A fresh look at the Frank-Wolfe algorithm with applications to sparse convex optimization

[J. 2011 PhD Thesis, summary: J. 2013 ICML Paper] Martin Jaggi

Ecole Polytechnique, Paris ICCOPT 2013 / 08 / 01 Constrained Convex Optimization



Constrained Convex Optimization

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









The Linearized Problem

$$\min_{\boldsymbol{s'}\in\mathcal{D}} f(\boldsymbol{x}) + \left\langle \boldsymbol{s'} - \boldsymbol{x}, \nabla f(\boldsymbol{x}) \right\rangle$$

f(x)

 \mathbf{x}

 $\mathcal{D} \subset \mathbb{R}^d$

S

Algorithm 1 Frank-Wolfe

Let
$$x^{(0)} \in \mathcal{D}$$

for $k = 0 \dots K$ do
Compute $s := \underset{s' \in \mathcal{D}}{\operatorname{arg min}} \langle s', \nabla f(x^{(k)}) \rangle$
Let $\gamma := \frac{2}{k+2}$
Update $x^{(k+1)} := (1 - \gamma)x^{(k)} + \gamma s$
end for

The Linearized Problem

$$\min_{\boldsymbol{s'} \in \mathcal{D}} f(\boldsymbol{x}) + \left\langle \boldsymbol{s'} - \boldsymbol{x}, \nabla f(\boldsymbol{x}) \right\rangle$$

Algorithm 1 Frank-Wolfe Let $x^{(0)} \in \mathcal{D}$ for $k = 0 \dots K$ do Compute $s := \underset{s' \in \mathcal{D}}{\operatorname{arg min}} \langle s' \cdot \nabla$ Let $\gamma := \frac{2}{k+2}$ Update $x^{(k+1)} := (1 - \frac{1}{2})$ end for



Two kinds of first-order methods



f(x)

 \mathbf{x}

 $\mathcal{D}\subset \mathbb{R}^d$

other Frank-Wolfe talks at ICCOPT '13:

Pradeep Ravikumar (Mon), Federico Pierucci (Tue), Guanghui Lan (Wed), Zaid Harchaoui (Wed), Simon Lacoste-Julien (Wed), Robert M. Freund (Wed), Paul Grigas (Wed), Ronny Luss (Wed), Dan Garber (Wed), (and summer school by Steve Wright)

	Frank-Wolfe	Gradient Descent
Iteration cost	(approx.) solve linearized problem on D	projection back to D
Iterates	sparse (in terms of used vertices)	dense 🗡

 \boldsymbol{x}

 \boldsymbol{x}

 $\mathcal{D}\subset\mathbb{R}^{d}$

Algorithm Variants

Predefined Stepsize

Algorithm 1 Frank-Wolfefor $k = 0 \dots K$ doCompute $s := \arg \min \langle s', \nabla f(x^{(k)}) \rangle$ Let $\gamma := \frac{2}{k+2}$ Update $x^{(k+1)} := (1 - \gamma)x^{(k)} + \gamma s$ end for

• Approximate Subproblems

[Dunn et al. 1978]

• Away-Steps

[Guélat et al. 1986]



Fully Corrective



What's new?

[J. 2011 PhD Thesis, summary: J. 2013 ICML Paper]

Primal-Dual Analysis

with **certificates** for approximation quality

- Approximate Subproblems and *inexact gradients* (and domains)
- Affine Invariance
- Lower Bound

and optimality in terms of sparsity

More Applications

Convergence Analysis

Primal Rate
$$f(\boldsymbol{x}^{(k)}) - f(\boldsymbol{x}^*) \leq \frac{2C_f}{k+2}$$

[Frank & Wolfe 1956, Dunn et al. 1978, 1980 : $O\left(rac{1}{k}
ight)$



efficient **certificates** for approximation quality

[Simplex domain: Clarkson 2008]

Here: For all algorithm variants, and affine invariant

Curvature constant $C_{f} := \sup_{\substack{\boldsymbol{x}, \boldsymbol{s} \in \mathcal{D}, \\ \gamma \in [0,1], \\ \boldsymbol{y} = \boldsymbol{x} + \gamma(\boldsymbol{s} - \boldsymbol{x})}} \frac{\frac{2}{\gamma^{2}} (f(\boldsymbol{y}) - f(\boldsymbol{x}) - \langle \boldsymbol{y} - \boldsymbol{x}, \nabla f(\boldsymbol{x}) \rangle)}{C_{f} \leq \operatorname{diam}_{\|.\|} (\mathcal{D})^{2} L}$

A Simple Duality Gap

Original Problem

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

The DualValue $\omega(\boldsymbol{x}) :=$ $\min_{\boldsymbol{s}' \in \mathcal{D}} f(\boldsymbol{x}) + \langle \boldsymbol{s}' - \boldsymbol{x}, \nabla f(\boldsymbol{x}) \rangle$

Weak Duality

 $\omega(\boldsymbol{x}) \leq \boldsymbol{f}(\boldsymbol{x}^*) \leq f(\boldsymbol{x}')$



Affine Invariance





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

Lower Bound



For an approximation **quality of 1/k**, any x must have **at least k non-zeros**. (use **at least k corners**)

Trade-Off: Approximation quality vs sparsity





Corollary: Obtain $O(\frac{1}{k})$ -approximate solution of sparsity k .

[Clarkson 2008]

Trade-Off: Approximation quality vs **sparsity** lower bound: $\Omega\!\left(\frac{1}{k}\right)$

 $f(\boldsymbol{x}) := \|\boldsymbol{x}\|_2^2$

Sparse Approximation

$$\mathcal{D} := \operatorname{conv}\left(\{\pm \mathbf{e}_i \mid i \in [n]\}\right)$$

Corollary: Obtain $O(\frac{1}{k})$ -approximate solution of sparsity k .

Trade-Off: Approximation quality vs **sparsity** lower bound: $\Omega\!\left(\frac{1}{k}\right)$

 $\min_{\|m{x}\|_1\leq 1}f(m{x})$

 ℓ_1 -ball

Greedy meets Frank-Wolfe



Low Rank Approximation

 $\min_{\|X\|_* \le 1} f(X)$

 $\mathcal{D} := \operatorname{conv}\left(\left\{\boldsymbol{u}\boldsymbol{v}^T \mid \boldsymbol{u} \in \mathbb{R}^n, \|\boldsymbol{u}\|_2 = 1 \\ \boldsymbol{v} \in \mathbb{R}^m, \|\boldsymbol{v}\|_2 = 1\right\}\right)$

trace-norm-ball

Corollary: Obtain $O\bigl(\frac{1}{k}\bigr)$ -approximate solution of rank k .

Trade-Off: Approximation quality vs rank

Projection:	FW-step:
Requires full SVD	approx. top singular vecto

lower bound: $\Omega\!\left(\frac{1}{k}\right)$

[J. & Sulovský 2010]



 $\min_{\|\boldsymbol{x}\|_p \leq 1} f(\boldsymbol{x})$

$$\mathcal{D} := \ \ell_p\text{-ball}$$



Examples of Atomic Domains Suitable for Frank-Wolfe

\mathcal{X}	Optimization Domain		Complexity of one Frank-Wolfe Iteration	
	Atoms \mathcal{A}	$\mathcal{D} = \operatorname{conv}(\mathcal{A})$	$\sup_{oldsymbol{s}\in\mathcal{D}}\langleoldsymbol{s},oldsymbol{y} angle$	Complexity
\mathbb{R}^n	Sparse Vectors	$\ .\ _1$ -ball	$\ oldsymbol{y} \ _{\infty}$	O(n)
\mathbb{R}^n	Sign-Vectors	$\ .\ _{\infty}$ -ball	$\ \boldsymbol{y} \ _1$	O(n)
\mathbb{R}^n	ℓ_p -Sphere	$\ .\ _p$ -ball	$\ \boldsymbol{y} \ _q$	O(n)
\mathbb{R}^n	Sparse Non-neg. Vectors	Simplex Δ_n	$\max_i \{ \boldsymbol{y}_i \}$	O(n)
\mathbb{R}^n	Latent Group Sparse Vec.	$\ .\ _{\mathcal{G}}$ -ball	$\left\ \max_{g\in\mathcal{G}}\left\ \boldsymbol{y}_{(g)}\right\ _{g}^{*}\right\ $	$\sum_{g \in \mathcal{G}} g $
$\boxed{\mathbb{R}^{m \times n}}$	Matrix Trace Norm	$\ .\ _{tr}$ -ball	$\ \boldsymbol{y}\ _{op} = \sigma_1(\boldsymbol{y})$	$\tilde{O}(N_f/\sqrt{\varepsilon'})$ (Lanczos)
$\mathbb{R}^{m imes n}$	Matrix Operator Norm	$\ .\ _{op}$ -ball	$\ \boldsymbol{y}\ _{tr} = \ (\sigma_i(\boldsymbol{y}))\ _1$	SVD
$\boxed{\mathbb{R}^{m \times n}}$	Schatten Matrix Norms	$\ (\sigma_i(.))\ _p$ -ball	$\ \ (\sigma_i(\boldsymbol{y})) \ _q$	SVD
$\boxed{\mathbb{R}^{m \times n}}$	Matrix Max-Norm	$\ .\ _{\max}$ -ball		$\tilde{O}(N_f(n+m)^{1.5}/\varepsilon'^{2.5})$
$\boxed{\mathbb{R}^{n \times n}}$	Permutation Matrices	Birkhoff polytope		$O(n^3)$
$\mathbb{R}^{n \times n}$	Rotation Matrices	8888888888888	88888888888888	SVD (Procrustes prob.)
$\mathbb{S}^{n \times n}$	Rank-1 PSD matrices of unit trace	$\{\boldsymbol{x} \succeq 0, \operatorname{Tr}(\boldsymbol{x}) = 1\}$	$\lambda_{ ext{max}}(oldsymbol{y})$	$\tilde{O}(N_f/\sqrt{\varepsilon'})$ (Lanczos)
$\mathbb{S}^{n \times n}$	PSD matrices of bounded diagonal	$\{ \boldsymbol{x} \succeq 0, \ \boldsymbol{x}_{ii} \leq 1 \}$		$\tilde{O}(N_f n^{1.5}/arepsilon'^{2.5})$

Table 1: Some examples of atomic domains suitable for optimization using the Frank-Wolfe algorithm. Here SVD refers to the complexity of computing a singular value decomposition, which is $O(\min\{mn^2, m^2n\})$. N_f is the number of non-zero entries in the gradient of the objective function f, and $\varepsilon' = \frac{2\delta C_f}{h}$ is the required accuracy for the knear subproblems. For any $\rho \in [1, \infty]$

, the conjugate value a is meant to satisfy $\frac{1}{2}$.

Factorized Matrix Domains

 $\mathcal{D} := \operatorname{conv}\left(\left\{ oldsymbol{u}oldsymbol{v}^T \; \left| oldsymbol{u} \in \mathbb{R}^n, \, \|oldsymbol{u}\|_2 = 1
ight\}
ight)$ (trace norm)



r	$\mathcal{A}_{\mathrm{left}} \subseteq \mathbb{R}^{m imes r}$	$\mathcal{A}_{\mathrm{right}} \subseteq \mathbb{R}^{n imes r}$	$\Omega_{\operatorname{conv}(\mathcal{A})}(M)$	$\Omega^*_{\mathcal{A}}(M)$	FW step
1	$\ .\ _2$ -sphere	$\ .\ _2$ -sphere	Trace norm $\ M\ _{tr}$	$\ M\ _{op}$	Lanczos, see Table 1
1	$\ .\ _1$ -sphere	$\ .\ _1$ -sphere	Vector ℓ_1 -norm $\ \vec{M}\ _1$	$\ ec{M}\ _{\infty}$	O(nm)
1	$\ .\ _{\infty}$ -sphere	$\ .\ _{\infty}$ -sphere		Cut-norm $\ .\ _{\infty \to 1}$	NP-hard (Alon & Naor, 2006)
n+m	$\ .\ _{2,\infty}$	$\ .\ _{2,\infty}$	Max-norm $\ M\ _{\max}$		SDP, see Table 1
1	$\left\ . \right\ _2 \cap \mathbb{R}^m_{\geq 0}$	$\left\ .\right\ _2 \cap \mathbb{R}^n_{\geq 0}$	"non-neg. trace norm"	366666666666	NP-hard (Murty & Kabadi, 1987)
1	Simplex Δ_m	Simplex Δ_n	"non-neg. matrix ℓ_1 -nor	m"	O(nm)

Table 2. Examples of some factorized matrix norms on $\mathbb{R}^{m \times n}$, each induced by two atomic norms (last two rows giving

Active Research

- Faster Convergence for Strongly Convex f under additional assumptions
 [Talk by Dan Garber, Guélat&Marcotte 1986, Beck&Teboulle 2004]
- Penalized Version [Talk by Zaid Harchaoui]

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) + \lambda \|\boldsymbol{x}\|_{\mathcal{A}}$$

- Non-Smooth f
 [Talks by Federico Pierucci and Guanghui Lan]
- Block-Wise Version
 [Talk by Simon Lacoste-Julien]
- Find More Applications!

$$\min_{oldsymbol{x}\in\mathcal{D}^{(1)} imes\dots imes\mathcal{D}^{(n)}}rac{f(oldsymbol{x})}{oldsymbol{x}=(oldsymbol{x}_{(1)},\dots,oldsymbol{x}_{(n)})}$$

