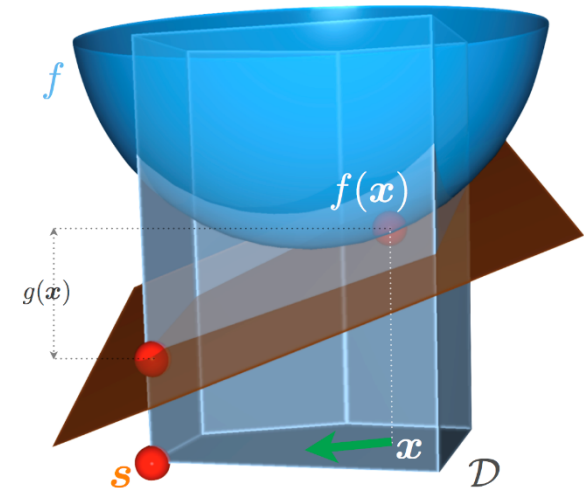
Conditional Gradients a.k.a. Frank-Wolfe Algorithm

Algorithm 1 Frank-Wolfe Algorithm [Frank and Wolfe, 1956]

Input: smooth convex f function with curvature C , $x_1 \in P$

Output: x_t points in P

- 1: **for** $t = 1$ **to** $T - 1$ **do**
 - 2: $v_t \leftarrow \text{LP}_P(\nabla f(x_t))$
 - 3: $x_{t+1} \leftarrow (1 - \gamma_t)x_t + \gamma_t v_t$ with $\gamma_t := \frac{2}{t+2}$
 - 4: **end for**
-



1. Advantages

1. **Extremely simple and robust:** no complicated data structures to maintain
2. **Easy to implement:** requires only a linear optimization oracle (first order method)
3. **Projection-free:** feasibility via linear optimization oracle
4. **Sparse distributions over vertices:** optimal solution is convex comb. (enables sampling)

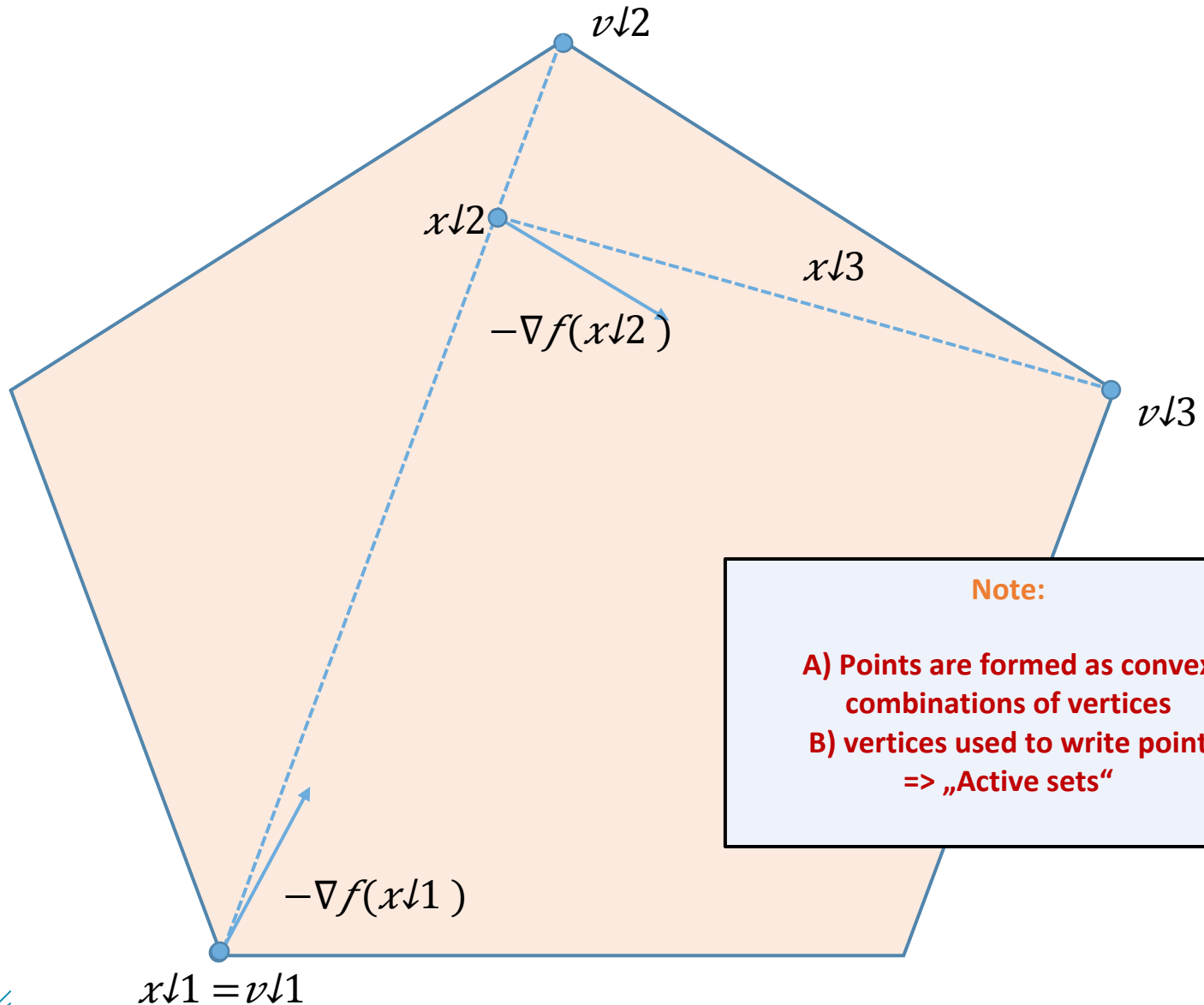
Source: [Jaggi 2013]

2. Disadvantages

1. Suboptimal convergence rate of $\mathcal{O}(1/T)$ in the worst-case

=> **Despite suboptimal rate often used because of simplicity**

Conditional Gradients a.k.a. Frank-Wolfe Algorithm



Note:

- A) Points are formed as convex combinations of vertices
- B) vertices used to write point \Rightarrow „Active sets“

Conditional Gradients a.k.a. Frank-Wolfe Algorithm

Algorithm 1 Frank-Wolfe Algorithm [Frank and Wolfe, 1956]

Input: smooth convex f function with curvature C , $x_1 \in P$

Output: x_t points in P

- 1: **for** $t = 1$ **to** $T - 1$ **do**
 - 2: $v_t \leftarrow \text{LP}_P(\nabla f(x_t))$
 - 3: $x_{t+1} \leftarrow (1 - \gamma_t)x_t + \gamma_t v_t$ with $\gamma_t := \frac{2}{t+2}$
 - 4: **end for**
-

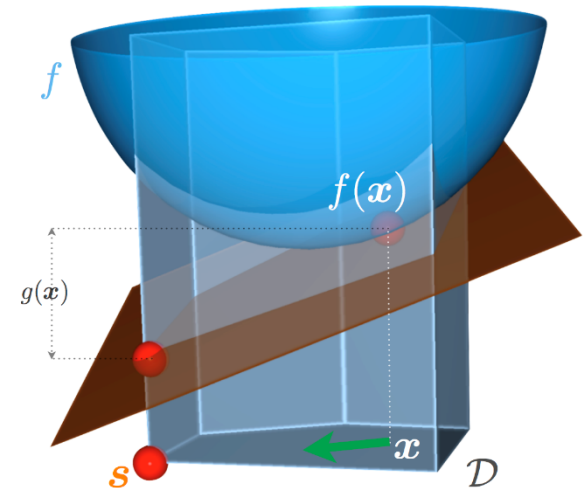
Establishing (Scaling).

FW algorithm takes direction $d \downarrow t = x \downarrow t - v \downarrow t$. Observe

$$\nabla f(x) \uparrow T (x \downarrow t - v \downarrow t) \geq \nabla f(x) \uparrow T (x \downarrow t - x \uparrow^*)$$

Hence with $\alpha \downarrow t = \|x \downarrow t - x \uparrow^*\| / D$ with D diameter of P :

$$\nabla f(x) \uparrow T (x \downarrow t - v \downarrow t) / \|x \downarrow t - v \downarrow t\| \geq \|x \downarrow t - x \uparrow^*\| / D \cdot \nabla f(x) \uparrow T (x \downarrow t - x \uparrow^*) / \|x \downarrow t - x \uparrow^*\|$$



Source: [Jaggi 2013]

The strongly convex case

Linear convergence in special cases

If f is strongly convex we would expect a linear rate of convergence.

Obstacle.

$$\nabla f(x)^\top (x_t - v_t) / \|x_t - v_t\| \geq \|x_t - x^*\| / D \cdot \nabla f(x)^\top (x_t - x^*) / \|x_t - x^*\|$$

Special case $x^* \in \text{rel.int}(P)$, say $B(x^*, 2r) \subseteq P$. Then:

Theorem [Marcotte, Guélat '86]. After a few iterations

$$\nabla f(x)^\top (x_t - v_t) / \|x_t - v_t\| \geq r/D \cdot \nabla f(x)^\top (x_t - x^*) / \|x_t - x^*\|$$

and linear convergence follows via (Scaling).

The strongly convex case Is linear convergence in general possible?

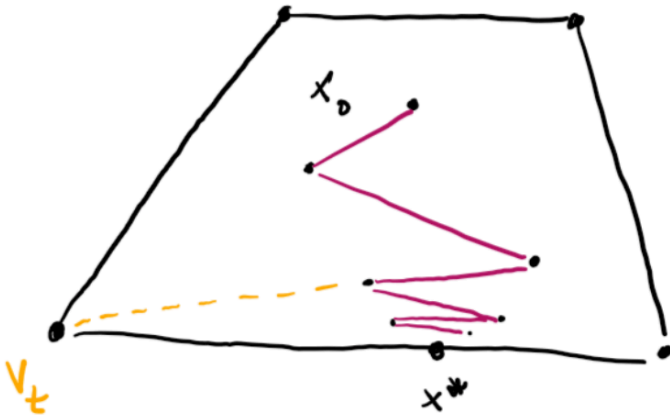
(Vanilla) Frank-Wolfe cannot achieve linear convergence in general:

Theorem [Wolfe '70]. x^* on boundary of P . For any $\delta > 0$ for infinitely many t :

$$f(x_t) - f(x^*) \geq 1/t^{1+\delta}$$

Issue: zig-zagging (b/c first order opt)

[Wolfe '70] proposed Away Steps



The strongly convex case

Linear convergence in general

First linear convergence result (in general)

[Garber, Hazan '13]

1. Simulating (theoretically efficiently) a stronger oracle rather using Away Steps
2. Involved constants are *extremely large* => algorithm unimplementable

Linear convergence for implementable variants

[Lacoste-Julien, Jaggi '15]

1. (Dominating) Away-steps are enough
2. Includes most known variants: Away-Step FW, Pairwise CG, Fully-Corrective FW, Wolfe's algorithm, ...
3. Key ingredient: There exists $w(P)$ (depending on **polytope** P (only!)) s.t.

$$\nabla f(x)^T (a \downarrow t - v \downarrow t) \geq w(P) \nabla f(x)^T (x \downarrow t - x^*) / \|x \downarrow t - x^*\|$$

($d \downarrow t = a \downarrow t - v \downarrow t$ is basically the direction that either variant dominates)

=> **Linear convergence via (Scaling)**

Many more variants and results...

Recently there has been a lot of work on Conditional Gradients, e.g.,

1. Linear convergence for conditional gradient sliding [Lan, Zhou '14]
2. Linear convergence for (some) non-strongly convex functions [Beck, Shtern '17]
3. Online FW [Hazan, Kale '12, Chen et al '18]
4. Stochastic FW [Reddi et al '16] and Variance-Reduced Stochastic FW [Hazan, Luo '16, Chen et al '18]
5. In-face directions [Freund, Grigas '15]

... and *many more!!*

=> **Very competitive and versatile in real-world applications**

Revisiting Conditional Gradients Lazification

Bottleneck 1: Cost of Linear Optimization

Drawbacks in the context of hard feasible regions

Basic assumption of conditional gradient methods:

Linear Optimization is cheap

As such accounted for as $\mathcal{O}(1)$. This assumption is **not warranted** if:

1. Linear Program of feasible region is huge
 1. Large shortest path problems
 2. Large scheduling problems
 3. Large-scale learning problems
2. Optimization over feasible region is NP-hard
 1. TSP tours
 2. Packing problems
 3. Virtually *every* real-world combinatorial optimization problem

Rethinking CG in the context of expensive oracle calls

Basic assumption for us:

*Linear Optimization is **not** cheap*

(Think: hard IP can easily require an hour to be solved => one call/it unrealistic)

1. Questions:

1. Is it necessary to call the oracle in each iteration?
2. Is it necessary to compute (*approximately*) optimal solutions?
3. Can we reuse information?

2. Theoretical requirements

1. Achieve identical convergence rates, otherwise any speedup will be washed out

3. Practical requirements

1. Make as few oracle calls as possible

Lazification approach using weaker oracle

Oracle 1 Weak Separation Oracle $\text{LPsep}_P(c, x, \Phi, K)$

Require: $c \in \mathbb{R}^n$ linear objective, $x \in P$ point, $K \geq 1$ accuracy, $\Phi > 0$ objective value;

Ensure: Either (1) $y \in P$ vertex with $c(x - y) > \Phi/K$, or (2) **false:** $c(x - z) \leq \Phi$ for all $z \in P$.

1. Interpretation of Weak Separation Oracle: *Discrete Gradient Directions*

1. Either a new point $y \in P$ that improves the current objective by at least Φ/K (positive call)
2. Or it asserts that all other points $z \in P$ improve no more than Φ (negative call)

2. Lazification approach

[Braun, P., Zink '17]

1. Use weaker oracle that allows for caching and early termination (no more expensive than LP)
2. Advantage: huge speedups in wall-clock time when LP is hard to solve
 1. For hard LPs speedups can be as large as 10^7
3. Disadvantage: weak separation oracle produces even weaker approx. than LP oracle
 1. Actual progress in iterations can be worse than with LP oracle
 2. Advantage vanishes if LP is very cheap and can be worse than original algorithm
 3. Caching is not "smart": it simply iterates over the already seen vertices

3. Optimal complexity for Weak Separation Oracle

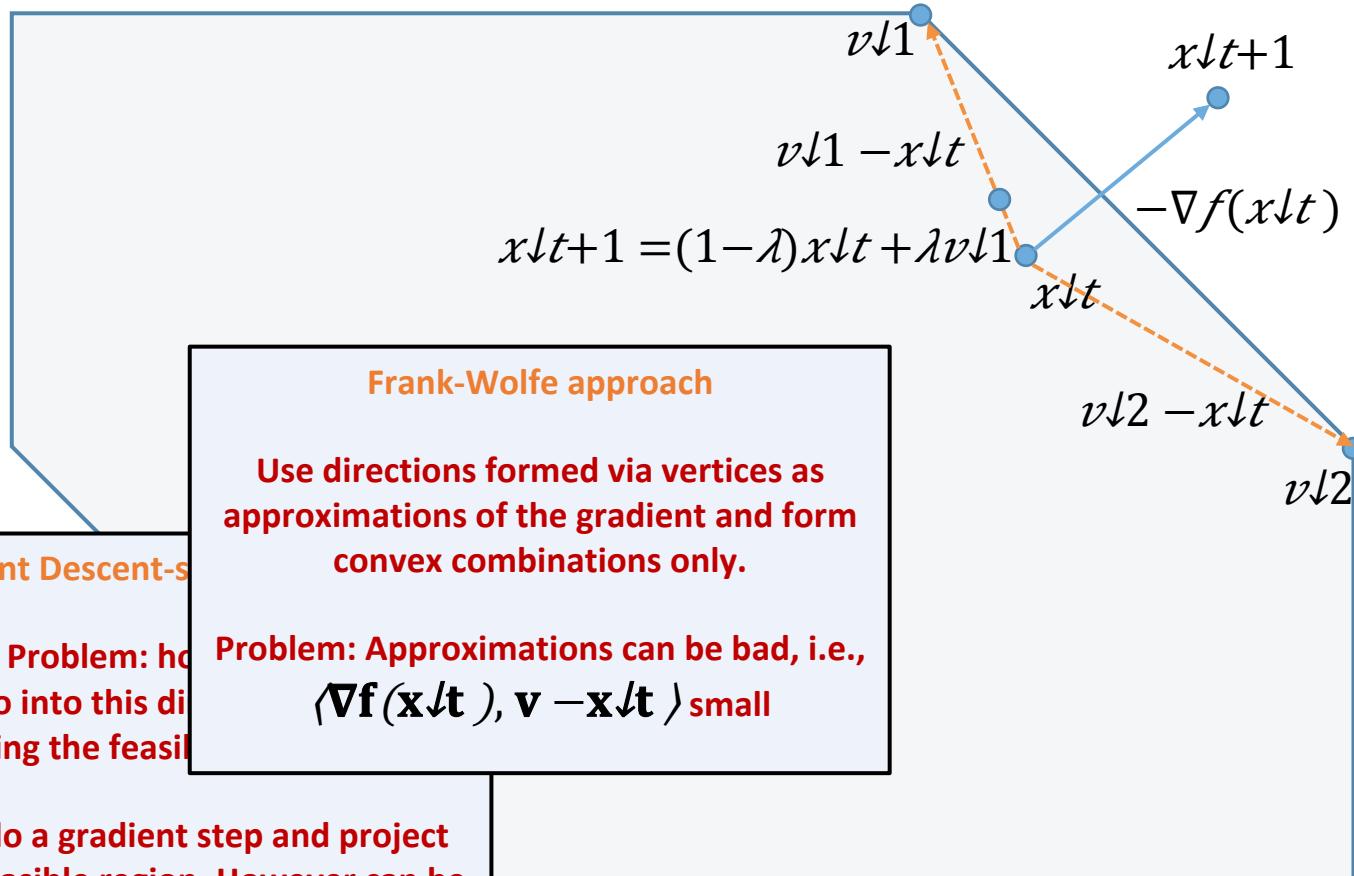
[Braun, Lan,

P., Zhou '17]

Revisiting Conditional Gradients Blending

Bottleneck 2: Quality of gradient approximation

Frank-Wolfe vs. Projected Gradient Descent



Frank-Wolfe approach

Use directions formed via vertices as approximations of the gradient and form convex combinations only.

Problem: Approximations can be bad, i.e., $\langle \nabla f(x_t), v - x_t \rangle$ small

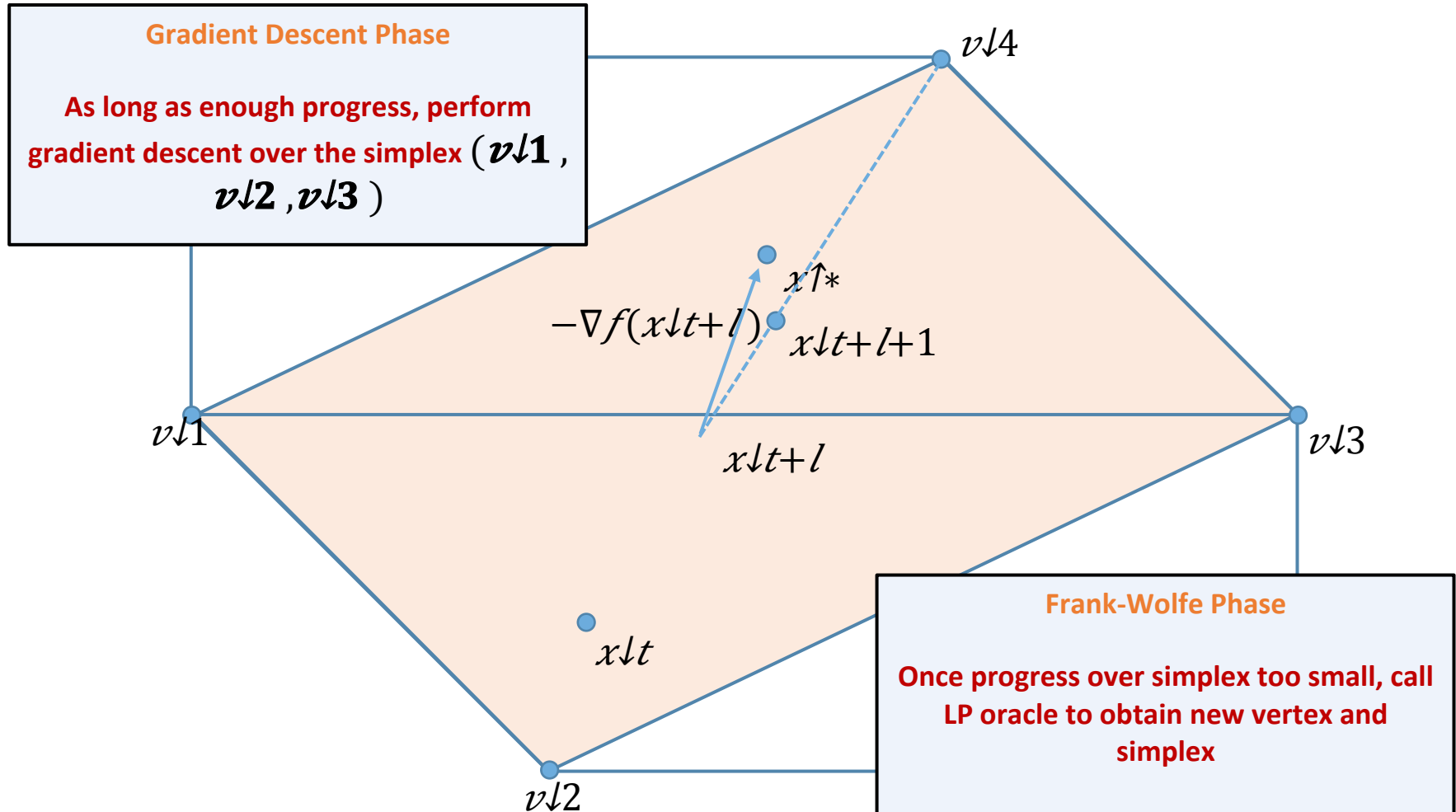
Gradient Descent-s

Problem: how can we go into this direction without leaving the feasible region?

Solution: do a gradient step and project back into feasible region. However can be very expensive

⇒ Tradeoff between ensured feasibility and quality of gradient approximations!

Blending of gradient steps and Frank-Wolfe steps



Main Theorem

You basically get what you expect.

Theorem. [Braun, P., Tu, Wright '18] Assume f is convex and smooth over the polytope P with curvature C and geometric strong convexity μ . Then Algorithm 1 ensures:

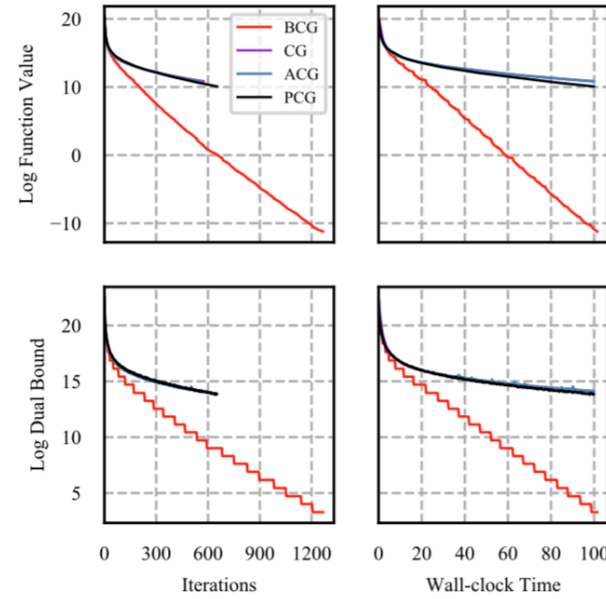
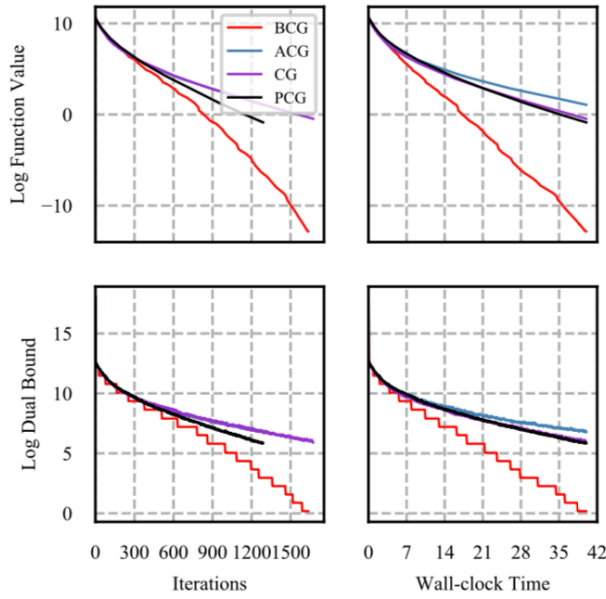
$$f(x_t) - f(x^*) \leq \varepsilon \quad \text{for } t \geq \Omega\left(\frac{C}{\mu} \log \frac{\Phi}{\varepsilon}\right),$$

where x^* is an optimal solution to f over P and $\Phi \geq f(x_0) - f(x^*)$.

(For previous empirical work with similar idea see also [Rao, Shah, Wright '15])

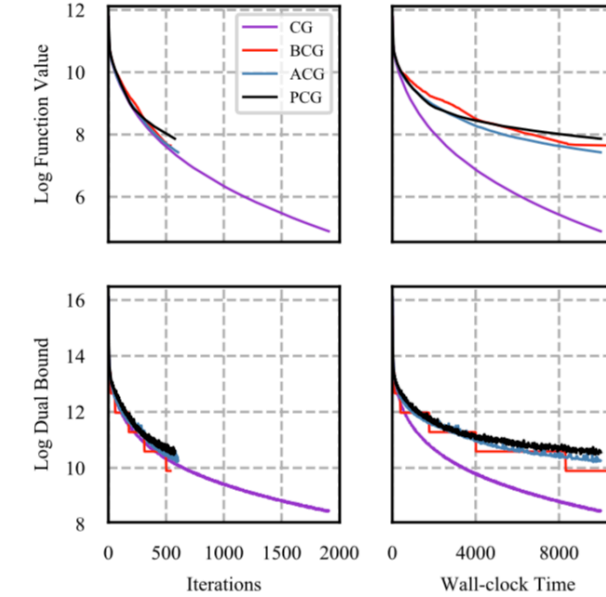
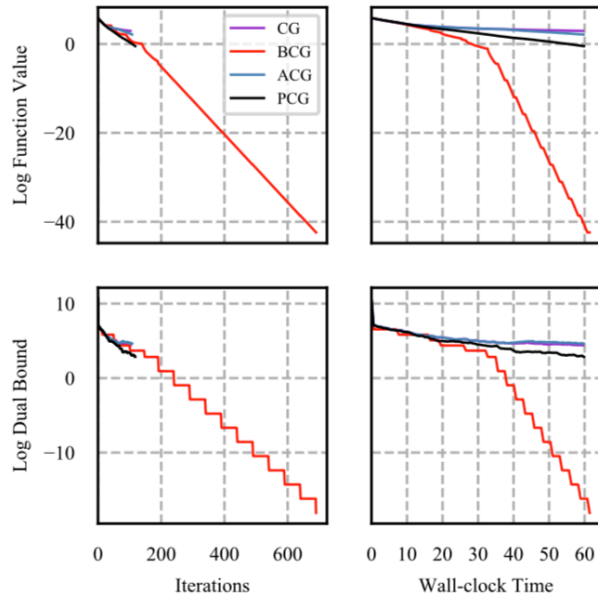
Computational Results

Lasso



Structured Regression

Sparse Signal Recovery



Matrix Completion

Revisiting Conditional Gradients Acceleration

How about acceleration?

The problem. Rates from standard proofs do not match known lower bounds:

1. Smooth convex case: $O(1/\varepsilon)$ vs. $\Omega(1/\sqrt{\kappa}\varepsilon)$
2. Smooth strongly convex case: $O(\mu/L \log 1/\varepsilon)$ vs. $\Omega(\sqrt{\kappa}\mu/L \log 1/\varepsilon)$

Acceleration closes this gap. Various approaches:

1. Polyak's Heavy Ball method
2. Nemirovski Acceleration with Line Search
3. Nesterov Acceleration
4.

Limits to Acceleration for LP-based Methods

Lower bound. Consider the optimization problem:

[Jaeggi 2013, Lan 2013]

$$\min_{x \in \Delta(n)} \|x\|_2^2$$

where $\Delta(n) = \{x \in \mathbb{R}^n_+ \mid \sum x_i = 1\}$ probability simplex.

Now, after k iterations the primal gap h_k is lower bounded as follows:

$$h_k \geq 1/k - 1/n$$

1. Smooth convex: After $n/2$ iterations $h_{n/2} \geq 1/n$
=> Vanilla FW rate is optimal (up to constant factors)
2. Smooth strongly convex: If $h_t \leq h_0 (1-r)^t$, then $r \leq 2 \log n/n$
=> Away-Step FW rate of $(1 - 1/8n)$ is optimal (up to log factors)

Acceleration Beyond the Dimension Threshold?

Basic idea: The lower bound limits acceleration only up to the dimension. However, if we seek an *accelerated global rate* of the form:

$$h \downarrow t \leq h \downarrow 0 (1-r)^{\uparrow t}, \text{ then } r \leq 2 \log n/n ,$$

i.e., the lower bound also limits rates *beyond* the dimension threshold.

In a nutshell: We can design an algorithm that runs a **constant number** $T \downarrow 0$ of unaccelerated steps and then has “true” acceleration kick in.

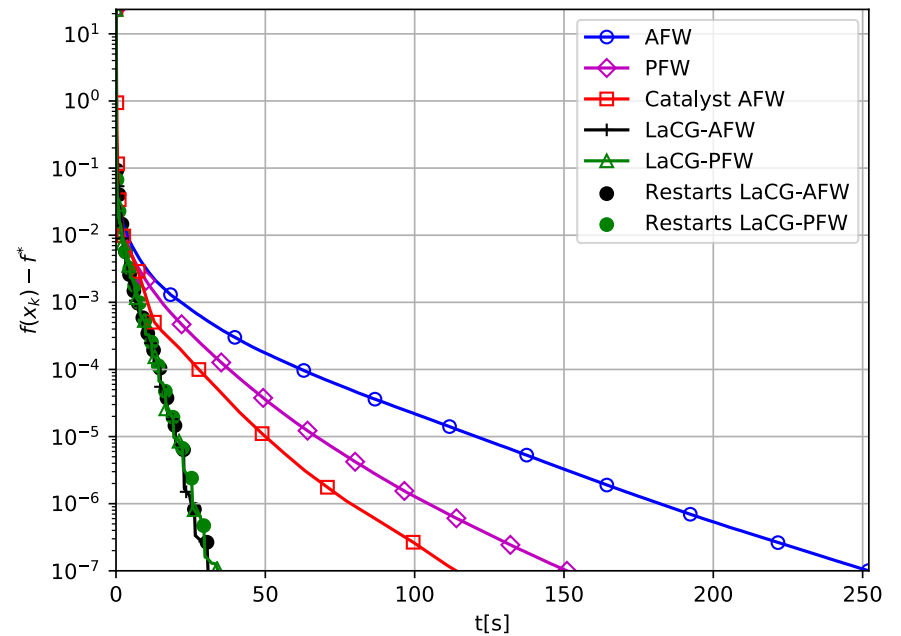
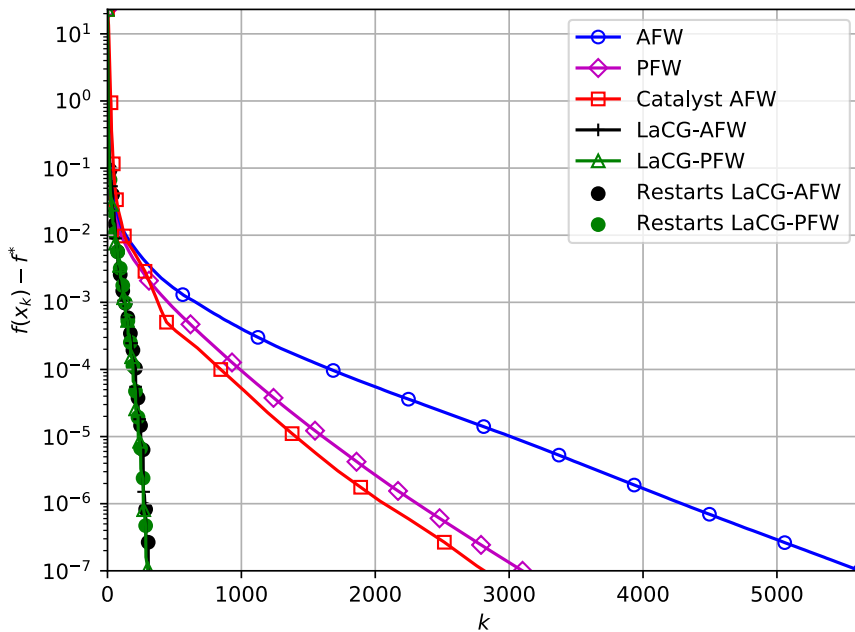
=> Asymptotically optimal rate. Roughly:

[Carderera, Diakonikolas, P. '19]

$$h \downarrow t \leq h \downarrow 0 (1 - \sqrt{\square \mu/L})^{\uparrow t - T \downarrow 0}$$

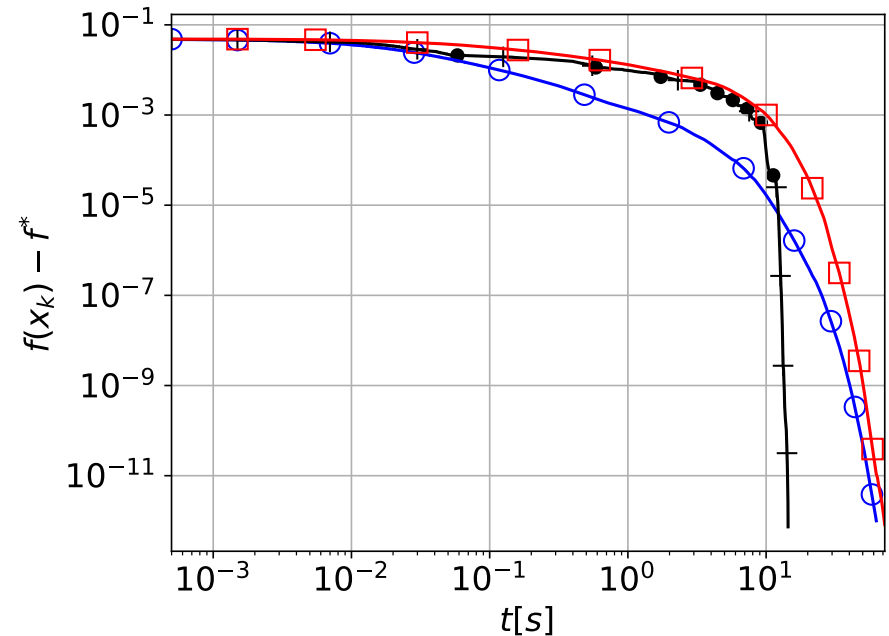
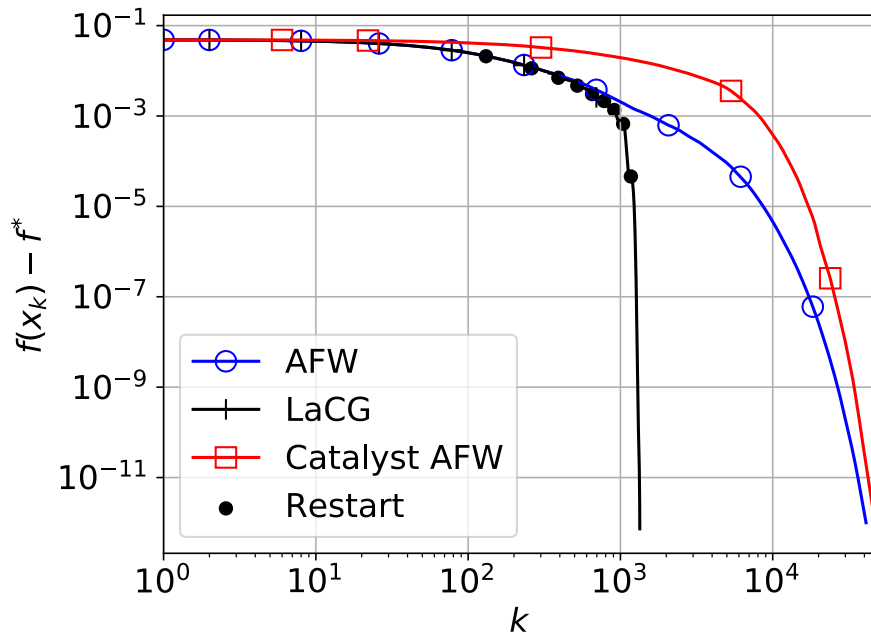
Preliminary Computational Results

Setup: Quadratic over Birkhoff Polytope
=> small dim-dependent term



Preliminary Computational Results

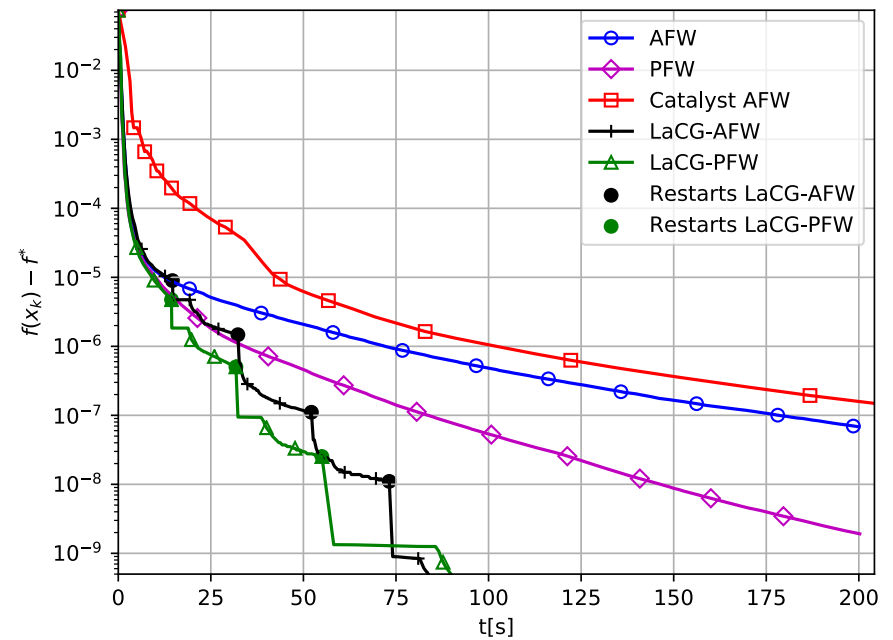
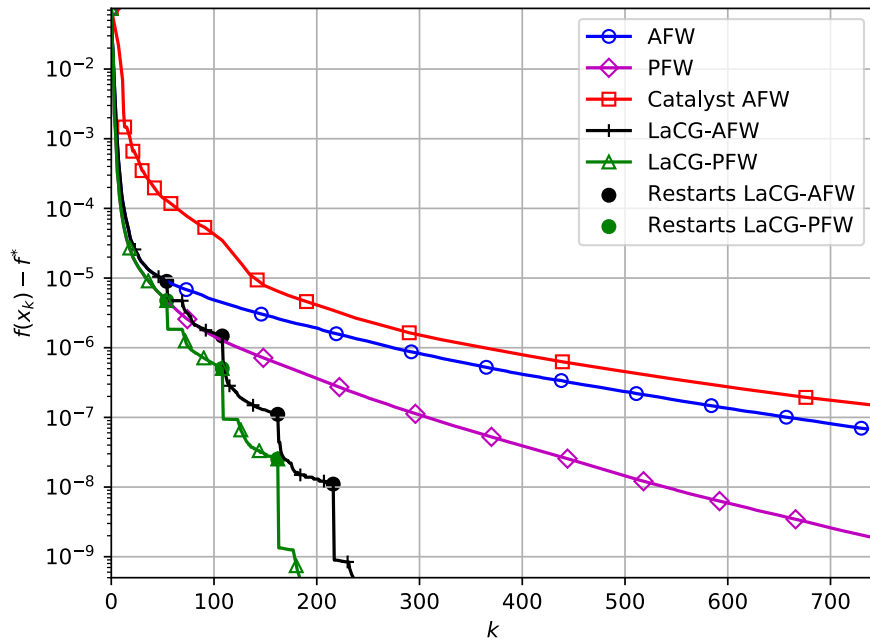
Setup: Quadratic over Probability Simplex (dim = 1000)
=> large dim-dependent term / lower bound instance



log-log scale

Preliminary Computational Results

Setup: Video Co-Localization



Want to know more?

Upcoming survey online in the next few weeks:

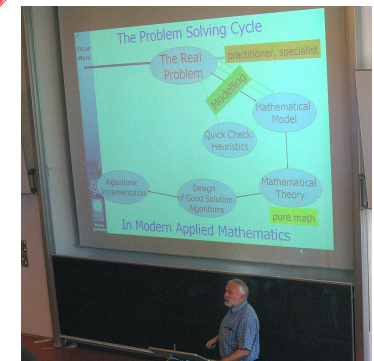
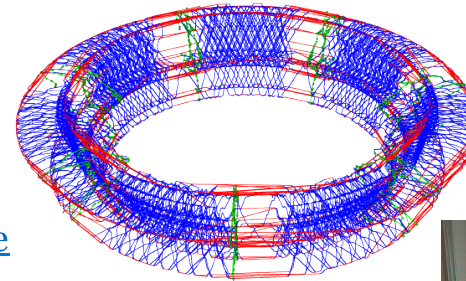
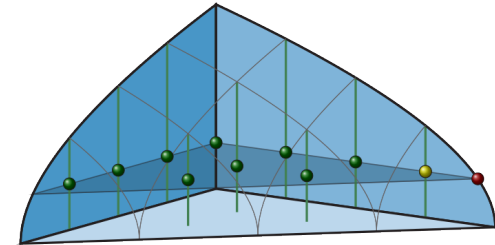
[Carderera, Combettes, P. “Conditional Gradients” ’19+]

Announcement: Combinatorial Optimization at Work

A summer school of TU Berlin in cooperation with MATH+ and the Berlin Mathematical School

Everything you always wanted to know about LP/MIP and real-world industrial applications (lectures and exercises)

- Dates of the course: **September 14 – 26, 2020**
- Language: **English**
- Location: **Zuse Institute Berlin**
- Application deadline: **June 14, 2020**
- Participation fee: **none**
- URL (info/application): **<http://co-at-work.zib.de>**
- Intended audience: **master/PhD students, Post-docs**
- Contact: **coaw@zib.de**
- Lectures by: **the SCIP team, developers of Xpress, Gurobi, Gams, and many more**





Smooth Constraint Convex Minimization via Conditional Gradients

Sebastian Pokutta

*Zuse Institute Berlin
Technische Universität Berlin*

Nice, 09/2019