

# Revisiting Frank-Wolfe:

## Projection-Free Sparse Convex Optimization

Martin Jaggi

CMAP - Ecole Polytechnique,  
ERC SIPA project,  
CNRS UMR 7641, Paris, France



**Introduction:** There are two types of **first-order methods** for constrained convex optimization. One of them became nearly forgotten in the last decades.

**Contributions:** Stronger and more general primal-dual convergence results for Frank-Wolfe methods, and a unified view on many variants and applications.

### The Frank-Wolfe Algorithm

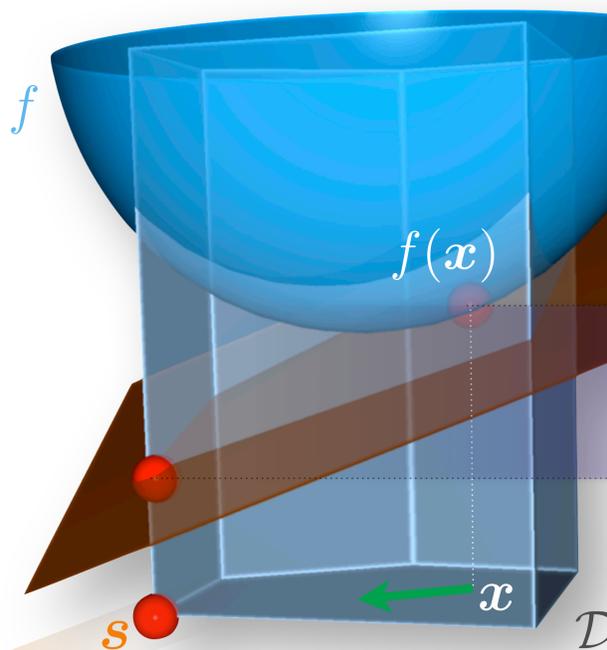
(or conditional gradient)

**Idea:** Minimize a linear approximation of  $f$

#### Algorithm 1 Frank-Wolfe

```

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



#### Setup

Constrained convex optim.

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

$\mathcal{D}$  compact,  $f$  differentiable

useful & efficient **certificate** for the approximation quality

$$g(\mathbf{x}) := \max_{\mathbf{s} \in \mathcal{D}} \langle \mathbf{x} - \mathbf{s}, \nabla f(\mathbf{x}) \rangle$$

#### Duality Gap

(for any bounded constrained problem)

### Convergence Results

#### Primal Rate

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

#### Primal-Dual Rate

$$g(\mathbf{x}^{(\hat{k})}) \leq \frac{7C_f}{k+2}$$

For some  $\hat{k} \leq k$

**For all algorithm variants!**

Curvature constant  $C_f$  is bounded by the Lipschitz constant of the gradient, times diameter:  
 $C_f \leq \text{diam}_{\|\cdot\|}(\mathcal{D})^2 L$  w.r.t. any choice of norm! (our algorithms/analysis are norm-free)

### Advantages

	Frank-Wolfe	Gradient Methods
Iterates	<b>Sparse</b> ✓ (using at most $k$ atoms after $k$ iterations)	<b>Dense</b> ✗
Iteration Cost	<b>Linear oracle</b> (can be much cheaper)	<b>Projection step</b>
Example:	top EV	full SVD

### Affine Invariance

The **algorithms** and our **analysis** are fully **invariant** under (affine) transformations of the input task

Contrasting gradient methods which use projections and norms

### Optimal Sparsity & Rate

The obtained **sparsity**  $k$  is **optimal** for an approximation **quality** of  $1/k$

No algorithm can do better in general  
Lower bounds for sparsity ( $l_1$ -domain) and low-rank (trace-norm domain)

### Algorithm Variants

- **Approximate subproblems** and **inexact gradients** (or inexact domain) (using approximate linear minimizers  $\mathcal{S}$ )
- **Line-search** for the optimal step-size  $\gamma \in [0, 1]$
- **Fully corrective** (re-optimizing over all used  $\mathcal{S}$ )
- **Away steps** (removing the worst of all used  $\mathcal{S}$ )

### Factorized Matrix Domains

For two sets  $\mathcal{A}_{\text{left}} \subseteq \mathbb{R}^{m \times r}$  and  $\mathcal{A}_{\text{right}} \subseteq \mathbb{R}^{n \times r}$  consider the **outer-product** matrices

$$\mathcal{A} := \left\{ LR^T \mid L \in \mathcal{A}_{\text{left}}, R \in \mathcal{A}_{\text{right}} \right\}$$

$\mathcal{D} := \text{conv}(\mathcal{A})$

→ every FW iteration is a **low-rank** update

Natural way to optimize over **matrix factorizations** (including **sparse** and **non-negative** ones)

$r$	$\mathcal{A}_{\text{left}} \subseteq \mathbb{R}^{m \times r}$	$\mathcal{A}_{\text{right}} \subseteq \mathbb{R}^{n \times r}$	$\Omega_{\text{conv}(\mathcal{A})}(M)$	$\Omega_{\mathcal{A}}(M)$	FW step
1	$\ \cdot\ _2$ -sphere	$\ \cdot\ _2$ -sphere	Trace norm $\ M\ _{\text{tr}}$	$\ M\ _{\text{op}}$	Lanczos, see Table 1
1	$\ \cdot\ _1$ -sphere	$\ \cdot\ _1$ -sphere	Vector $\ell_1$ -norm $\ M\ _1$	$\ M\ _{\infty}$	$O(nm)$
1	$\ \cdot\ _{\infty}$ -sphere	$\ \cdot\ _{\infty}$ -sphere	Cut-norm $\ M\ _{\infty \rightarrow 1}$		NP-hard
$n+m$	$\ \cdot\ _{2,\infty}$	$\ \cdot\ _{2,\infty}$	Max-norm $\ M\ _{\text{max}}$		SDP, see Table 1
1	$\ \cdot\ _2 \cap \mathbb{R}_{\geq 0}^m$	$\ \cdot\ _2 \cap \mathbb{R}_{\geq 0}^n$	"non-neg. trace norm"		NP-hard [MKS7]
1	Simplex $\Delta_m$	Simplex $\Delta_n$	"non-neg. matrix $\ell_1$ -norm"		$O(nm)$

### Applications to sparse and low-rank optimization

**Generalized sparsity problems:** Usually in machine learning and signal processing applications, we optimize over a domain  $\mathcal{D} := \text{conv}(\mathcal{A})$ , that is the convex hull of a simple set of things/atoms (**atomic norm** idea).

For such problems, Frank-Wolfe methods are **particularly suitable**:

$\mathcal{X}$	Optimization Domain Atoms $\mathcal{A}$	$\mathcal{D} = \text{conv}(\mathcal{A})$	Complexity of one Frank-Wolfe Iteration $\Omega_{\mathcal{D}}^*(\mathbf{y}) = \sup_{\mathbf{s} \in \mathcal{D}} \langle \mathbf{s}, \mathbf{y} \rangle$ Complexity
$\mathbb{R}^n$	Sparse vectors	$\ \cdot\ _1$ -ball	$\ \mathbf{y}\ _{\infty}$ $O(n)$
$\mathbb{R}^n$	Sign-vectors	$\ \cdot\ _{\infty}$ -ball	$\ \mathbf{y}\ _1$ $O(n)$
$\mathbb{R}^n$	$\ell_p$ -Sphere	$\ \cdot\ _p$ -ball	$\ \mathbf{y}\ _q$ $O(n)$
$\mathbb{R}^n$	Sparse non-neg. vectors	Simplex $\Delta_n$	$\max_i \{y_i\}$ $O(n)$
$\mathbb{R}^n$	Latent group sparse vect.	$\ \cdot\ _G$ -ball	$\max_{g \in \mathcal{G}} \ \mathbf{y}(g)\ _q$ $\sum_{g \in \mathcal{G}}  g $
$\mathbb{R}^{m \times n}$	Matrix trace norm	$\ \cdot\ _{\text{tr}}$ -ball	$\ \mathbf{y}\ _{\text{op}} = \sigma_1(\mathbf{y})$ $\tilde{O}(N_f / \sqrt{\varepsilon'})$ (Lanczos)
$\mathbb{R}^{m \times n}$	Matrix operator norm	$\ \cdot\ _{\text{op}}$ -ball	$\ \mathbf{y}\ _{\text{tr}} = \ \sigma_i(\mathbf{y})\ _1$ SVD
$\mathbb{R}^{m \times n}$	Schatten matrix norms	$\ (\sigma_i(\cdot))\ _p$ -ball	$\ (\sigma_i(\mathbf{y}))\ _q$ SVD
$\mathbb{R}^{m \times n}$	Matrix max-norm	$\ \cdot\ _{\text{max}}$ -ball	$\tilde{O}(N_f(n+m)^{1.5}/\varepsilon'^{2.5})$
$\mathbb{R}^{n \times n}$	Permutation matrices	Birkhoff polytope	$O(n^3)$
$\mathbb{R}^{n \times n}$	Rotation matrices		SVD (Procrustes prob.)
$\mathbb{S}^{n \times n}$	Rank-1 PSD matrices of unit trace	$\{\mathbf{x} \succeq 0, \text{Tr}(\mathbf{x}) = 1\}$	$\lambda_{\text{max}}(\mathbf{y})$ $\tilde{O}(N_f / \sqrt{\varepsilon'})$ (Lanczos)
$\mathbb{S}^{n \times n}$	PSD matrices of bounded diagonal	$\{\mathbf{x} \succeq 0, x_{ii} \leq 1\}$	$\tilde{O}(N_f n^{1.5}/\varepsilon'^{2.5})$

**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  $O(\min\{m^2, n^2\})$ .  $N_f$  is the number of non-zero entries in the gradient of the objective function  $f$ , and  $\varepsilon' = \frac{\delta C_f}{k+2}$  is the required accuracy for the linear subproblems. For any  $p \in [1, \infty]$ , the conjugate value  $q$  is meant to satisfy  $\frac{1}{p} + \frac{1}{q} = 1$ , allowing  $q = \infty$  for  $p = 1$  and vice versa.

Many other application such as optimizing over **structured atomic norms**, **matrix factorizations**, and **submodular** optimization