

# Optimization for Machine Learning

## CS-439

Lecture 8: Frank-Wolfe algorithm, and Muon

**Martin Jaggi**

EPFL – [github.com/epfml/OptML\\_course](https://github.com/epfml/OptML_course)

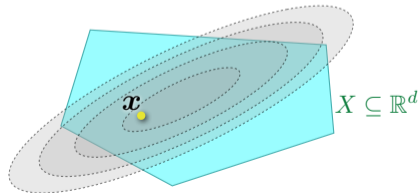
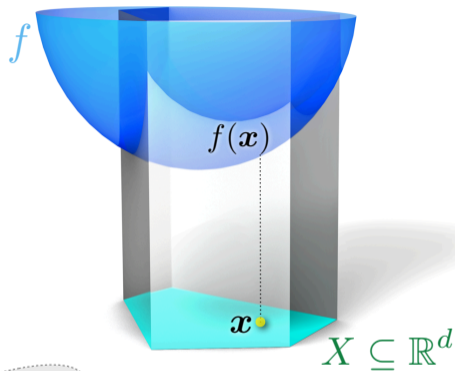
April 24, 2026

# Frank-Wolfe

# Constrained Optimization

## Constrained Optimization Problem

$$\begin{array}{ll} \text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{X} \end{array}$$



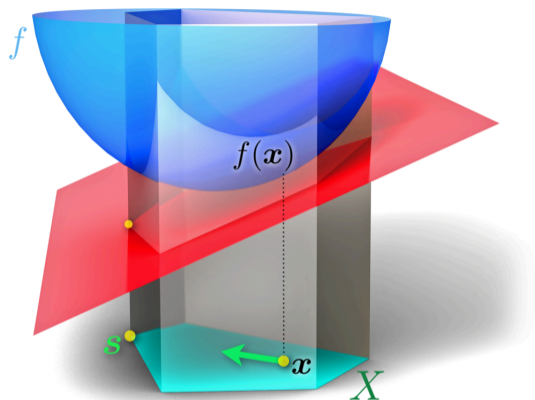
# Frank-Wolfe Algorithm

Frank-Wolfe Algorithm:

$$\mathbf{s} := \text{LMO}(\nabla f(\mathbf{x}_t)),$$

$$\mathbf{x}_{t+1} := (1 - \gamma)\mathbf{x}_t + \gamma\mathbf{s},$$

for timesteps  $t = 0, 1, \dots$ , and  
stepsize  $\gamma := \frac{2}{t+2}$ .



**Linear Minimization Oracle:**

$$\text{LMO}(\mathbf{g}) := \underset{\mathbf{s} \in \mathcal{X}}{\text{argmin}} \langle \mathbf{s}, \mathbf{g} \rangle$$

# Properties

- ▶ **Always feasible:**  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_t \in \mathcal{X}$ .  
 $\mathbf{x}_{t+1}$  is on line segment  $[\mathbf{s}, \mathbf{x}_t]$ , for  $\gamma \in [0, 1]$ .
- ▶ **Reduces** non-linear to linear optimization
- ▶ **Projection-free**
- ▶ **Sparse iterates** (in terms of corners  $\mathbf{s}$  used)

Invented and analyzed 1956 by Marguerite Frank and Philip Wolfe.

# Example

## Lasso Regression

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|^2 \quad s.t. \quad \|\mathbf{x}\|_1 \leq 1$$

L1-ball is the convex hull of the unit basis vectors:

$$\mathcal{X} = \{\mathbf{x} \mid \|\mathbf{x}\|_1 \leq 1\} = \text{conv}(\{\pm \mathbf{e}_1, \dots, \pm \mathbf{e}_n\}).$$

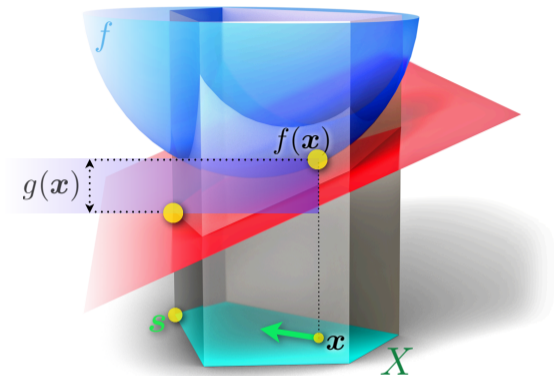
- ▶  $\nabla f(\mathbf{x}) = \mathbf{g} := \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b})$
- ▶  $\text{LMO}(\mathbf{g}) = -\text{sign}(g_i) \mathbf{e}_i$  with  $i := \underset{i \in [n]}{\text{argmax}} |g_i|$

simpler than projection onto L1-ball !

# Duality Gap

## Duality Gap

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



Certificate for optimization quality:

$$\begin{aligned} g(\mathbf{x}) &= \max_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{x} - \mathbf{s}, \nabla f(\mathbf{x}) \rangle \\ &\geq \langle \mathbf{x} - \mathbf{x}^*, \nabla f(\mathbf{x}) \rangle \\ &\geq f(\mathbf{x}) - f(\mathbf{x}^*) \end{aligned}$$

## Stepsize variants

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

$$\gamma_t := \operatorname{argmin}_{\gamma \in [0,1]} f((1-\gamma)\mathbf{x}_t + \gamma\mathbf{s}), \quad (\text{line-search})$$

$$\gamma_t := \min \left\{ \frac{g(\mathbf{x}_t)}{L \|\mathbf{s} - \mathbf{x}_t\|^2}, 1 \right\}, \quad (\text{gap-based})$$

## Convergence in $\mathcal{O}(1/\varepsilon)$ steps

### Theorem

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and *smooth* with parameter  $L$ , and  $\mathbf{x}_0 \in \mathcal{X}$ . Then choosing any of the above stepsizes, the Frank-Wolfe algorithm yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{2L \operatorname{diam}(\mathcal{X})^2}{T + 1}$$

Where  $\operatorname{diam}(\mathcal{X}) := \max_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \|\mathbf{x} - \mathbf{y}\|$  is the diameter of  $\mathcal{X}$ .

## Proof of Convergence in $\mathcal{O}(1/\varepsilon)$ steps

### Lemma

For a step  $\mathbf{x}_{t+1} := \mathbf{x}_t + \gamma(\mathbf{s} - \mathbf{x}_t)$  with arbitrary step-size  $\gamma \in [0, 1]$ , it holds that

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma g(\mathbf{x}_t) + \frac{\gamma^2}{2} L \operatorname{diam}(\mathcal{X})^2 ,$$

if  $\mathbf{s} = \operatorname{LMO}(\nabla f(\mathbf{x}_t))$ .

### Proof.

We write  $\mathbf{x} := \mathbf{x}_t$ ,  $\mathbf{y} := \mathbf{x}_{t+1} = \mathbf{x} + \gamma(\mathbf{s} - \mathbf{x})$ . From the definition of smoothness of  $f$ , we have

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x} + \gamma(\mathbf{s} - \mathbf{x})) \\ &\leq f(\mathbf{x}) + \gamma \langle \mathbf{s} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{\gamma^2}{2} L \operatorname{diam}(\mathcal{X})^2 . \end{aligned}$$

The lemma follows by definition of the duality gap. □

## Proof of Convergence in $\mathcal{O}(1/\varepsilon)$ steps

From the Lemma we know that for every step of FW, it holds that

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma g(\mathbf{x}_t) + \gamma^2 C,$$

if we chose  $\gamma := \frac{2}{t+2}$  and write  $C := \frac{1}{2}L \text{diam}(\mathcal{X})^2$ . This bound holds also for all mentioned line-search variants (*different LHS, same RHS*).

Writing  $h(\mathbf{x}) := f(\mathbf{x}) - f(\mathbf{x}^*)$  for the (unknown) objective error at any point  $\mathbf{x}$ , this implies that

$$\begin{aligned} h(\mathbf{x}_{t+1}) &\leq h(\mathbf{x}_t) - \gamma g(\mathbf{x}_t) + \gamma^2 C \\ &\leq h(\mathbf{x}_t) - \gamma h(\mathbf{x}_t) + \gamma^2 C \\ &= (1 - \gamma)h(\mathbf{x}_t) + \gamma^2 C, \end{aligned}$$

by the certificate property  $h(\mathbf{x}) \leq g(\mathbf{x})$  of the duality gap.  
The theorem then follows by induction (Lab 8, Exercise 1). □

# Affine Invariance

## Curvature Constant

$$C_f := \sup_{\substack{\mathbf{x}, \mathbf{s} \in \mathcal{X}, \gamma \in [0,1] \\ \mathbf{y} = \mathbf{x} + \gamma(\mathbf{s} - \mathbf{x})}} \frac{1}{\gamma^2} (f(\mathbf{y}) - f(\mathbf{x}) - \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle)$$

Algorithm is invariant to scaling (affine transformations) of the input problem.

So is  $C_f$ .

(same as Newton, but here for **constrained** problems)

$$C_f \leq \frac{L}{2} \text{diam}(\mathcal{X})^2 \quad \text{for any norm!}$$

All proofs hold for  $C_f$ , instead of picking a particular  $L \text{diam}(\mathcal{X})^2$ .

# Convergence in Duality Gap

## Theorem

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\mathcal{X}$  be convex with  $C_{(f,\mathcal{X})} < \infty$ , and  $\mathbf{x}_0 \in \mathcal{X}$ ,  $T \geq 2$ . Then choosing any of the above stepsizes, the Frank-Wolfe algorithm yields a  $t$ ,  $1 \leq t \leq T$  s.t.

$$g(\mathbf{x}_t) \leq \frac{27/2 C_f}{T + 1}$$

## Proof.

Idea: not all gaps can be small (use Lemma again). □

# Extensions and Use Cases

## Extensions:

- ▶ **Approximate** LMO (of additive or multiplicative accuracy)
- ▶ LMO applied to stochastic gradient → **stochastic FW** (only works with momentum)
- ▶ **Randomized** LMO
- ▶ unconstrained problems (Matching Pursuit variants)

## Use cases:

Whenever projection is more costly than solving a linear problem

- ▶ **Lasso** and other L1-constrained problems
- ▶ **Matrix Completion**: scalable algorithm
- ▶ Relaxation of **combinatorial problems** (e.g. matchings, network flows etc)
- ▶ Muon optimizer (highly performant for training deep neural networks)

recall:  $\text{LMO}(\mathbf{g}) := \underset{\mathbf{s} \in \mathcal{X}}{\text{argmin}} \langle \mathbf{s}, \mathbf{g} \rangle$

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

Examples	$\mathcal{A}$	$ \mathcal{A} $	$d$	LMO ( $\mathbf{g}$ )
L1-ball	$\{\pm \mathbf{e}_i\}$	$2d$	$d$	$\pm \mathbf{e}_i$ with $\text{argmax}_i  g_i $ (attains the $\ \mathbf{g}\ _\infty$ norm)
Simplex	$\{\mathbf{e}_i\}$	$d$	$d$	$\mathbf{e}_i$ with $\text{argmin}_i g_i$
Norms	$\{\mathbf{x}, \ \mathbf{x}\  \leq 1\}$	$\infty$	$d$	$\text{argmin}_{\mathbf{s}, \ \mathbf{s}\  \leq 1} \langle \mathbf{s}, \mathbf{g} \rangle$
Nuclear norm	$\{\mathbf{X}, \ \mathbf{X}\ _* \leq 1\}$	$\infty$	$d^2$	(attains the Operator norm $\ \mathbf{g}\ _{\text{op}}$ )
Operator norm	$\{\mathbf{X}, \ \mathbf{X}\ _{\text{op}} \leq 1\}$	$\infty$	$d^2$	(attains the Nuclear norm $\ \mathbf{g}\ _*$ )
Wavelets	..	$\infty$	$\infty$	..

# Muon

# Muon: Momentum Um Orthogonalized by Newton-schulz

**Algorithm.** At each step  $t$ , with stochastic gradient  $\mathbf{g}_t$ :

$$\begin{aligned}\mathbf{M}_t &:= \mu \mathbf{M}_{t-1} + (1-\mu) \mathbf{g}_t && \text{(momentum buffer)} \\ \mathbf{X}_{t+1} &:= \mathbf{X}_t - \eta \text{orth}(\mathbf{M}_t) - \eta \lambda \mathbf{X}_t && \text{(update, with optional weight decay } \lambda \text{)}\end{aligned}$$

**Orthogonalization.** (Matrix-sign operation) If  $\mathbf{M}_t = \mathbf{U}\Sigma\mathbf{V}^\top$  is the SVD, then

$$\text{orth}(\mathbf{M}_t) = \mathbf{U}\mathbf{V}^\top.$$

All singular values are replaced by 1. Computed efficiently via [Newton–Schulz iteration](#) (5 steps, matmuls only — GPU-friendly, no SVD).

**Scope in neural network training.** Applied only to 2D weight matrices. Embeddings, biases&gains, and routers (in case of MoE) use AdamW. Stepsize  $\eta \approx 0.02$ . [JJB<sup>+</sup>24].

## Why orthogonalize? Balancing per-layer update magnitudes

**A layer's job:** map input  $\mathbf{a} \in \mathbb{R}^n$  to output activations  $\mathbf{b} = \mathbf{X}\mathbf{a} \in \mathbb{R}^m$ .

**What matters for stability:** How much an update  $\Delta\mathbf{X}$  changes the output for a typical input:

$$\|\Delta\mathbf{b}\| = \|\Delta\mathbf{X} \cdot \mathbf{a}\| \leq \|\Delta\mathbf{X}\|_{\text{op}} \cdot \|\mathbf{a}\|.$$

The **operator norm** (spectral norm)  $\|\Delta\mathbf{X}\|_{\text{op}} = \sigma_{\max}(\Delta\mathbf{X})$  is the right quantity — it controls the worst-case output change per unit input.

**Problem with raw gradients.** A (stochastic or deterministic) gradient  $\mathbf{M}_t$  typically has singular values of very different magnitudes. A classical gradient step thus changes outputs of a layer differently than its input. Imbalance can amplify through depth.

**Muon's key change.** Replacing  $\mathbf{M}_t$  with  $\mathbf{UV}^\top$  sets **every** singular value to 1: the update has spectral norm exactly 1 and acts uniformly across all directions. Each layer update contributes a controlled, bounded change to its output activations.

# Muon as Frank–Wolfe

**Frank–Wolfe** over a convex set  $\mathcal{X}$ :

$$\mathbf{S}_t = \text{LMO}(\mathbf{M}_t) := \arg \min_{\mathbf{S} \in \mathcal{X}} \langle \mathbf{S}, \mathbf{M}_t \rangle, \quad \mathbf{X}_{t+1} = (1 - \gamma)\mathbf{X}_t + \gamma \mathbf{S}_t.$$

**Muon update** (rearranged) is identical to FW with  $\gamma = \eta\lambda$ :

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \mathbf{U}\mathbf{V}^\top - \eta\lambda\mathbf{X}_t = (1 - \eta\lambda)\mathbf{X}_t + \eta\lambda \left( -\frac{1}{\lambda} \mathbf{U}\mathbf{V}^\top \right).$$

**Key identity.** Orthogonalization is the LMO over the operator-norm unit ball. For  $\mathbf{M}_t = \mathbf{U}\Sigma\mathbf{V}^\top$ :

$$-\mathbf{U}\mathbf{V}^\top = \arg \min_{\|\mathbf{S}\|_{\text{op}} \leq 1} \langle \mathbf{S}, \mathbf{M}_t \rangle$$




The extreme points of  $\{S : \|S\|_{\text{op}} \leq 1\}$  are the (semi-)orthogonal matrices.

**Corollary:** **Muon** with weight decay  $\lambda$  is identical to **Stochastic Frank–Wolfe** on

$$\min_{\mathbf{X}} F(\mathbf{X}) \quad \text{s.t.} \quad \|\mathbf{X}\|_{\text{op}} \leq 1/\lambda.$$

Weight decay sets constraint radius; orthogonalization solves the LMO [PXA<sup>+</sup>25, SW25]

# Bibliography

-  Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy Bernstein.  
Muon: An optimizer for hidden layers in neural networks.  
*blogpost* <https://kellerjordan.github.io/posts/muon/>, 2024.
-  Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and Volkan Cevher.  
Training deep learning models with norm-constrained lmos.  
*arXiv preprint* [arXiv:2502.07529](https://arxiv.org/abs/2502.07529), 2025.
-  Maria-Eleni Sfyraiki and Jun-Kun Wang.  
Lions and Muons: Optimization via Stochastic Frank-Wolfe.  
*arXiv preprint* [arXiv:2506.04192](https://arxiv.org/abs/2506.04192), 2025.