

Frank-Wolfe algorithm - iterative first order optimization algorithm for constrained convex optimization.

M. Frank , P. Wolfe [1956]

a.k.a. Conditional gradients method,  
reduced gradient method,  
convex combination algorithm.

$$\min f(x)$$

s.t.  $x \in D$

$D$ : compact convex set

$f: D \rightarrow \mathbb{R}$  convex differentiable real-valued function

Algorithm: let  $k \leftarrow 0$ , and  $x_0$  be any point in  $D$

Step 1: (Direction finding subproblem): Find  $s_k$  solving

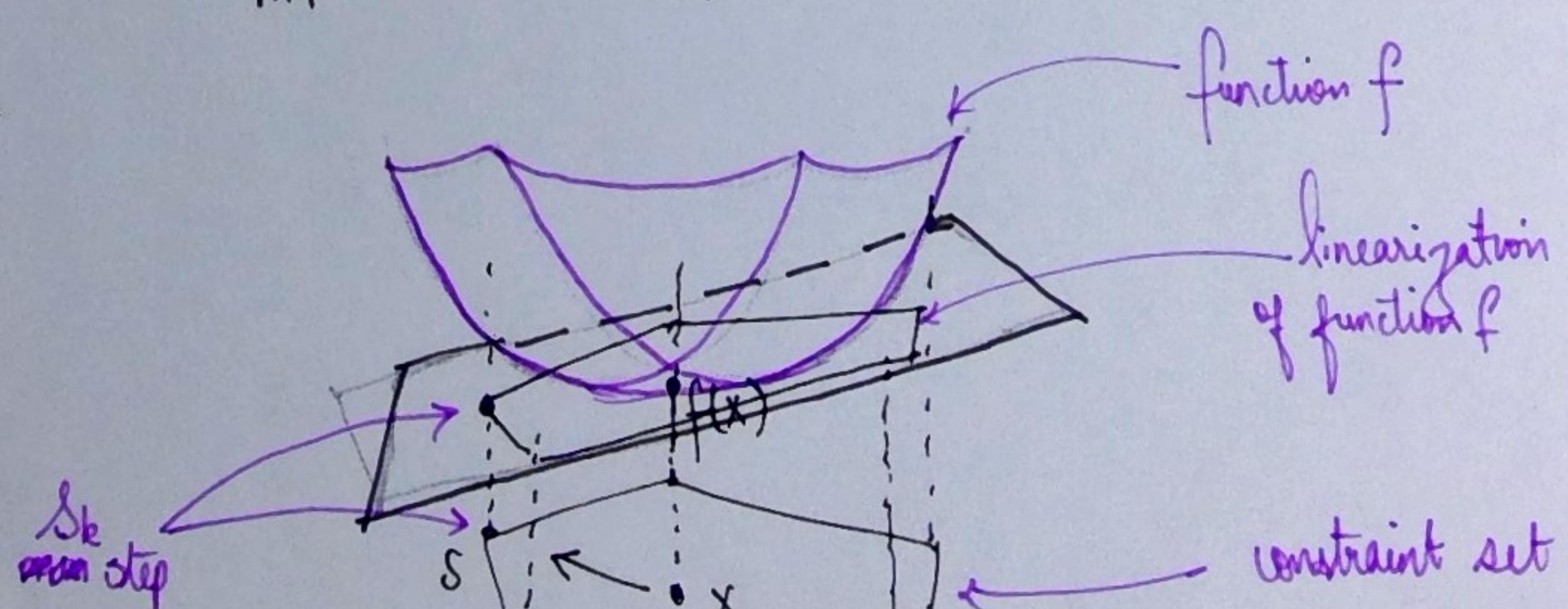
$$\min_{s \in D} s^T \nabla f(x_k)$$

Step 2: Determine step size  $\gamma_k$ :  $(\gamma_k \leftarrow \frac{2}{k+1}; \text{fixed step size fixed ahead})$   
 $\gamma_k$ :  $\gamma$  minimizes  $f(x_k + \gamma(s_k - x_k))$   
 st.  $0 \leq \gamma \leq 1$   
 (line search)

Step 3: Update:  $x_{k+1} \leftarrow x_k + \gamma_k(s_k - x_k)$ ;  $k \leftarrow k+1$

Iterate.

What's happening?



- No projection step!
- convergence is sublinear in general. Error is  $O(1/k)$  after  $k$  iterations as long as the gradient is Lipschitz continuous.
- Iterates are sparse combinations of extremal points of the feasible set. When the constraints are linear, the <sup>sub</sup>problem becomes an LP! (linear program)

What happens when  $C = \{x : \|x\| \leq t\}$  for a norm  $\|\cdot\|$ ?

$$\begin{aligned} \text{Then } s &\in \arg\min_{\|s\| \leq t} \nabla f(x_{k+1})^T s \\ &= -t \left( \arg\max_{\|s\| \leq 1} \nabla f(x_{k+1})^T s \right) \\ &= -t \circ \|\nabla f(x_{k+1})\|_* \quad \|\cdot\|_* \text{ is the corresponding dual norm.} \end{aligned}$$

### Convergence Analysis

$$\begin{aligned} M &= \max_{\substack{y, x, s \in C \\ y = (1-\gamma)x + \gamma s}} \frac{2}{\gamma^2} (f(y) - f(x) - \nabla f(x)^T (y - x)) \\ (\text{Bounded curvature}) \quad M &< \infty \end{aligned}$$

Note:  $\gamma \in [0, 1]$ . What happens when  $f(\cdot)$  is linear? ( $M=0$ )

$f(y) - f(x) - \nabla f(x)^T (y - x)$  is called Bregman divergence defined by  $f$ .

Theorem: Conditional gradient method using fixed step sizes  $\gamma_k = \frac{2}{(k+1)}$   
 $k=1, 2, \dots$  satisfies

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

What does this mean?

# iteration to achieve  $\epsilon$  error is  $O(1/\epsilon)$

Convergence rate matches that of projected gradient when  $\nabla f$  is Lipschitz.  
But are the assumptions comparable?

fact: If  $\nabla f$  is Lipschitz with Lipschitz constant  $L$ , then  
 $M \leq \text{diam}(C) \times L$  where

$$\text{diam}(C) = \max_{x, s \in C} \|x - s\|_2$$

Recall:

$$f(y) - f(x) - \nabla f(x)^T(y-x) \leq \frac{L}{2} \|y-x\|_2^2 \quad \leftarrow \begin{array}{l} L \text{ Lipschitz condition} \\ (\text{one of the several equivalent}) \end{array}$$

Max. over all  $y = ((1-\gamma)x + \gamma s)$  and mul by  $2/\gamma^2$

$$M \leq \max_{\substack{x, s, y \in C \\ y = (1-\gamma)x + \gamma s}} \left( \frac{2}{\gamma^2} \right) \cdot \frac{L}{2} \|y-x\|_2^2 = \max_{x, s \in C} L \|x-s\|_2^2$$

Essentially assuming bounded curvature is no stronger than assuming Lipschitz.

Convergence

key inequality  $f(x_k) \leq f(x_{k-1}) - \gamma_k g(x_{k-1}) + \frac{\gamma_k^2 M}{2}$

$$g(x) = \max_{s \in C} \nabla f(x)^T(s-x)$$

$$\begin{aligned} x^+ &\leftarrow x_k \\ x &\leftarrow x_{k-1} \end{aligned}$$

$$\begin{aligned} f(x^+) &= f(x + \gamma(s-x)) \\ &\leq f(x) + \gamma \nabla f(x)^T(s-x) + \frac{\gamma^2 M}{2} \end{aligned}$$

$$\text{Let } h(x) = f(x) - f(x^*)$$

$$= f(x) - \gamma g(x) + \frac{\gamma^2 M}{2}$$

claim  $g(x) \geq h(x)$  always.

aside

$$\boxed{\begin{aligned} f(y) &\geq f(x) + \nabla f(x)^T(y-x) \\ f(y) - f(x) &\geq \nabla f(x)^T(y-x) \\ f(x^*) - f(x^*) &\geq \cancel{\nabla f(x)^T(x^*-x)} \\ f(x) - f(x^*) &\leq \nabla f(x)^T(x-x^*) \leq \max_g \nabla f(x)^T(x-s) \end{aligned}}$$

$$h(x^+) \leq h(x) - \gamma g(x) + \frac{\gamma^2}{2} M$$

$$h(x^+) \leq h(x)(1-\gamma) + \frac{\gamma^2}{2} M \quad \underline{\text{Proof is left as homework.}}$$

$$\Rightarrow h(x_k) \leq \frac{2M}{k+2}$$