

# Back to Optimization with constraints

$$\begin{aligned} \min & f(x) \\ \text{s.t.} & g_i(x) \leq 0 \quad i=1 \dots m \end{aligned}$$

Temporarily absorb  $h_j$ 's as 2 neg constraints

$$h_j(x) \leq 0 \quad j=1 \dots l$$

We need  $-h_j$  &  $h_j$  both convex & affine  $h_j$

① From midsem question

Define:  $T_{g_i}(x) = 0$  if  $g_i(x) \leq 0$  &  $= \infty$  o/w

$$\min_x f(x) + \sum_i T_{g_i}(x) \quad \text{:- convex function but not diff.}$$

Solve by either analysing optimality conditions in terms of subgradients or employ subgradient descent...

② Write it equivalently as a one program (yet to be analysed)... MIDSEM

③ Replace  $I g_i(x)$  with a more "graceful" penalty function

$$\min_x f(x) - \sum_{i=1}^m \lambda_i \log(-g_i(x))$$

iteratively decrease  $\lambda_i \geq 0$

④ Instead consider the Lagrangian fn

$$L(x, \lambda) = f(x) + \sum \lambda_i g_i(x)$$

We will briefly visit ① & then ④  
& later ② & ③

⑤ Recall gradient descent & Newton:

$$x^{k+1} = \min_x f(x^k) + \nabla^T f(x^k) (x - x^k) + \frac{1}{2} (x - x^k)^T M (x - x^k)$$

$M = I$  for gradient desc &  $\nabla^2 f(x^k)$  for Newton

"Proximal" / "Mirror descent" / Projection  
algorithms treat problem of finding  
 $x^{k+1}$  as that of locating next iterate  
as close as possible to  $x^k$

↓  
in the sense of an approximation  
or in the sense of minimizing  
constraint violation etc

Recap from Midsem for (1)

(10 Marks)

5. Consider a constrained convex optimization problem:

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && g_i(\mathbf{x}) \leq 0 \quad \text{for } i = 1 \dots m \end{aligned} \tag{2}$$

where  $f$  and  $g_i$ 's are closed convex functions.

We will discuss two ways of reformulating this problem:

Let  $S_{g_i} = \{y \mid g_i(y) \leq 0\}$   
 $S_{g_i}$  is a convex set  
 (sublevel set of convex fn)

- (a) Consider indicator function  $I_{g_i}(\mathbf{x})$  associated with  $g_i(\mathbf{x})$  for each  $i = 1 \dots m$  such that  $I_{g_i}(\mathbf{x}) = 0$  iff  $g_i(\mathbf{x}) \leq 0$  and  $I_{g_i}(\mathbf{x}) = 1$  otherwise.
- Prove that  $\partial I_{g_i}(\mathbf{x})$  is a convex cone. Is it closed?
  - Pose (2) as an equivalent unconstrained convex optimization problem making use of  $I_{g_i}(\mathbf{x})$ .
  - Now derive a necessary and sufficient condition for global constrained optimality of (2) at a point  $\mathbf{x}^*$ .

(5 Marks)

(a) (i) We can prove that  $\partial I_{g_i}(\mathbf{x}) = \{g \mid g^T y \leq g^T x \quad \forall y \in S_{g_i}\}$

$$I_{g_i}(y) \geq I_{g_i}(x) + \text{subgrad}^T(x)(y-x) \quad \forall y$$

If  $g_i(x) \leq 0$  then  $0 \geq \text{subgrad}^T(x)(y-x)$  is necessary & sufficient  
 $\leq \partial I_{g_i}(x) = \{g \mid g^T y \leq g^T x \quad \forall y \in S_{g_i}\} = \bigcap_{y \in S_{g_i}} \{g \mid g^T (y-x) \leq 0\}$  = intersection of half spaces of infinite hyperplane through origin = closed convex cone

$\therefore I_{g_i}(x)$  is a closed convex cone  $\forall x: g_i(x) \leq 0$

$I_{g_i}(x) = \emptyset$  if  $x: g_i(x) > 0$

(ii) There are multiple ways to achieve it:-

$$\lim_{\lambda \rightarrow \infty} \min_x f(x) + \lambda \sum_i I_{g_i}(x) \dots \text{Equivalent to redefining}$$

$$I_{g_i}(x) = 0 \quad \text{if } g_i(x) \leq 0$$

$$= \infty \quad \text{o/w}$$

& using  $\min_x f(x) + \sum_i \lambda_i I_{g_i}(x)$

OR

$$\min_x \max_{\lambda_i} f(x) + \sum_i \lambda_i I_{g_i}(x)$$

(iii) Necessary & sufficient conditions for optimality

(from Q4) : Let  $I_{g_i}(x) = 0$  if  $g_i(x) \leq 0$   
 $= \infty$  o/w

$$0 \in \partial(f(x) + \sum I_{g_i}(x))$$

$$\partial I_{g_i}(x) = \{g \mid g^T y \leq g^T x \quad \forall y: g_i(y) \leq 0\}$$

$$\text{if } g_i(x) \leq 0$$

&  $\partial I_{g_i}(x) = \emptyset$  o/w

# Subgradient Methods

- subgradient method and stepsize rules
- convergence results and proof
- optimal step size and alternating projections
- speeding up subgradient methods

$g_x$  is a subgradient at  $x$  if  
 $f(y) \geq f(x) + g_x^T(y-x) \quad \forall y \in \text{dom}(f)$

# Subgradient method

**subgradient method** is simple algorithm to minimize nondifferentiable convex function  $f$

$$x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)}$$

- $x^{(k)}$  is the  $k$ th iterate
- $g^{(k)}$  is **any** subgradient of  $f$  at  $x^{(k)}$
- $\alpha_k > 0$  is the  $k$ th step size

not a descent method, so we keep track of best point so far

$$f_{\text{best}}^{(k)} = \min_{i=1, \dots, k} f(x^{(i)})$$

$f(x^k) \neq f_{\text{best}}^{(k)}$   
not necessarily

## Step size rules

step sizes are fixed ahead of time

- *constant step size*:  $\alpha_k = \alpha$  (constant)
- *constant step length*:  $\alpha_k = \gamma / \|g^{(k)}\|_2$  (so  $\|x^{(k+1)} - x^{(k)}\|_2 = \gamma$ )
- *square summable but not summable*: step sizes satisfy

$$\sum_{k=1}^{\infty} \alpha_k^2 < \infty, \quad \sum_{k=1}^{\infty} \alpha_k = \infty$$

- *nonsummable diminishing*: step sizes satisfy

$$\lim_{k \rightarrow \infty} \alpha_k = 0, \quad \sum_{k=1}^{\infty} \alpha_k = \infty$$



## Assumptions

- $f^* = \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$
- $\|g\|_2 \leq G$  for all  $g \in \partial f$  (equivalent to Lipschitz condition on  $f$ )
- $\|x^{(1)} - x^*\|_2 \leq R$

(H/w)

$$|f(y) - f(x)| \leq G \|x - y\|$$

these assumptions are stronger than needed, just to simplify proofs

## Convergence results

define  $\bar{f} = \lim_{k \rightarrow \infty} f_{\text{best}}^{(k)}$

- *constant step size*:  $\bar{f} - f^* \leq G^2\alpha/2$ , *i.e.*,  
**converges to  $G^2\alpha/2$ -suboptimal**  
(converges to  $f^*$  if  $f$  differentiable,  $\alpha$  small enough)
- *constant step length*:  $\bar{f} - f^* \leq G\gamma/2$ , *i.e.*,  
**converges to  $G\gamma/2$ -suboptimal**
- *diminishing step size rule*:  $\bar{f} = f^*$ , *i.e.*, **converges**

(Recall: If  $f$  is Lipschitz, gradient descent gives  $O(1/\epsilon^2)$  convergence: That proof comes from here)

## Convergence proof

**key quantity:** *Euclidean distance to the optimal set*, not the function value

let  $x^*$  be any minimizer of  $f$

$$\begin{aligned}\|x^{(k+1)} - x^*\|_2^2 &= \|x^{(k)} - \alpha_k g^{(k)} - x^*\|_2^2 \\ &= \|x^{(k)} - x^*\|_2^2 - 2\alpha_k g^{(k)T} (x^{(k)} - x^*) + \alpha_k^2 \|g^{(k)}\|_2^2 \\ &\leq \|x^{(k)} - x^*\|_2^2 - 2\alpha_k \underbrace{(f(x^{(k)}) - f^*)}_{\geq 0} + \alpha_k^2 \|g^{(k)}\|_2^2\end{aligned}$$

using  $f^* = f(x^*) \geq f(x^{(k)}) + g^{(k)T}(x^* - x^{(k)})$

apply recursively to get

$$\begin{aligned} \|x^{(k+1)} - x^*\|_2^2 &\leq \|x^{(1)} - x^*\|_2^2 - 2 \sum_{i=1}^k \alpha_i (f(x^{(i)}) - f^*) + \sum_{i=1}^k \alpha_i^2 \|g^{(i)}\|_2^2 \\ &\leq R^2 - 2 \sum_{i=1}^k \alpha_i (f(x^{(i)}) - f^*) + G^2 \sum_{i=1}^k \alpha_i^2 \end{aligned}$$

now we use

$$\sum_{i=1}^k \alpha_i (f(x^{(i)}) - f^*) \geq (f_{\text{best}}^{(k)} - f^*) \left( \sum_{i=1}^k \alpha_i \right)$$

to get

$$f_{\text{best}}^{(k)} - f^* \leq \frac{R^2 + G^2 \sum_{i=1}^k \alpha_i^2}{2 \sum_{i=1}^k \alpha_i}.$$

Note of why summability of  $\sum \alpha_i$  becomes important

$$f_{\text{best}}^{(k)} = \min_{i=0 \dots k} f(x^{(i)})$$

**constant step size:** for  $\alpha_k = \alpha$  we get

*we always had  $\|x^{(i)} - x^*\|_2 = R$*

$$f_{\text{best}}^{(k)} - f^* \leq \frac{R^2 + G^2 k \alpha^2}{2k\alpha}$$

righthand side converges to  $G^2 \alpha / 2$  as  $k \rightarrow \infty$

**constant step length:** for  $\alpha_k = \gamma / \|g^{(k)}\|_2$  we get

$$f_{\text{best}}^{(k)} - f^* \leq \frac{R^2 + \sum_{i=1}^k \alpha_i^2 \|g^{(i)}\|_2^2}{2 \sum_{i=1}^k \alpha_i} \leq \frac{R^2 + \gamma^2 k}{2\gamma k / G},$$

righthand side converges to  $G\gamma / 2$  as  $k \rightarrow \infty$

**square summable but not summable step sizes:**

suppose step sizes satisfy

$$\sum_{k=1}^{\infty} \alpha_k^2 < \infty, \quad \sum_{k=1}^{\infty} \alpha_k = \infty$$

then

$$f_{\text{best}}^{(k)} - f^* \leq \frac{R^2 + G^2 \sum_{i=1}^k \alpha_i^2}{2 \sum_{i=1}^k \alpha_i}$$

as  $k \rightarrow \infty$ , numerator converges to a finite number, denominator converges to  $\infty$ , so  $f_{\text{best}}^{(k)} \rightarrow f^*$

## Stopping criterion

- terminating when  $\frac{R^2 + G^2 \sum_{i=1}^k \alpha_i^2}{2 \sum_{i=1}^k \alpha_i} \leq \epsilon$  is really, really, slow
- optimal choice of  $\alpha_i$  to achieve  $\frac{R^2 + G^2 \sum_{i=1}^k \alpha_i^2}{2 \sum_{i=1}^k \alpha_i} \leq \epsilon$  for smallest  $k$ :

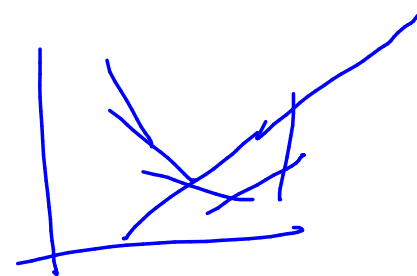
$$\alpha_i = (R/G)/\sqrt{k}, \quad i = 1, \dots, k$$

number of steps required:  $k = (RG/\epsilon)^2$

- the truth: there really isn't a good stopping criterion for the subgradient method . . .

## Example: Piecewise linear minimization

$$\text{minimize } f(x) = \max_{i=1, \dots, m} (a_i^T x + b_i)$$



to find a subgradient of  $f$ : find index  $j$  for which

$$a_j^T x + b_j = \max_{i=1, \dots, m} (a_i^T x + b_i)$$

and take  $g = a_j$

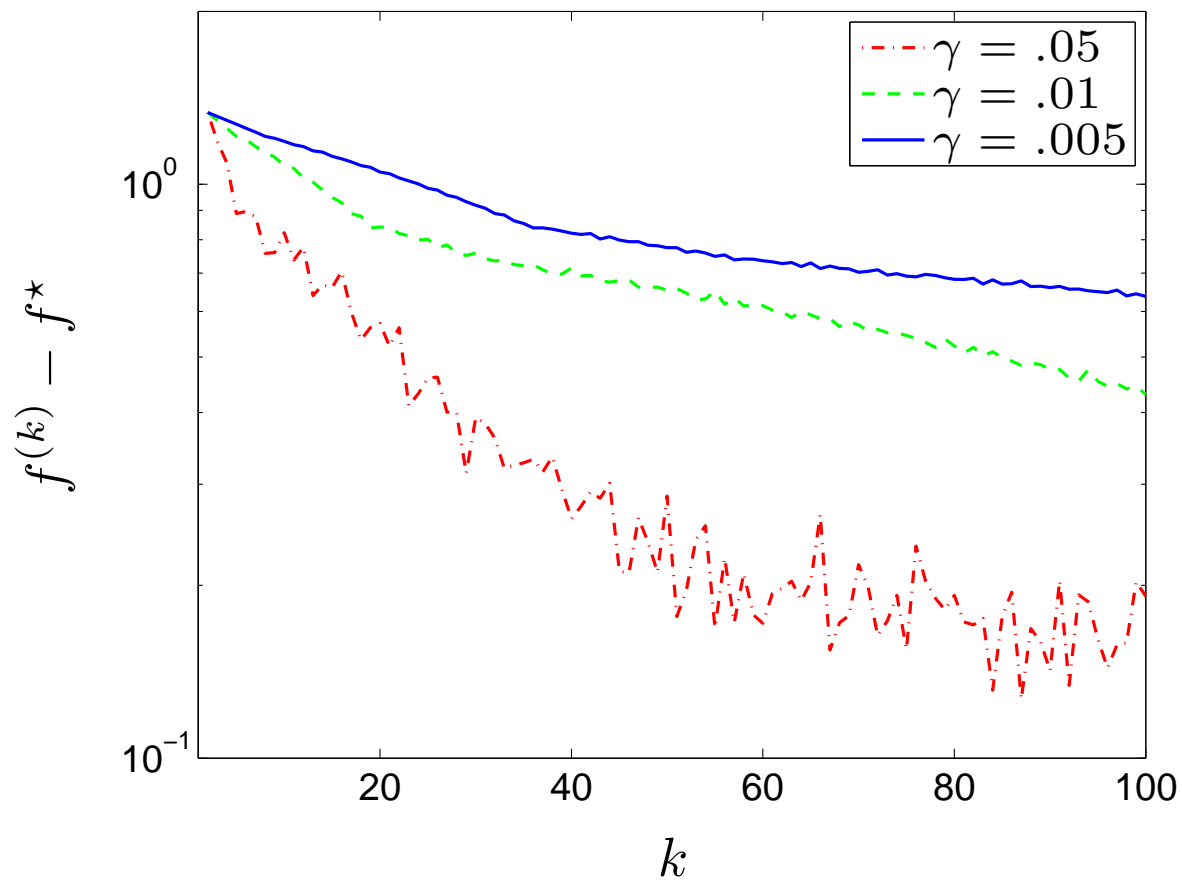
subgradient method:  $x^{(k+1)} = x^{(k)} - \alpha_k a_j$

Note: This problem is equivalent to the following constrained opt problem:

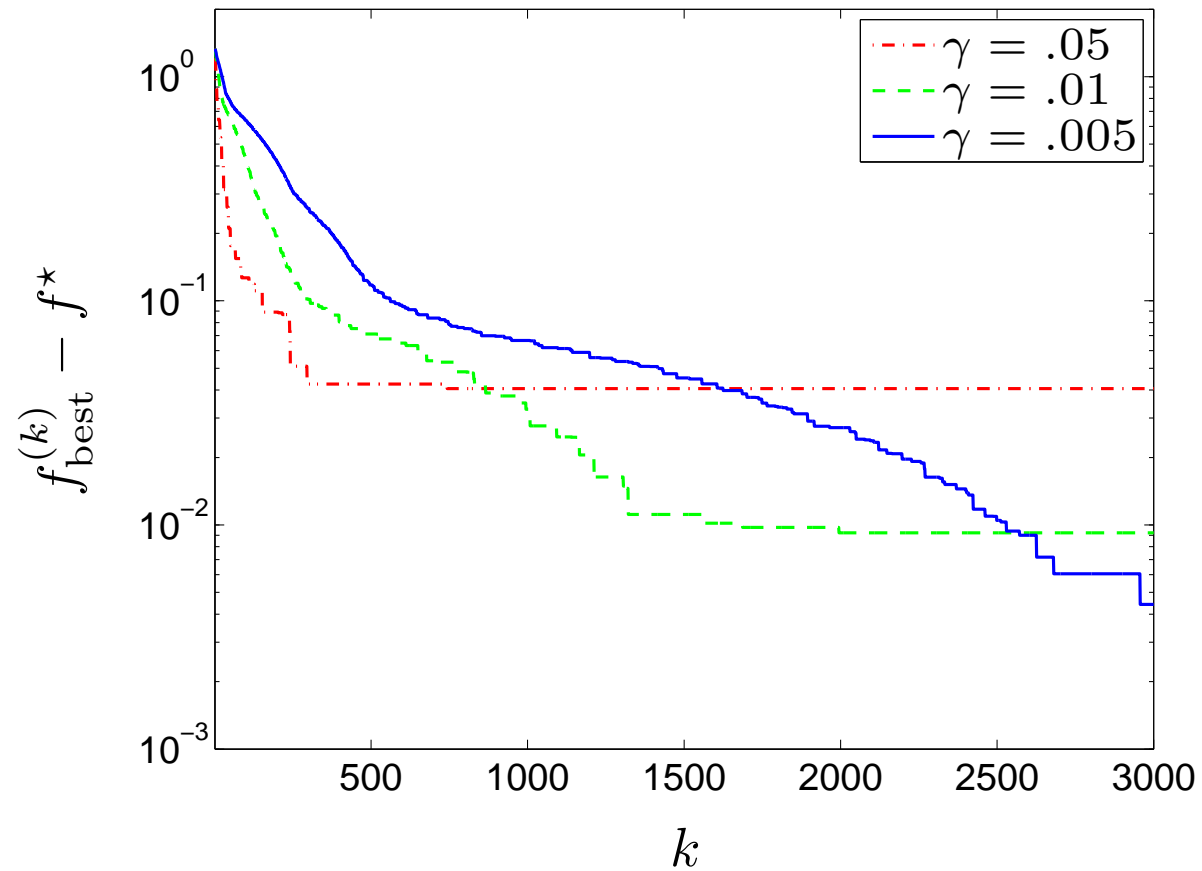
$$\min t$$
$$a_i^T x + b_i - t \leq 0 \quad \forall i$$



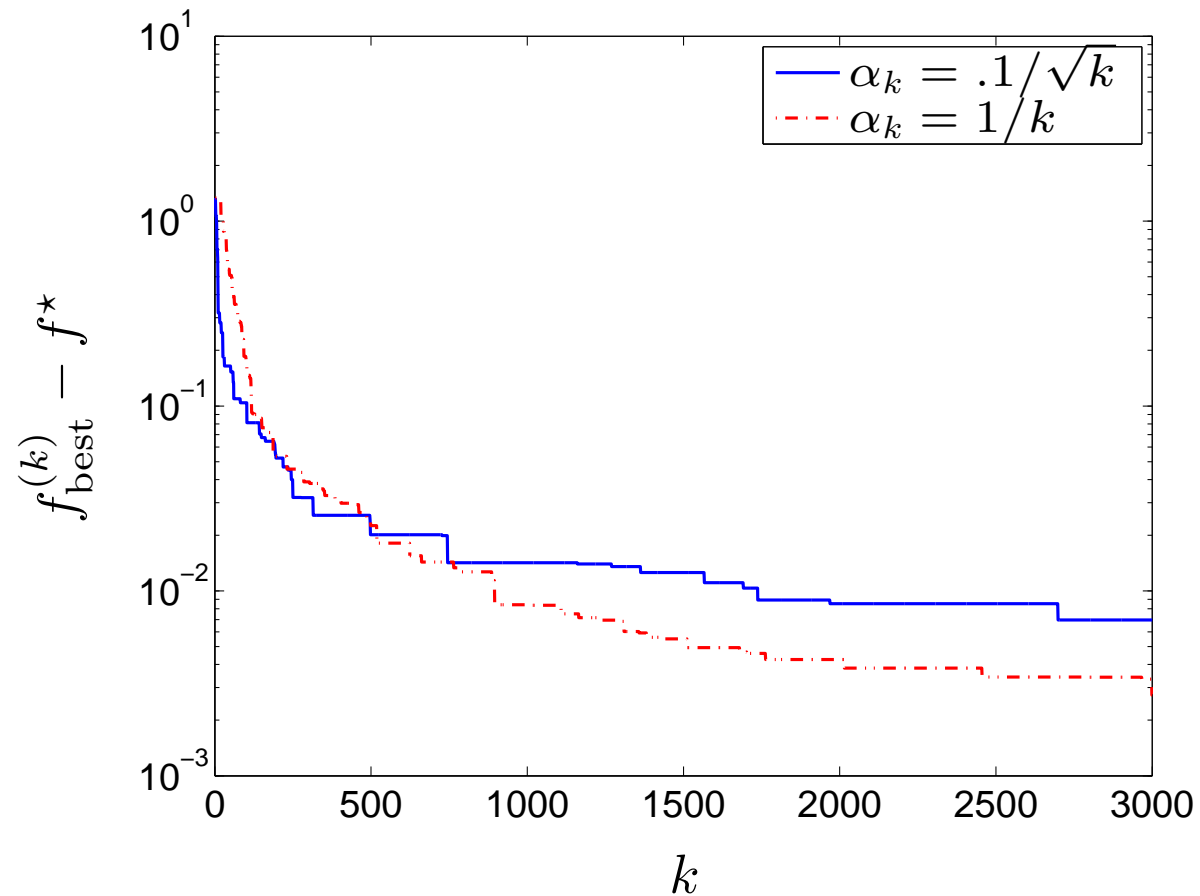
problem instance with  $n = 20$  variables,  $m = 100$  terms,  $f^* \approx 1.1$   
constant step length,  $\gamma = 0.05, 0.01, 0.005$ , first 100 iterations



$f_{\text{best}}^{(k)} - f^*$ , constant step length  $\gamma = 0.05, 0.01, 0.005$



diminishing step rule  $\alpha_k = 0.1/\sqrt{k}$  and square summable step size rule  $\alpha_k = 1/k$



Suppose!

$$\begin{aligned} \min & f(x) \\ \text{s.t.} & g_i(x) \leq 0 \end{aligned}$$

$$\min f(x) + \eta \max_i g_i(x)$$

(we let  $\eta$  iteratively tend to  $\infty$ )

You need to find the formulation of the constrained opt problem for which the subgradient can be discovered easily.

H/W

Eg Lasso:  $\min_x \|Ax - b\|_2^2 \rightarrow \text{Regression loss/error}$   
 $\|x\|_1 \leq \theta$

## Finding a point in the intersection of convex sets

$$C_i = \{x \mid g_i(x) \leq 0\}$$

$C = C_1 \cap \dots \cap C_m$  is nonempty,  $C_1, \dots, C_m \subseteq \mathbf{R}^n$  closed and convex

find a point in  $C$  by minimizing

$$f(x) = \max\{\mathbf{dist}(x, C_1), \dots, \mathbf{dist}(x, C_m)\}$$

with  $\mathbf{dist}(x, C_j) = f(x)$ , a subgradient of  $f$  is

$$g = \nabla \mathbf{dist}(x, C_j) = \frac{x - P_{C_j}(x)}{\|x - P_{C_j}(x)\|_2}$$

*x could be  
an iterative  
obtained  
using gradient descent*

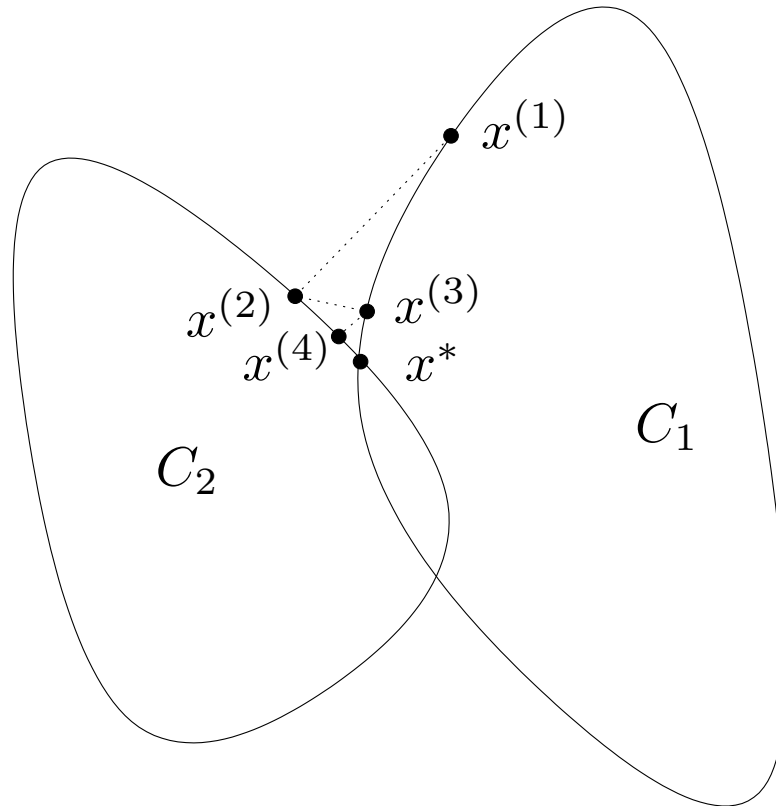
subgradient update with optimal step size:

$$\begin{aligned}x^{(k+1)} &= x^{(k)} - \alpha_k g^{(k)} \\ &= x^{(k)} - f(x^{(k)}) \frac{x^{(k)} - P_{C_j}(x^{(k)})}{\|x^{(k)} - P_{C_j}(x^{(k)})\|_2} \\ &= P_{C_j}(x^{(k)})\end{aligned}$$

- a version of the famous *alternating projections* algorithm
- at each step, project the current point onto the farthest set
- for  $m = 2$  sets, projections alternate onto one set, then the other
- convergence:  $\mathbf{dist}(x^{(k)}, C) \rightarrow 0$  as  $k \rightarrow \infty$

# Alternating projections

first few iterations:



...  $x^{(k)}$  eventually converges to a point  $x^* \in C_1 \cap C_2$

## Speeding up subgradient methods

- subgradient methods are very slow
- often convergence can be improved by keeping memory of past steps

$$x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)} + \beta_k (x^{(k)} - x^{(k-1)})$$

(heavy ball method)

keeping track of several previous iterates

**other ideas:** localization methods, conjugate directions, . . .

Only approximately conjugate for non quadratic problems.



Now visiting

④

$$L(x, \lambda, \mu) = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

for  $\min f(x)$   
s.t.  $g_i(x) \leq 0 \quad i=1 \dots m$   
 $h_j(x) = 0 \quad j=1 \dots l$

$\min L(x, \lambda, \mu)$  --- you should ideally have  $\lambda_i \geq 0$   
to penalize  $g_i(x) > 0$

$$\min_x f(x) \geq \min_x f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

$g_i(x) \leq 0$   
 $h_j(x) = 0$

$g_i(x) \leq 0, \lambda_i \geq 0$   
 $h_j(x) = 0$

$$\geq \min_{x, \lambda_i \geq 0} L(x, \lambda, \mu)$$

$$\begin{aligned} \min_x f(x) \\ \text{s.t. } g_i(x) \leq 0 \\ h_j(x) = 0 \end{aligned}$$

$$\geq \max_{\lambda \geq 0} \quad \text{(circled in red)$$

$$\min_x L(x, \lambda, u)$$

Pushes up the lower bound from previous inequality.

$$\begin{cases} \min f(x) \\ \text{s.t. } g_i(x) \leq 0 \quad i=1 \dots m \\ h_j(x) = 0 \quad j=1 \dots k \end{cases}$$

We will generalize the inequalities & equalities

$$\begin{aligned} \min_x f(x) &\geq \min_x \max_{\lambda, \mu} \underbrace{f(x) + \sum_{i=1}^m \lambda_i g_i(x) + \sum_{j=1}^k \mu_j h_j(x)}_{L(x, \lambda, \mu)} \\ \text{s.t. } g_i(x) &\leq 0 \\ h_j(x) &= 0 \\ \lambda_i &\geq 0 \quad \mu_j \in \mathbb{R} \end{aligned}$$

$$\geq \min_x \max_{\lambda, \mu} L(x, \lambda, \mu) \quad \left\{ \begin{array}{l} \text{under strong} \\ \text{duality} \\ \lambda_i^* g_i(x^*) = 0 \\ \forall i \\ \lambda_i^* \geq 0 \\ \mu_j^* h_j(x^*) = 0 \\ \forall j \end{array} \right.$$

$$\geq \max_{\lambda, \mu} \min_x L(x, \lambda, \mu) \quad \left\{ \begin{array}{l} \lambda_i \geq 0 \\ \mu_j \in \mathbb{R} \end{array} \right.$$

General weak duality result

$L^*(\lambda, \mu)$  or Lagrange dual fn.

$$= \max_{\lambda, \mu, \lambda \geq 0} L^*(\lambda, \mu)$$

Dual opt problem

$$\min_x f(x) \leq \max_{\substack{\lambda_i \geq 0 \\ \mu \in \mathbb{R}}} L^*(\lambda, \mu)$$

s.t.  $g_i(x) \leq 0$   
 $h_j(x) = 0$

Q1: Did we require  $f$ ,  $g_i$ 's &  $h_j$ 's to be convex or affine? ANS: No

Q2: Is  $L^*$  concave irrespective of  $f$ ,  $g_i$ 's &  $h_j$ 's? Note:  $L(x, \lambda, \mu)$  is affine in  $\lambda, \mu$

$$L^* = \min_x \underbrace{L(x, \lambda, \mu)}_{L_x(\lambda, \mu)}$$



min of affine fns is concave

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

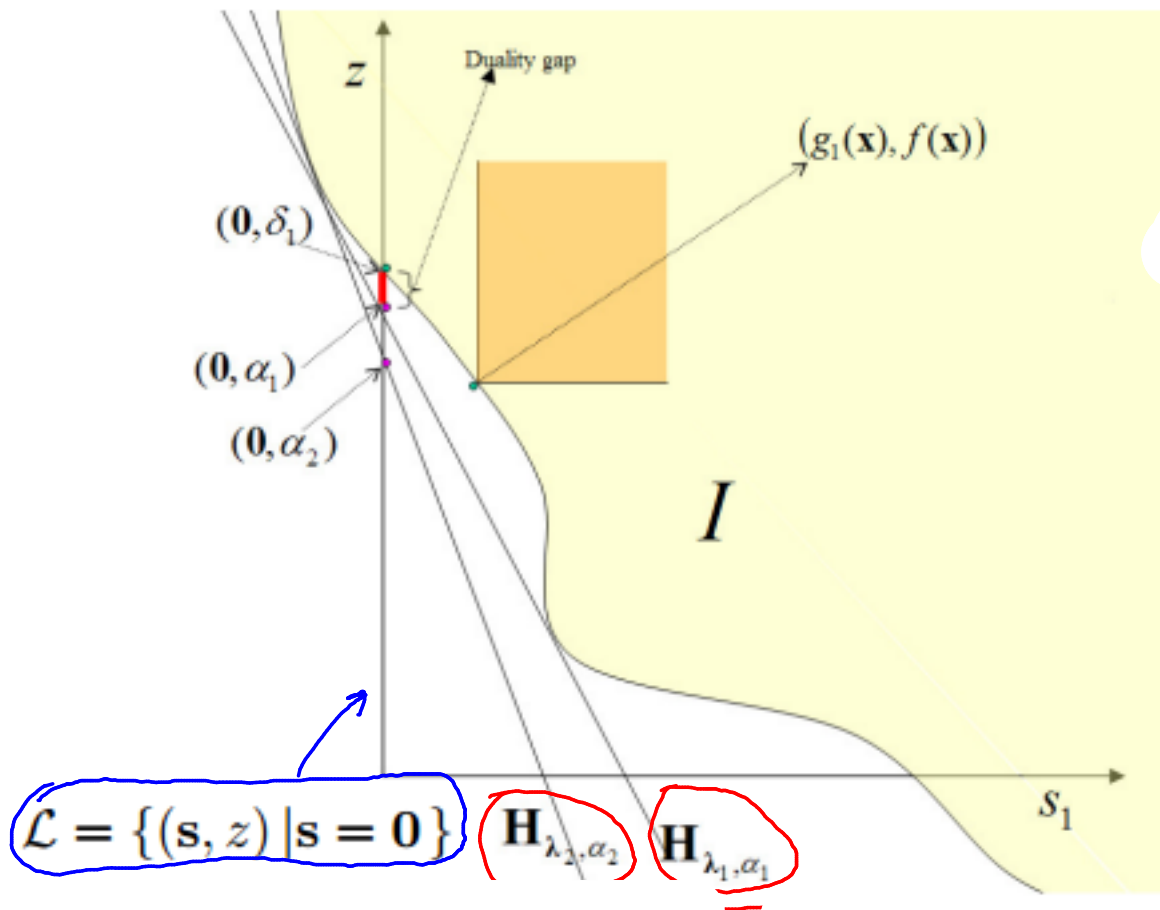
$$\text{s.t. } g_i(x) \leq 0 \quad i=1, \dots, m$$

(The general dual problem & its geometric interpretation)

pg 292, sec 4.4.3 of <http://www.cse.iitb.ac.in/~cs709/notes/BasicsOfConvexOptimization.pdf>

Consider the set:

$$\mathcal{I} = \{(s, z) \mid s \in \mathbb{R}^m, z \in \mathbb{R}, \exists x \in \mathcal{D} \text{ with } g_i(x) \leq s_i \forall 1 \leq i \leq m, f(x) \leq z\}$$



$$\mathcal{H}_{\lambda, \alpha} = \{(s, z) \mid \lambda^T \cdot s + z = \alpha\}$$