

Lecture #14 The Frank-Wolfe method

E1 260

(Conditional gradient method)

- Algorithm
 - Geometric understanding
 - Examples
- Convergence analysis
 - Smooth convex functions

The FW algorithm:

$$\begin{array}{l} \text{minimize} \quad f(\underline{x}) \\ \text{s.t.} \quad \underline{x} \in C \end{array}$$

- f is differentiable, L -smooth convex function
- $C \subseteq \text{Dom}(f)$ is convex and closed set

FW has two steps:

① Direction finding:

Solves a linear optimization over a convex set

$$\underline{s}_t = \arg \min_{\underline{s} \in C} \nabla f^T(\underline{x}_t) \underline{s}$$

local linear approximation

$$f(\underline{x}) \approx f(\underline{x}_t) + \nabla f^T(\underline{x}_t) (\underline{s} - \underline{x}_t)$$

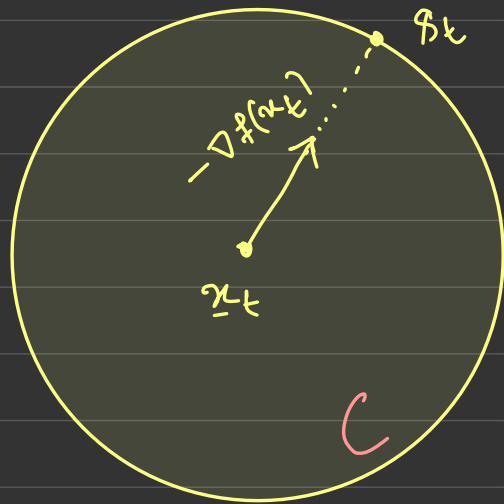
② update:

$$\underline{x}_{t+1} = (1 - \gamma_t) \underline{x}_t + \gamma_t \underline{s}_t \in C$$

③ Step size : $\alpha_t = \frac{2}{t+1}$

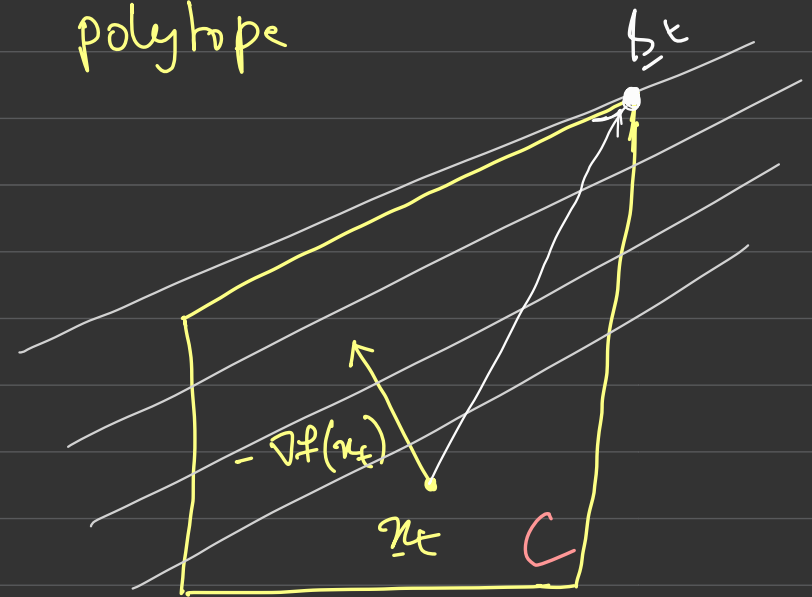
Examples:

① C is a sphere centered at origin



$$s_t \in \arg \min_{A \in C} \|\nabla f(x_t)\|$$

② polytope

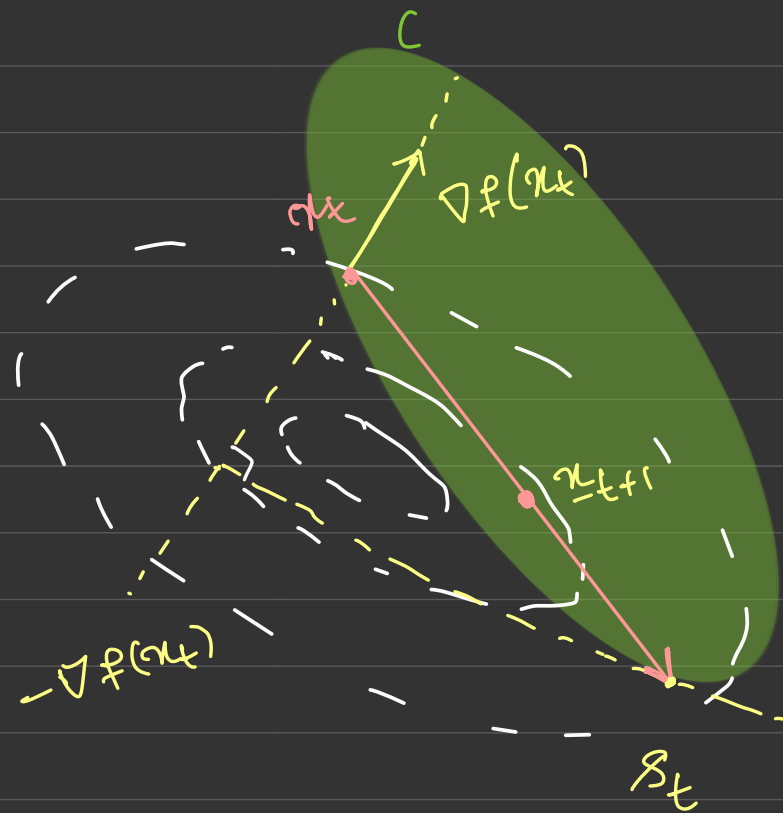


Linear program

In this example, FW follows GD trajectory

FW does n't always follow -ve gradient

Another example:



Projected Gradient descent vs. FW

- ① Projection is replaced by linear optimization
- ② Both require gradient computation

Norm constraints:

$C = \{ \underline{x} : \|\underline{x}\| \leq k \}$ for an arbitrary norm $\|\cdot\|$

$$\underline{s}_t \in \arg \min_{\|\underline{s}\| \leq k} \nabla f^T(\underline{x}_t) \underline{s}$$

$$= -k \arg \max_{\|\underline{s}\| \leq k} \nabla f^T(\underline{x}_t) \underline{s}$$

$$= -k \underbrace{\partial \|\nabla f(\underline{x}_t)\|_*}_{\text{Subgradient of its dual norm}}$$

Subgradient of its dual norm

l_1 - norm:

$$\min f(x)$$

$$\text{s.t. } \|x\|_1 \leq k$$

$$s_t \in -k \partial \|\nabla f(x_t)\|_\infty$$

$$i_t \in \arg \max_{i=1, \dots, d} |\nabla_i f(x_t)|$$

FW update:

$$x_{t+1} = (1 - \gamma_t) x_t - \gamma_t \cdot k \operatorname{sgn}(\nabla_{i_t} f(x_t)) e_{i_t}$$

- Greedy co-ordinate descent (Simpler than projection onto l_1 - ball)

Can be applied for non-convex problems:

minimize $-\underline{x}^T Q \underline{x}$ s.t. $\|\underline{x}\|_2 \leq 1$
with $Q > 0$

$$\underline{\delta}_t \leftarrow - \frac{\nabla f(\underline{x}_t)}{\|\nabla f(\underline{x}_t)\|_2}$$

$$= - \frac{\nabla f(\underline{x}_t)}{\|\nabla f(\underline{x}_t)\|_2} = \frac{Q \underline{x}_t}{\|Q \underline{x}_t\|_2}$$

$$\Rightarrow \underline{x}_{t+1} = (1 - \alpha_t) \underline{x}_t + \alpha_t \frac{Q \underline{x}_t}{\|Q \underline{x}_t\|_2}$$

Set $\alpha_t = 1$ $\left[= \arg \min_{0 \leq \alpha \leq 1} f((1-\alpha)\underline{x}_t + \alpha \frac{Q \underline{x}_t}{\|Q \underline{x}_t\|_2}) \right]$

$\Rightarrow \underline{x}_{t+1} = \frac{Q \underline{x}_t}{\|Q \underline{x}_t\|_2}$: power method to find the leading eigen vector of Q

Convergence result:

Let $f: \text{Dom}(f) \rightarrow \mathbb{R}$ be convex and L -smooth and

$$D = \text{diameter}(C) = \sup_{x, y \in C} \|x - y\|. \quad \text{with } \gamma_t = \frac{2}{t+1},$$

FW satisfies

$$f(x_T) - f(x^*) \leq \frac{2LD^2}{T+1}$$

Sublinear convergence : $O\left(\frac{1}{T}\right)$ [same as projected gradient descent]
 ϵ -accuracy : $O\left(\frac{1}{\epsilon}\right)$

Proof: $f(\underline{y}) \leq f(\underline{x}) + \nabla f^\top(\underline{x})(\underline{y} - \underline{x}) + \frac{L}{2} \|\underline{x} - \underline{y}\|^2$

$\underline{y} = \underline{x}_{t+1} ; \underline{x} = \underline{x}_t$

$$f(\underline{x}_{t+1}) - f(\underline{x}_t) \leq \nabla f^\top(\underline{x}_t)(\underline{x}_{t+1} - \underline{x}_t) + \frac{L}{2} \|\underline{x}_{t+1} - \underline{x}_t\|^2$$

$\underline{x}_{t+1} = \underline{x}_t + \gamma_t (\underline{s}_t - \underline{x}_t)$

$$f(\underline{x}_{t+1}) - f(\underline{x}_t) \leq \underbrace{\gamma_t \nabla f^\top(\underline{x}_t)(\underline{s}_t - \underline{x}_t)}_{\leq \nabla f^\top(\underline{x}_t) \underline{x}^*} + \frac{L}{2} \gamma_t^2 \underbrace{\|\underline{s}_t - \underline{x}_t\|^2}_{\leq \frac{L}{2} \gamma_t^2 \mathcal{D}^2}$$

$\underline{s}_t \in \arg \min_{\underline{s} \in \mathcal{C}} \nabla f^\top(\underline{x}_t) \underline{s}$

$$f(\underline{x}_{t+1}) - f(\underline{x}_t) \leq \underbrace{\gamma_t \nabla f^\top(\underline{x}_t) [\underline{x}^* - \underline{x}_t]}_{\text{convexity}} + \frac{L}{2} \gamma_t^2 \mathcal{D}^2$$

$$\leq \gamma_t [f(\underline{x}^*) - f(\underline{x}_t)] + \frac{L}{2} \gamma_t^2 \mathcal{D}^2$$

$$f(\underline{x}_{t+1}) - f(\underline{x}^*) \leq (1 - \gamma_t) [f(\underline{x}_t) - f(\underline{x}^*)] + \frac{L}{2} \gamma_t^2 \mathcal{D}^2$$

$$\Delta_{t+1} \leq (1 - \gamma_t) \Delta_t + \frac{L}{2} \gamma_t^2 D^2$$

Our claim:

$$\Delta_t \leq \frac{2L D^2}{T+1}$$

$$\gamma_t = \frac{2}{t+1}$$

Proof by induction:

① Base case $t = 1$: $\Delta_2 \leq 0 + \frac{L}{2} D^2 \leq \frac{2L D^2}{3}$ [$\gamma_1 = 1$]

② Inductive hypothesis: assume the upper bound is true for all τ

③ Need to show it holds for t

$$\textcircled{1} \quad \Delta_{t+1} \leq (1 - \gamma_t) \Delta_t + \frac{L}{2} \gamma_t^2 D^2$$

$$\leq \left(1 - \frac{2}{t+1}\right) \frac{2L D^2}{t+1} + \frac{L}{2} \cdot \frac{4}{(t+1)^2} D^2$$

$$= \frac{t-1}{t+1} \cdot \frac{2L D^2}{t+1} + \frac{2L D^2}{(t+1)^2}$$

$$= \frac{2LD^2}{(t+1)^2} [t-1+1]$$

$$= \frac{2LD^2}{(t+1)} \cdot \frac{t}{(t+1)} \leq \frac{2LD^2}{(t+1)}$$



• FW updates are affine invariant:

Suppose $\underline{x} = A\underline{x}'$ and $F(\underline{x}') = \mathbb{F}(A\underline{x}')$

for non singular A . Then

$$\underline{s}' = \arg \min_{z \in A^{-1}C} \nabla^T F(\underline{x}')$$

$$(\underline{x}')^+ = (1-\alpha)\underline{x}^+ + \alpha \underline{s}'$$

multiplying by A produces same updates as that from \mathbb{F} .

• In general, stronger convexity does not improve convergence of FW.

- Additional conditions on the constraint set.

- μ -strongly convex set yield linear convergence.

Duality gap: (not covered; only for reference)

Constrained problem: minimize $f(x) + \underbrace{\mathbb{I}_C(x)}$

Indicator function

$$\mathbb{I}_C(x) = \begin{cases} 0, & \text{if } x \in C \\ +\infty, & \text{otherwise} \end{cases}$$

Recall: $f^*(y) = \sup_x \{ \underline{x}^T \underline{y} - f(x) \}$

$$\text{minimize}_x \text{ maximize}_u [\underline{x}^T \underline{u} - f^*(u)] + \mathbb{I}_C(x)$$

$$= \text{maximize}_u \text{ minimize}_x [\mathbb{I}_C(x) + \underline{x}^T \underline{u}] - f^*(u)$$

$$= \text{maximize}_u \left. \begin{aligned} & - \mathbb{I}_C^*(-u) - f^*(u) \end{aligned} \right\} \text{dual problem}$$

Duality gap between \underline{x} and \underline{u} :

$$f(\underline{x}) + f^*(\underline{u}) + \mathbb{I}_C^*(-\underline{u}) \geq \underline{x}^\top \underline{u} + \mathbb{I}_C^*(-\underline{u})$$

At $\underline{x} = \underline{x}_k$ and $\underline{u} = \nabla f(\underline{x}_k)$

$$\nabla f^\top(\underline{x}_k) \underline{x}_k + \max_{\underline{s} \in C} -\nabla f^\top(\underline{x}_k) \underline{s} = \nabla f^\top(\underline{x}_k) (\underline{x}_k - \underline{s}_k)$$

In fact,

$$f(\underline{x}_k) - f(\underline{x}^*) \geq \nabla f^\top(\underline{x}_k) (\underline{x}_k - \underline{s}_k)$$