

Subgradient Methods

Wing-Kin (Ken) Ma

Department of Electronic Engineering,
The Chinese University Hong Kong, Hong Kong

Lesson 13, ELEG5481

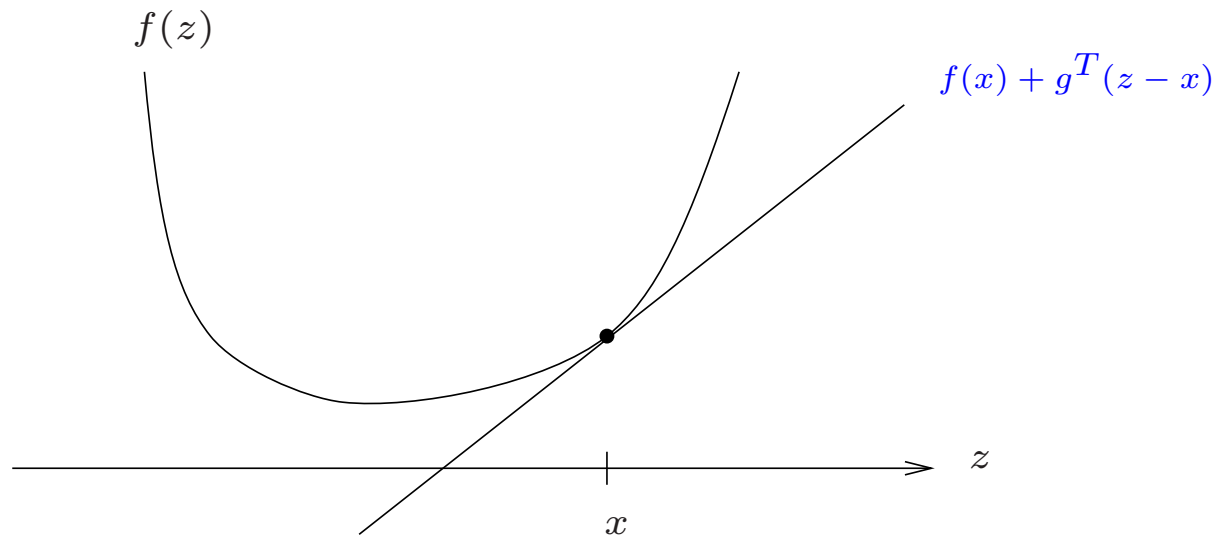
Subgradient Methods

- Subgradient methods are a class of simple methods for solving convex problems, including those with nondifferentiable functions.
- developed in the Soviet Union in the 60's and 70's by Shor and others.
- can be slow (perhaps very slow) in convergence.
- can be applied to many different problems, including those where interior-point methods cannot be used.
- can be used to decouple or decompose a large problem into many smaller ones. This has played a significant role in internet optimization, network utility max., and dynamic spectrum management in multiuser multicarrier systems.

Definition of Subgradient

- A vector $\mathbf{g} \in \mathbb{R}^n$ is said to be a *subgradient* of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at $\mathbf{x} \in \text{dom} f$ if, for all $\mathbf{z} \in \text{dom} f$,

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x})$$



- If f is convex and differentiable, then its gradient $\nabla f(\mathbf{x})$ at \mathbf{x} is a subgradient.
- A subgradient can exist even when f is nondifferentiable at \mathbf{x} .

Subdifferential

- A function f is called subdifferentiable at x if at least one subgradient of f exists at x .
- The set of all subgradients at x is called the *subdifferential* of f at x , and is denoted as

$$\partial f(x)$$

- A function f is called subdifferentiable if it is subdifferentiable at all $x \in \text{dom } f$.

Example: Absolute value

- Consider $f(x) = |x|$.
- A subgradient of f at x , denoted as g here, is

$$g = \begin{cases} 1, & x > 0 \\ -1, & x < 0 \\ \text{any value between } -1 \text{ and } 1, & x = 0 \end{cases}$$

- The subdifferential is

$$\partial f(x) = \begin{cases} \{1\}, & x > 0 \\ \{-1\}, & x < 0 \\ [-1, 1], & x = 0 \end{cases}$$

- Note that $|x|$ is not differentiable; the derivative does not exist at $x = 0$.

Basic Properties of Subgradients

- $\partial f(\mathbf{x})$ is a closed convex set, even for nonconvex f .
- If f is convex and $\mathbf{x} \in \text{int dom } f$, then $\partial f(\mathbf{x})$ is nonempty and bounded. (that means a convex f is *usually* subdifferentiable)
- If f is convex and differentiable, then

$$\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}.$$

- If f is convex and $\partial f(\mathbf{x}) = \{\mathbf{g}\}$, then f is differentiable at \mathbf{x} .
- \mathbf{x}^* is a minimizer of a convex f if and only if f is subdifferentiable at \mathbf{x}^* and

$$\mathbf{0} \in \partial f(\mathbf{x}^*).$$

Calculus of Subgradients

- nonnegative scaling: for $\alpha \geq 0$,

$$\partial(\alpha f)(\mathbf{x}) = \alpha \partial f(\mathbf{x})$$

- sum: Suppose $f = f_1 + \dots + f_m$, f_i all being convex.

$$\partial f(\mathbf{x}) = \partial f_1(\mathbf{x}) + \dots + \partial f_m(\mathbf{x})$$

The same property applies to integrals.

- affine transformation of domain: Suppose f is convex, and let $h(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$.

$$\partial h(\mathbf{x}) = \mathbf{A}^T \partial f(\mathbf{A}\mathbf{x} + \mathbf{b}).$$

- pointwise max.: Suppose f_1, \dots, f_m are convex, and let $f(\mathbf{x}) = \max_{i=1, \dots, m} f_i(\mathbf{x})$.

$$\partial f(\mathbf{x}) = \text{conv} \cup \{ \partial f_i(\mathbf{x}) \mid f_i(\mathbf{x}) = f(\mathbf{x}) \}$$

Example: Pointwise Linear Function

- Consider

$$f(\mathbf{x}) = \max_{i=1,\dots,m} \mathbf{a}_i^T \mathbf{x} + b_i$$

- Let $f_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{x} + b_i$. We have $\partial f_i(\mathbf{x}) = \{\mathbf{a}_i\}$.

- Let $\mathcal{K}(\mathbf{x}) = \left\{ j \mid \mathbf{a}_j^T \mathbf{x} + b_j = \max_{i=1,\dots,m} \mathbf{a}_i^T \mathbf{x} + b_i \right\}$.

$$\partial f(\mathbf{x}) = \text{conv} \bigcup_{j \in \mathcal{K}(\mathbf{x})} \{\mathbf{a}_j\}$$

- In particular, when $\mathcal{K}(\mathbf{x}) = \{k\}$, we have $\partial f(\mathbf{x}) = \{\mathbf{a}_k\}$.

Example: 1-norm

- Consider

$$f(\mathbf{x}) = \|\mathbf{x}\|_1 = \underbrace{|x_1|}_{f_1} + \dots + \underbrace{|x_n|}_{f_n}$$

- Its subdifferential is

$$\begin{aligned}\partial f(\mathbf{x}) &= \partial f_1(\mathbf{x}) + \dots + \partial f_m(\mathbf{x}) \\ &= \{\mathbf{g} \mid g_i = 1 \text{ if } x_i > 0, g_i = -1 \text{ if } x_i < 0, g_i \in [-1, 1] \text{ if } x_i = 0\}\end{aligned}$$

- Alternatively,

$$f(\mathbf{x}) = \max_{\mathbf{s} \in \{-1, 1\}^n} \underbrace{\mathbf{s}^T \mathbf{x}}_{f_{\mathbf{s}}(\mathbf{x})}$$

and

$$\begin{aligned}\partial f(\mathbf{x}) &= \text{conv} \bigcup \{\partial f_{\mathbf{s}}(\mathbf{x}) \mid \mathbf{s}^T \mathbf{x} = \|\mathbf{x}\|_1, \mathbf{s} \in \{-1, 1\}^n\} \\ &= \{\mathbf{s} \mid \mathbf{s}^T \mathbf{x} = \|\mathbf{x}\|_1, \mathbf{s} \in [-1, 1]^n\}\end{aligned}$$

- To put it simple, $\text{sign}(\mathbf{x})$ is a subgradient of f at \mathbf{x} .

Supremum

- The pointwise max. result can be extended to supremum. Suppose

$$f(\mathbf{x}) = \sup_{\alpha \in \mathcal{A}} f_{\alpha}(\mathbf{x})$$

where f_{α} are subdifferentiable and \mathcal{A} is compact.

$$\partial f(\mathbf{x}) = \text{conv} \cup \{ \partial f_{\alpha}(\mathbf{x}) \mid f_{\alpha}(\mathbf{x}) = f(\mathbf{x}) \}$$

- **Example:** Consider $f(\mathbf{x}) = \lambda_{\max}(\mathbf{A}(\mathbf{x}))$, $\mathbf{A}(\mathbf{x}) = \mathbf{A}_0 + \sum_{i=1}^n x_i \mathbf{A}_i$. Since

$$\lambda_{\max}(\mathbf{A}(\mathbf{x})) = \sup_{\|\mathbf{y}\|_2=1} f_{\mathbf{y}}(\mathbf{x}), \quad f_{\mathbf{y}}(\mathbf{x}) = \mathbf{y}^T \mathbf{A}(\mathbf{x}) \mathbf{y}$$

we have

$$\partial f(\mathbf{x}) = \text{conv} \cup \{ (\mathbf{y}^T \mathbf{A}_1 \mathbf{y}, \dots, \mathbf{y}^T \mathbf{A}_n \mathbf{y}) \mid \mathbf{y} \text{ a principal eigenvector of } \mathbf{A}(\mathbf{x}) \}$$

In particular, if the max. eigenvector of $\mathbf{A}(\mathbf{x})$, \mathbf{y} , is unique,

$$\partial f(\mathbf{x}) = \{ (\mathbf{y}^T \mathbf{A}_1 \mathbf{y}, \dots, \mathbf{y}^T \mathbf{A}_n \mathbf{y}) \}.$$

The Subgradient Method for Unconstrained Opt.

- The goal is to solve

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

- A basic subgradient method:

given $\{\alpha_k\}$, a step size sequence; & an initial point $\mathbf{x}^{(0)}$.

$k := 0$; $i_{\text{best}} := 0$.

repeat

$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$, where $\mathbf{g}^{(k)}$ is any subgradient of f at $\mathbf{x}^{(k)}$.

$k := k + 1$.

$f_{\text{best}}^{(k)} = \min\{f_{\text{best}}^{(k-1)}, f(\mathbf{x}^{(k)})\}$. If $f(\mathbf{x}^{(k)}) = f_{\text{best}}^{(k)}$, then $i_{\text{best}} := k$.

until a stopping criterion is satisfied.

output $\mathbf{x}^{(i_{\text{best}})}$.

- Look similar to the gradient descent method (for differentiable f), but not the same.
- choose the best point among the generated sequence $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$

Step Size Rules

There are many different choices for the step sizes. Some typical rules are

- Constant step size: $\alpha_k = \alpha$.
- Constant step length: $\alpha_k = \gamma / \|\mathbf{g}^{(k)}\|_2$, where $\gamma > 0$.
- Square summable but not summable: the step sizes satisfy

$$\alpha_k \geq 0, \quad \sum_{k=1}^{\infty} \alpha_k^2 < \infty, \quad \sum_{k=1}^{\infty} \alpha_k = \infty$$

An example is $\alpha_k = a/(b+k)$, where $a, b > 0$.

- Nonsummable diminishing: The step sizes satisfy

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

An example is $\alpha_k = a/\sqrt{k}$, where $a > 0$.

Convergence

Let $f^* = \inf_{\mathbf{x}} f(\mathbf{x})$, and G be such that $\|\mathbf{g}^{(k)}\|_2 \leq G$ for all k .

- Constant step size $\alpha_k = \alpha$:

$$\lim_{k \rightarrow \infty} f_{\text{best}}^{(k)} - f^* \leq G^2 \alpha / 2$$

- Constant step length $\alpha_k = \gamma / \|\mathbf{g}^{(k)}\|_2$, $\gamma > 0$:

$$\lim_{k \rightarrow \infty} f_{\text{best}}^{(k)} - f^* \leq G \gamma / 2$$

- Square summable but not summable; and nonsummable diminishing:

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

Given a solution precision ϵ , the number of iterates k for achieving $f_{\text{best}}^{(k)} - f^* < \epsilon$ can also be proven.

Example: Minimizing a piece-wise linear function

- Consider

$$\min_{\mathbf{x}} \left(\max_{i=1, \dots, m} \mathbf{a}_i^T \mathbf{x} + b_i \right)$$

- At the k iteration, find an (any) index for which

$$\mathbf{a}_j^T \mathbf{x}^{(k)} + b_j = \max_{i=1, \dots, m} \mathbf{a}_i^T \mathbf{x}^{(k)} + b_i$$

and we have

$$\mathbf{g}^{(k)} = \mathbf{a}_j.$$

Example: Solving SDPs

- The basic subgradient method may be used to solve SDPs (are you sure?)
- For simplicity, consider

$$\begin{aligned} \min_{\mathbf{X}} \quad & \text{Tr}(\mathbf{C}\mathbf{X}) \\ \text{s.t.} \quad & X_{ii} = 1, \quad i = 1, \dots, n \\ & \mathbf{X} \succeq \mathbf{0} \end{aligned}$$

which has well-known applications in approximating MAXCUT & ML MIMO detection.

Example: Solving SDPs (cont'd)

- Let us add a redundant equality to the SDP

$$\begin{aligned} \min_{\mathbf{X}} \quad & \text{Tr}(\mathbf{C}\mathbf{X}) \\ \text{s.t.} \quad & \mathbf{X} \succeq \mathbf{0}, \quad X_{ii} = 1, \quad i = 1, \dots, n \\ & \text{Tr}(\mathbf{X}) = n \end{aligned}$$

The dual of the SDP above is

$$\begin{aligned} \max_{\boldsymbol{\mu}, \nu} \quad & -\boldsymbol{\mu}^T \mathbf{1} - n\nu \\ \text{s.t.} \quad & \mathbf{C} + \mathbf{D}(\boldsymbol{\mu}) + \nu \mathbf{I} \succeq \mathbf{0} \end{aligned}$$

- Since

$$\mathbf{C} + \mathbf{D}(\boldsymbol{\mu}) + \nu \mathbf{I} \succeq \mathbf{0} \iff \lambda_{\min}(\mathbf{C} + \mathbf{D}(\boldsymbol{\mu})) \geq -\nu$$

we can rewrite the dual problem as an unconstrained problem

$$\max_{\boldsymbol{\mu}} \quad -\boldsymbol{\mu}^T \mathbf{1} + n\lambda_{\min}(\mathbf{C} + \mathbf{D}(\boldsymbol{\mu}))$$

Example: Solving SDPs (cont'd)

- Now we deal with the dual problem

$$\max_{\boldsymbol{\mu}} d(\boldsymbol{\mu}) \triangleq -\boldsymbol{\mu}^T \mathbf{1} + n\lambda_{\min}(\mathbf{C} + \mathbf{D}(\boldsymbol{\mu}))$$

by subgradient.

- A subgradient of $-d(\boldsymbol{\mu})$ at $\boldsymbol{\mu}$ is

$$\mathbf{g} = \mathbf{1} - n\mathbf{q}_{\min}^2$$

where the superscript 2 denotes the elementwise square, and \mathbf{q}_{\min} is a minimum eigenvector of $\mathbf{C} + \mathbf{D}(\boldsymbol{\mu})$.

Example: Solving SDPs (cont'd)

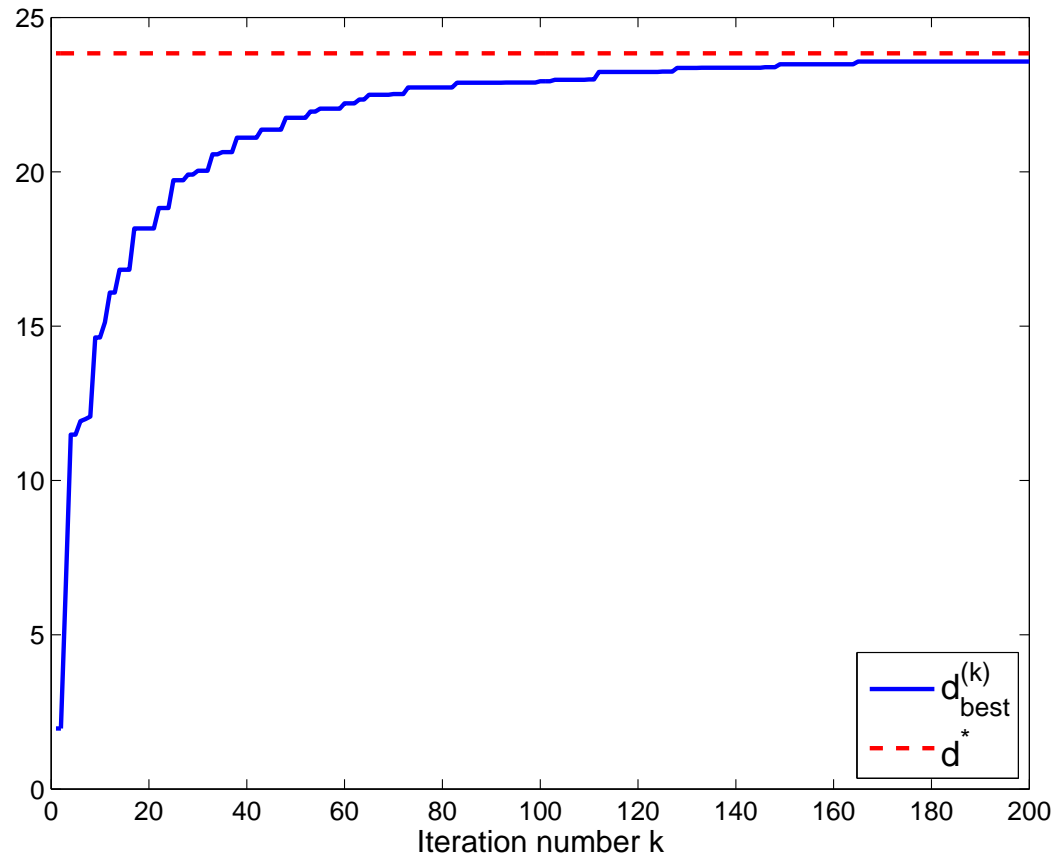


Figure 1: The value $d_{\text{best}}^{(k)}$ versus the iteration number k , for the subgradient method for SDP. The problem size is $n = 20$, and the step size rule is $\alpha_k = 1/\sqrt{k}$.

The Projected Subgradient Method

- The goal is to solve

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

where \mathcal{C} is a convex set.

- In the projected subgradient method, the iterates are obtained by

$$\mathbf{x}^{(k+1)} = \mathcal{P}_{\mathcal{C}} \left(\mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)} \right),$$

where $\mathcal{P}_{\mathcal{C}}$ is the Euclidean projection on \mathcal{C} ; i.e.,

$$\mathcal{P}_{\mathcal{C}}(\mathbf{x}) = \arg \min_{\mathbf{y} \in \mathcal{C}} \|\mathbf{y} - \mathbf{x}\|_2^2$$

- The convergence result is similar to that of the basic subgradient method.

Example: 1-norm minimization

- Consider

$$\begin{aligned} \min \quad & \|\mathbf{x}\|_1 \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b} \end{aligned}$$

where \mathbf{A} is fat.

- We have $\text{sign}(\mathbf{x}) \in \partial f(\mathbf{x})$
- We have $\mathcal{C} = \{\mathbf{x} \mid \mathbf{Ax} = \mathbf{b}\}$, and

$$\mathcal{P}(\mathbf{y}) = \mathbf{A}^\dagger \mathbf{b} + (\mathbf{I} - \mathbf{AA}^\dagger) \mathbf{y},$$

where $\mathbf{A}^\dagger = \mathbf{A}^T (\mathbf{AA}^T)^{-1}$.

- The corresponding projected gradient update is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k (\mathbf{I} - \mathbf{AA}^\dagger) \text{sign}(\mathbf{x})$$

Example: 1-norm minimization (cont'd)

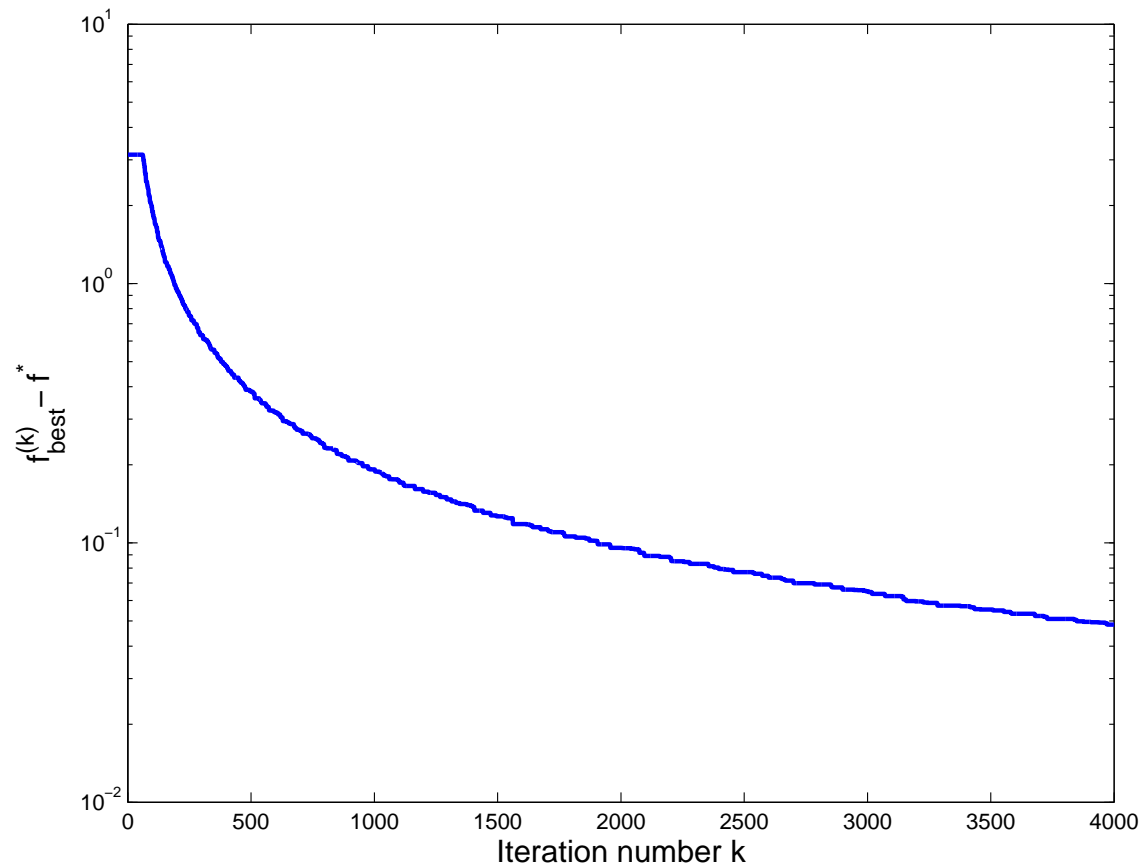


Figure 2: The gap $f_{\text{best}}^{(k)} - f^*$ versus the iteration number k , for the projected subgradient method for 1-norm minimization. The problem size is $m = 50$, $n = 1000$, and the step size rule is $\alpha_k = 0.5/k$.

The projected subgradient method is efficient only when the projection on \mathcal{C} can be easily computed; e.g.,

- An affine set: linear projection
- A halfspace: similar to affine sets
- The set of non-ve nos. $\mathcal{C} = \mathbb{R}_+^n$, a box $\mathcal{C} = \{\mathbf{x} \mid -1 \leq x_i \leq 1, i = 1, \dots, n\}$: projection is truncation
- A 2-norm ball $\mathcal{C} = \{\mathbf{x} \mid \|\mathbf{x}\|_2 \leq 1\}$: projection is rescaling.
- An ellipsoid: no closed form, but can be easily computed.
- Simplex $\mathcal{C} = \{\mathbf{x} \succeq \mathbf{0} \mid \mathbf{x}^T \mathbf{1} \leq 1\}$: no closed form, but can be easily computed.
- The cone of PSD matrices: projection is to discard eigen-components that are -ve.

Projected Subgradient for Dual Problems

- We consider a constrained, not necessarily convex, problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \end{aligned}$$

- We focus on dealing with its dual

$$\begin{aligned} \max_{\boldsymbol{\lambda}} \quad & d(\boldsymbol{\lambda}) \\ \text{s.t.} \quad & \boldsymbol{\lambda} \succeq \mathbf{0} \end{aligned}$$

where $d(\boldsymbol{\lambda}) = \inf_{\mathbf{x}} \left(f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) \right)$.

- Recall that $d(\boldsymbol{\lambda})$ is always concave, even when the primal problem is nonconvex.
- The projected subgradient method can be applied, if we can compute the subgradients of $d(\boldsymbol{\lambda})$.

- Let

$$\mathbf{x}^*(\boldsymbol{\lambda}) = \arg \min_{\mathbf{x}} (f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}))$$

denote a minimizer that attains $d(\boldsymbol{\lambda})$. We can write

$$d(\boldsymbol{\lambda}) = f_0(\mathbf{x}^*(\boldsymbol{\lambda})) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}^*(\boldsymbol{\lambda}))$$

- A subgradient of $-d$ at $\boldsymbol{\lambda}$ is then

$$-(f_1(\mathbf{x}^*(\boldsymbol{\lambda})), \dots, f_m(\mathbf{x}^*(\boldsymbol{\lambda}))) \in \partial(-d)(\boldsymbol{\lambda})$$

- The updates of projected subgradient applied to dual max. is

$$\text{solve for } \mathbf{x}^{(k)} = \arg \min_{\mathbf{x}} \left(f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i^{(k)} f_i(\mathbf{x}) \right)$$

$$\lambda_i^{(k+1)} = \left(\lambda_i^{(k)} + \alpha_k f_i(\mathbf{x}^{(k)}) \right)_+, \quad i = 1, \dots, m$$

where $(\cdot)_+$ is the projection on \mathbb{R}_+ . That is say, we are solving a sequence of unconstrained Lagrangian minimization.

- The projected subgradient method is a Lagrangian dual relaxation, in general. The generated points $\mathbf{x}^{(k)}$ may not be primal feasible.
- Suppose strong duality holds (e.g., convex problems with the Slater condition), and each $\mathbf{x}^{(k)}$ is a unique minimizer. Then the limit point of $\mathbf{x}^{(k)}$ is primal feasible (in fact, optimal).
- This dual max. approach, a.k.a. dual decomposition in some applications, plays a significant role.

Application: Dynamic Spectrum Management (DSM)

- Scenario: A multiuser subcarrier system, with K users and N subcarriers.
- Goal: joint power allocation for sum rate maximization

$$\max \sum_{n=1}^N \sum_{k=1}^K \log(1 + \text{SINR}_k^n(s_1^n, \dots, s_K^n)) \quad (\text{rates of all subcarriers \& users})$$

$$\text{s.t.} \quad \sum_{n=1}^N s_k^n \leq P_k, \quad k = 1, \dots, K \quad (\text{per-user total power constraint})$$

$$0 \leq s_k^n \leq S_{\max}, \quad \forall k, n \quad (\text{per-subcarrier power limit})$$

where the opt. variable s_k^n is the power of k th user at subcarrier n , and

$$\text{SINR}_k^n(s_1^n, \dots, s_K^n) = \frac{\alpha_{kk}^n s_k^n}{\sigma_k^n + \sum_{j \neq k} \alpha_{kj}^n s_j^n}$$

is the SINR of user k at subcarrier n .

Why DSM is hard?

- The DSM sum rate max. problem

$$\begin{aligned} \max \quad & \sum_{n=1}^N \sum_{k=1}^K \log(1 + \text{SINR}_k^n(s_1^n, \dots, s_K^n)) \\ \text{s.t.} \quad & \sum_{n=1}^N s_k^n \leq P_k, \quad k = 1, \dots, K \\ & 0 \leq s_k^n \leq S_{\max}, \quad k = 1, \dots, K, n = 1, \dots, N \end{aligned}$$

is nonconvex, even for $N = 1$.

- It is NP-hard in general.
- The no. of subcarriers N can be large; e.g., $N = 256$, $N = 1024$, ..., and per-user power constraints make the rate max. coupled w.r.t. subcarriers.

- The dual of the DSM problem is

$$\begin{aligned} \min \quad & \mathbf{p}^T \boldsymbol{\lambda} + \varphi(\boldsymbol{\lambda}) \\ \text{s.t.} \quad & \boldsymbol{\lambda} \succeq \mathbf{0} \end{aligned}$$

where $\mathbf{p} = (P_1, \dots, P_K)$,

$$\begin{aligned} \varphi(\boldsymbol{\lambda}) = \max \quad & \left(\sum_{n=1}^N \sum_{k=1}^K \log(1 + \text{SINR}_k^n(s_1^n, \dots, s_K^n)) - \lambda_k s_k^n \right) \\ \text{s.t.} \quad & 0 \leq s_k^n \leq S_{\max}, k = 1, \dots, K, n = 1, \dots, N \end{aligned}$$

- An important result is that

$$\varphi(\boldsymbol{\lambda}) = \sum_{n=1}^N \max_{\substack{s_1^n, \dots, s_K^n, \\ 0 \leq s_k^n \leq S_{\max}}} \left(\sum_{k=1}^K \log(1 + \text{SINR}_k^n(s_1^n, \dots, s_K^n)) - \lambda_k s_k^n \right)$$

i.e., $\varphi(\boldsymbol{\lambda})$ decomposes to many per-subcarrier power allocation problems.

- What remains is that we need to solve the per-subcarrier problems

$$\varphi_n(\boldsymbol{\lambda}) = \max_{\substack{s_1^n, \dots, s_K^n, \\ 0 \leq s_k^n \leq S_{\max}}} \left(\sum_{k=1}^K \log(1 + \text{SINR}_k^n(s_1^n, \dots, s_K^n)) - \lambda_k s_k^n \right)$$

for $n = 1, \dots, N$.

- The problem above is still nonconvex.
- For $K = 2$, exhaustive search was used (OSB [**Cendrillon et al.'06**]).
- For $K > 2$, some approximation methods should be used.
- For the OFDMA variation (one subcarrier can only be occupied by one user), there is a simple way of solving the per-subcarrier problem [**Luo-Zhang'09**].
- (there are many more refs. & nice results in DSM that I have no time to mention here)

Optimal value of a convex opt. problem

- Consider the optimal value of a convex optimization problem

$$\begin{aligned}\phi(\mathbf{x}, \mathbf{y}) &= \min_{\mathbf{z}} f_0(\mathbf{z}) \\ \text{s.t. } & f_i(\mathbf{z}) \leq x_i, \quad i = 1, \dots, m, \quad \mathbf{A}\mathbf{z} = \mathbf{y}\end{aligned}$$

where f_0, f_1, \dots, f_m are convex. Its dual is

$$\begin{aligned}\phi(\mathbf{x}, \mathbf{y}) &= \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}} g(\boldsymbol{\lambda}, \boldsymbol{\mu}) - \mathbf{x}^T \boldsymbol{\lambda} - \mathbf{y}^T \boldsymbol{\mu} \\ \text{s.t. } & \boldsymbol{\lambda} \succeq \mathbf{0}\end{aligned}$$

where $g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{z}} (f_0(\mathbf{z}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{z}) + \boldsymbol{\mu}^T \mathbf{A}\mathbf{z})$.

- Suppose that strong duality holds at (\mathbf{x}, \mathbf{y}) , & let $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ be a dual opt. solution fixing (\mathbf{x}, \mathbf{y}) .

$$-(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \in \partial\phi(\mathbf{x}, \mathbf{y})$$

- This property is useful, e.g., in primal decomposition methods.

Application: MIMO BC Capacity

- Scenario: A multiuser MIMO broadcast channel (BC).
- Goal: Solve the MIMO BC capacity, which has been shown to be

$$\begin{aligned} & \max_{\mathbf{Q}_1, \dots, \mathbf{Q}_K} \log \det \left(\mathbf{I} + \sum_{k=1}^K \mathbf{H}_k \mathbf{Q}_k \mathbf{H}_k^H \right) \\ & \text{s.t.} \quad \sum_{k=1}^K \text{Tr}(\mathbf{Q}_k) \leq P_{\text{total}} \\ & \quad \mathbf{Q}_1, \dots, \mathbf{Q}_K \succeq \mathbf{0} \end{aligned}$$

where \mathbf{H}_k is MIMO channel from the basestation to user k .

- This problem is convex (CVX can do the job).
- Can we derive a simple algorithm by using the subgradient concepts?

A Related Problem: MIMO MAC Capacity

- To solve MIMO BC, let us look at a related problem— MIMO multiple access channel (MAC).

$$\begin{aligned} \max_{\mathbf{Q}_1, \dots, \mathbf{Q}_K \succeq \mathbf{0}} \quad & \log \det \left(\mathbf{I} + \sum_{k=1}^K \mathbf{H}_k \mathbf{Q}_k \mathbf{H}_k^H \right) \\ \text{s.t.} \quad & \text{Tr}(\mathbf{Q}_k) \leq P_k, \quad k = 1, \dots, K \end{aligned}$$

where P_k is the total power limit of user k .

- The MIMO MAC capacity is convex.
- A convenient way of solving the MIMO MAC capacity is to use the iterative water filling algorithm (IWFA):
 - at each iteration, maximize the objective fn. w.r.t. a \mathbf{Q}_k while fixing the other $\{\mathbf{Q}_l\}_{l \neq k}$.
 - the maximization at each iteration is a single-user water filling problem.

Projected Subgradient for MIMO BC Capacity

- We can write the MIMO BC capacity as

$$\max_{P_1, \dots, P_K} \underbrace{\left(\begin{array}{l} \max_{\mathbf{Q}_1, \dots, \mathbf{Q}_K \succeq \mathbf{0}} \quad \log \det \left(\mathbf{I} + \sum_{k=1}^K \mathbf{H}_k \mathbf{Q}_k \mathbf{H}_k^H \right) \\ \text{s.t.} \quad \text{Tr}(\mathbf{Q}_k) \leq P_k, \quad k = 1, \dots, K \end{array} \right)}_{\triangleq \phi(\mathbf{p})}$$

s.t. $\mathbf{p}^T \mathbf{1} \leq P_{\text{total}}$

that is, we use the subgradients of $-\phi(\mathbf{p})$ to solve the MIMO BC capacity.

- Specifically, the updates in projected subgradient are

find $\boldsymbol{\lambda}^{(k)}$ that is an optimal dual solution of $\phi(\mathbf{p}^{(k)})$, by IWFA.

$$\mathbf{p}^{(k+1)} = \left(\mathbf{p}^{(k)} + \alpha_k \boldsymbol{\lambda}^{(k)} \right)_{\mathcal{S}}$$

where \mathcal{S} is the projection on the simplex $\mathcal{S} = \{\mathbf{p} \mid \mathbf{p}^T \mathbf{1} \leq P_{\text{total}}\}$ (no closed form, but can be easily computed [hint: it's like water filling]).

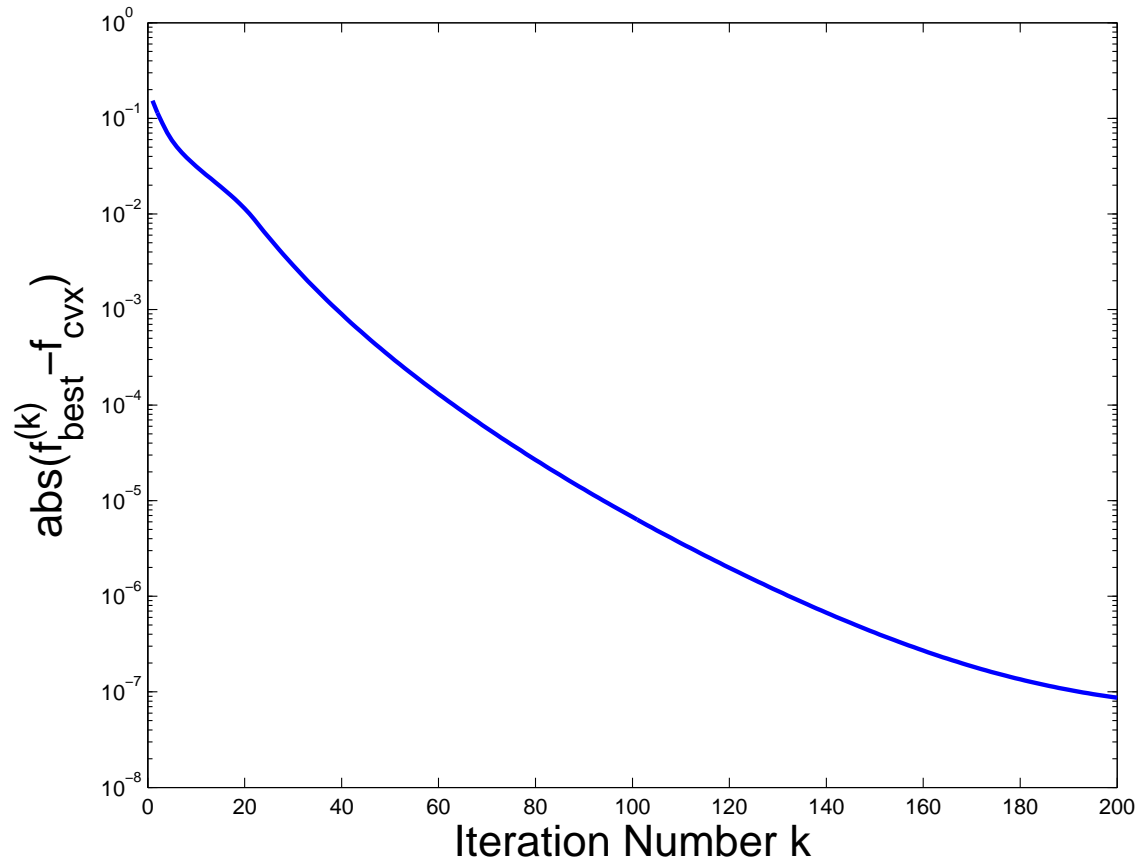


Figure 3: The gap $f_{\text{best}}^{(k)} - f^*$ versus the iteration number k , for the projected subgradient method applied to MIMO BC capacity computations. $M_t = 12$, $M_r = 4$, $K = 3$, $P_{\text{total}} = 100$, and the step size rule is $\alpha_k = 9/\sqrt{k}$.

Subgradient Method for Constrained Optimization

- Consider a convex problem

$$\begin{aligned} \min f_0(\mathbf{x}) \\ \text{s.t. } f_i(\mathbf{x}) \leq 0, i = 1, \dots, m \end{aligned}$$

We have seen that by applying the subgradient method to the dual, the problem can be solved.

- The subgradient method can also be applied directly to the primal.
- The method takes the same form

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

where

$$\mathbf{g}^{(k)} \in \begin{cases} \partial f_0(\mathbf{x}^{(k)}), & f_i(\mathbf{x}^{(k)}) \leq 0, i = 1, \dots, m \\ \partial f_j(\mathbf{x}^{(k)}), & f_j(\mathbf{x}^{(k)}) > 0 \end{cases}$$

Discussion

- Subgradient methods may provide low-complexity implementations to certain problems, but possibly with low accuracy.
- There are approaches that can speed up convergence; e.g., ellipsoid methods, cutting plane methods, bundle methods, ... They are more complex requiring more computations to carry out the update.

References on Subgradient Methods

N. Shor, *Minimization Methods for Non-Differentiable Functions*, Springer Series in Computational Mathematics, Springer, 1985.

D. Bertsekas, *Nonlinear Programming*. Athena Scientific, 1999.

S. Boyd and A. Mutapcic, *Subgradient Methods*, Notes for EE364b, Stanford University, 2006. Available online.

References on Applications

[Cendrillon *et al.*'06] R. Cendrillon, W. Yu, M. Moonen, J. Verlinden, and T. Bostoen, "Optimal multiuser spectrum balancing for digital subscriber lines," *IEEE Trans. Commun.*, vol. 54, no. 5, 2006.

[Luo-Zhang'09] Z.-Q. Luo and S. Zhang, "Duality gap estimation and polynomial time approximation for optimal spectrum management," *IEEE Trans. Signal Process.*, vol. 57, no. 7, 2009.

[Palomar'05] D. P. Palomar, "Convex Primal Decomposition for Multicarrier Linear MIMO Transceivers," *IEEE Trans. on Signal Processing*, vol. 53, no. 12, Dec. 2005.

[Yu'03] W. Yu, "A Dual Decomposition Approach to the Sum Power Gaussian Vector Multiple Access Channel Sum Capacity Problem," in *Proc. Conf. on Information and Systems (CISS)*, The Johns Hopkins Univ., March 12-14, 2003.