

Subgradients

Ryan Tibshirani
Convex Optimization 10-725

Last time: gradient descent

Consider the problem

$$\min_x f(x)$$

for f convex and differentiable, $\text{dom}(f) = \mathbb{R}^n$. **Gradient descent:**
choose initial $x^{(0)} \in \mathbb{R}^n$, repeat

$$x^{(k)} = x^{(k-1)} - t_k \cdot \nabla f(x^{(k-1)}), \quad k = 1, 2, 3, \dots$$

Step sizes t_k chosen to be fixed and small, or by backtracking line search

If ∇f is Lipschitz, gradient descent has convergence rate $O(1/\epsilon)$.
Downsides:

- Requires f differentiable
- Can be slow to converge

• Review.

(Gradient Descent: $x^{(k)} = x^{(k-1)} - t \nabla f(x)$

f : Diff, Convex, and Lipsitz conti with L , then for $t < \frac{1}{L}$,

$$|f(x^{(k)}) - f^*| < \frac{1}{2tk} \|x^{(0)} - x^*\|_2^2$$

error: $O(\frac{1}{k})$, step: $O(\frac{1}{\epsilon}) \leadsto \frac{1}{0.0001} : 10000 \text{ steps}$

(G.D with Strong Convexity

f : Diff, Strong Convex with m and Lipsitz conti with L , then for $t < \frac{2}{L+m}$,

$$|f(x^{(k)}) - f^*| < \gamma^n L \|x^{(0)} - x^*\|_2^2$$

error: $O(e^{-k})$, step: $O(\log \frac{1}{\epsilon}) \leadsto \log \frac{1}{0.0001} : 4 \text{ steps}$

$$mI < \nabla^2 f(x) < LI.$$

$$\Rightarrow h_1 < x^T \nabla^2 f(x) x < h_n$$

• but it assume ^{not} differentiable f .

\Rightarrow Cannot use First Order Condition
(Gradient Descent)

very important

Outline

Today: crucial mathematical underpinnings!

- Subgradients
- Examples
- Properties
- Optimality characterizations

Subgradients : 새로운 개념

: 미분이 불가능해나까 대신 사용.

Recall that for convex and differentiable f ,

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) \quad \text{for all } x, y$$

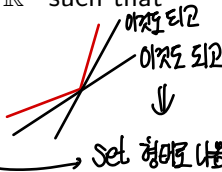
That is, linear approximation always underestimates f



A subgradient of a convex function f at x is any $g \in \mathbb{R}^n$ such that

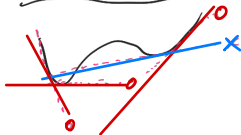
$$f(y) \geq f(x) + \boxed{g^T}(y - x) \quad \text{for all } y$$

$\nabla f(x)$ 를 대신하는 "기울기"



- Always exists¹
- If f differentiable at x , then $(g = \nabla f(x))$ uniquely
- Same definition works for nonconvex f (however, subgradients need not exist)

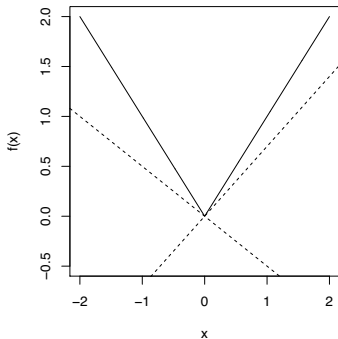
convex 아니어도
subgradient가 성립한다!!



¹On the relative interior of $\text{dom}(f)$

Examples of subgradients

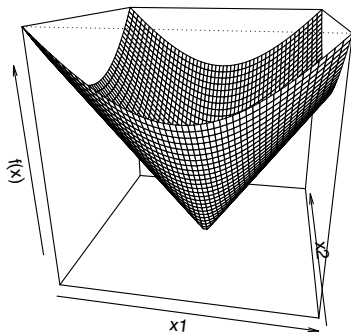
Consider $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = |x|$



- For $x \neq 0$, unique subgradient $g = \text{sign}(x)$
- For $x = 0$, subgradient g is any element of $[-1, 1]$

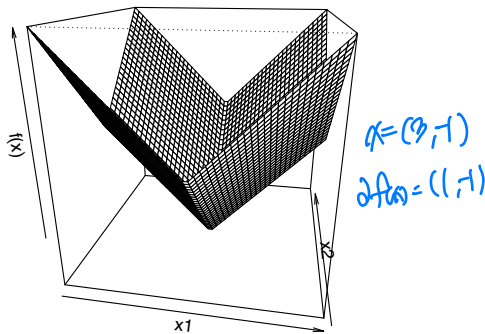
Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f(x) = \|x\|_2 = \sqrt{x^T x}$

$x^T x$ 는 항상 0,
 $\sqrt{x^T x}$ 는 $x \neq 0$ 에서 안됨.



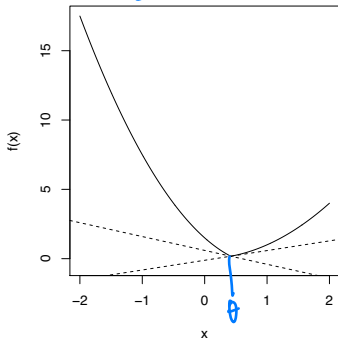
- For $x \neq 0$, unique subgradient $g = x / \|x\|_2$
- For $x = 0$, subgradient g is any element of $\{z : \|z\|_2 \leq 1\}$

Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f(x) = \|x\|_1$



- For $x_i \neq 0$, unique i th component $g_i = \text{sign}(x_i)$
- For $x_i = 0$, i th component g_i is any element of $[-1, 1]$

Consider $f(x) = \max\{f_1(x), f_2(x)\}$, for $f_1, f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ convex, differentiable



- For $f_1(x) > f_2(x)$, unique subgradient $g = \nabla f_1(x)$
- For $f_2(x) > f_1(x)$, unique subgradient $g = \nabla f_2(x)$
- For $f_1(x) = f_2(x)$, subgradient g is any point on line segment between $\nabla f_1(x)$ and $\nabla f_2(x)$

$$[\nabla f_1(0), \nabla f_2(0)]$$

$\nabla f(x)$: full gradient : normal (Subdifferential)

$\partial f(x)$: full subgradient
여러할 할 수 있는 집합 : partial.

subgradient
집합

Set of all subgradients of convex f ^{on x .} is called the subdifferential:

$$\partial f(x) = \{g \in \mathbb{R}^n : g \text{ is a subgradient of } f \text{ at } x\}$$

- Nonempty (only for convex f)
- $\partial f(x)$ is closed and convex (even for nonconvex f)
- If f is differentiable at x , then $\partial f(x) = \{\nabla f(x)\}$ (일일이 2개 이하)
- If $\partial f(x) = \{g\}$, then f is differentiable at x and $\nabla f(x) = g$
하나라면 다행히도.

Connection to convex geometry

Convex set $C \subseteq \mathbb{R}^n$, consider indicator function $I_C : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$I_C(x) = I\{x \in C\} = \begin{cases} 0 & \text{if } x \in C \\ \infty & \text{if } x \notin C \end{cases}$$

min $f(x) + I_C(x)$
 \downarrow
 call ∞ not feasible

For $x \in C$ ($\partial I_C(x) = \mathcal{N}_C(x)$), the **normal cone** of C at x is, recall

$$\mathcal{N}_C(x) = \{g \in \mathbb{R}^n : g^T x \geq g^T y \text{ for any } y \in C\}$$

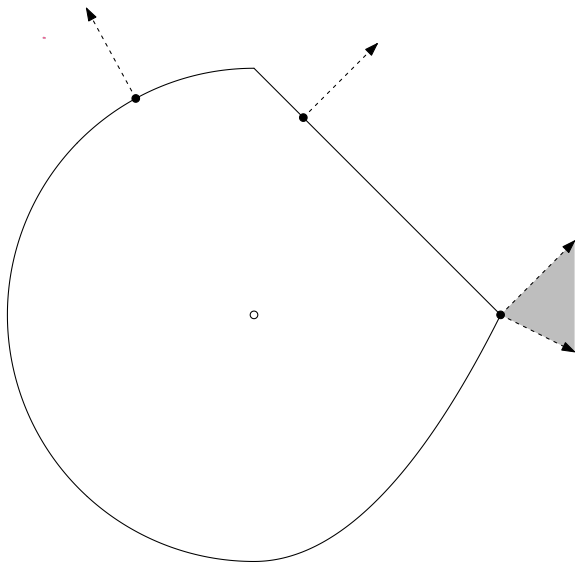
$0 \geq g^T(x-y) = \|g\| \|x-y\| \cos \theta$
 θ is the angle between g and $(x-y)$
 $\theta \leq 90^\circ$



Why? By definition of subgradient g ,

$$I_C(y) \geq I_C(x) + g^T(y - x) \quad \text{for all } y$$

- For $y \notin C$, $I_C(y) = \infty$
- For $y \in C$, this means $0 \geq g^T(y - x)$



Subgradient calculus

subgradient 계산에 쓰는 기법·규칙.

Basic rules for convex functions:

- Scaling: $\partial(a f) = a \cdot \partial f$ provided $a > 0$
- Addition: $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$
- Affine composition: if $g(x) = f(Ax + b)$, then

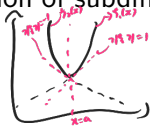
$$\partial g(x) = A^T \partial f(Ax + b)$$

- Finite pointwise maximum: if $f(x) = \max_{i=1, \dots, m} f_i(x)$, then

$$\partial f(x) = \text{conv} \left(\bigcup_{i: f_i(x) = f(x)} \partial f_i(x) \right)$$

convex hull

convex hull of union of subdifferentials of active functions at x



$$\partial f_1(a) = [-1, 1]$$

$$\partial f_2(a) = [-1, 1]$$

$$\partial f_1(a) \cup \partial f_2(a) = [-1, 1] \cup [-1, 1] = [-1, 1]$$

$$\text{conv}(\partial f_1(a) \cup \partial f_2(a)) = \text{conv}([-1, 1] \cup [-1, 1]) = [-1, 1]$$

- **General composition:** if

$$f(x) = h(g(x)) = h(g_1(x), \dots, g_k(x))$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$, $h : \mathbb{R}^k \rightarrow \mathbb{R}$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ h is convex and nondecreasing in each argument, g is convex, then

$$\partial f(x) \subseteq \left\{ p_1 q_1 + \dots + p_k q_k : \right. \\ \left. p \in \partial h(g(x)), q_i \in \partial g_i(x), i = 1, \dots, k \right\}$$

- **General pointwise maximum:** if $f(x) = \max_{s \in S} f_s(x)$, then

$$\partial f(x) \supseteq \text{cl} \left\{ \text{conv} \left(\bigcup_{s: f_s(x) = f(x)} \partial f_s(x) \right) \right\}$$

Under some regularity conditions (on S, f_s), we get equality

- **Norms:** important special case. To each norm $\|\cdot\|$, there is a **dual norm** $\|\cdot\|_*$ such that

$$\|x\| = \max_{\|z\|_* \leq 1} z^T x$$

≈ holder's inequality

(For example, $\|\cdot\|_p$ and $\|\cdot\|_q$ are dual when $1/p + 1/q = 1$.)
In fact, for $f(x) = \|x\|$ (and $f_z(x) = z^T x$), we get equality:

$$\partial f(x) = \text{cl} \left\{ \text{conv} \left(\bigcup_{z: f_z(x) = f(x)} \partial f_z(x) \right) \right\}$$

Note that $\partial f_z(x) = z$. And if z_1, z_2 each achieve the max at x , which means that $z_1^T x = z_2^T x = \|x\|$, then by linearity, so will $tz_1 + (1-t)z_2$ for any $t \in [0, 1]$. Thus

$$\partial f(x) = \text{argmax}_{\|z\|_* \leq 1} z^T x$$

Optimality condition

Recall, first order optimality condition
 For any f (convex or not), f is diff & convex $\Rightarrow f(x) \geq \nabla f(y)(x-y), \forall y$

$$f(x^*) = \min_x f(x) \iff 0 \in \partial f(x^*)$$

That is, x^* is a minimizer if and only if 0 is a subgradient of f at x^* . This is called the **subgradient optimality condition**

Why? Easy: $g = 0$ being a subgradient means that for all y

$$f(y) \geq f(x^*) + 0^T(y - x^*) = f(x^*)$$

Note the implication for a convex and differentiable function f ,
 with $\partial f(x) = \{\nabla f(x)\}$

convex & diff 함수 for 어떤 $f'=0$ 해서 찾아내, subgradient 하면 non-convex not diff 이셔도 가능.

Derivation of first-order optimality

Example of the power of subgradients: we can use what we have learned so far to derive the **first-order optimality condition**. Recall

$$\min_x f(x) \quad \text{subject to } x \in C$$

is solved at x , for f convex and differentiable, if and only if

$$\nabla f(x)^T (y - x) \geq 0 \quad \text{for all } y \in C$$

Intuitively: says that gradient increases as we move away from x .
How to prove it? First recast problem as

$$\min_x f(x) + I_C(x)$$

Now apply subgradient optimality: $0 \in \partial(f(x) + I_C(x))$



Observe

$$0 \in \partial(f(x) + I_C(x))$$

$$\iff 0 \in \{\nabla f(x)\} + \mathcal{N}_C(x)$$

$$\iff -\nabla f(x) \in \mathcal{N}_C(x)$$

$$\iff -\nabla f(x)^T x \geq -\nabla f(x)^T y \text{ for all } y \in C$$

$$\iff \nabla f(x)^T (y - x) \geq 0 \text{ for all } y \in C$$

as desired

Note: the condition $0 \in \partial f(x) + \mathcal{N}_C(x)$ is a **fully general** condition for optimality in convex problems. But it's not always easy to work with (KKT conditions, later, are easier)

Example: lasso optimality conditions

Given $y \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times p}$, **lasso** problem can be parametrized as

$$\min_{\beta} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

where $\lambda \geq 0$. Subgradient optimality: $= \frac{1}{2} \partial \|y - X\beta\|_2^2 + \lambda \partial \|\beta\|_1$

$$0 \in \partial \left(\frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right) \quad = -X^T(y - X\beta) + \lambda \partial \|\beta\|_1$$

$$\iff 0 \in -X^T(y - X\beta) + \lambda \partial \|\beta\|_1$$

$$\iff X^T(y - X\beta) = \lambda v$$

for some $v \in \partial \|\beta\|_1$, i.e.,

$$\beta = (-4, 3, 2, 0)$$

$$\downarrow$$

$$(-1, 1, 1, 1, 1)$$

$$v_i \in \begin{cases} \{1\} & \text{if } \beta_i > 0 \\ \{-1\} & \text{if } \beta_i < 0 \\ [-1, 1] & \text{if } \beta_i = 0 \end{cases}, \quad i = 1, \dots, p$$

$$\begin{aligned} \frac{\partial}{\partial x} f(g(x)) &= f'(g(x)) \cdot g'(x) \\ \text{if } f(x) = x^T a, &\rightarrow f'(x) = 2x \\ g(x) = y - X\beta &\rightarrow g'(x) = -X \\ f(g(x)) &= (y - X\beta)^T (y - X\beta) \\ \frac{\partial}{\partial x} f(g(x)) &= f'(g(x)) \cdot g'(x) \\ &= \end{aligned}$$

Penalized Regression.

· (LSE + 페널티)를 통해 다중공선성/오버피팅을 억제 하는 기법

$$\min \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \|\beta\|_p^p \quad (L_p\text{-Norm: } \|x\|_p = (\sum_{i=1}^n x_i^p)^{1/p}, x \in \mathbb{R}^n)$$

벡터의 크기를 재는 단위 ex) $(\frac{3}{5}) < (\frac{4}{5}) \Rightarrow \|(3,4)\|_2 = (3^2+4^2)^{1/2} = 5$
 $\|(5,5)\|_2 = (2^2+5^2)^{1/2} = \sqrt{29}$

$\|\cdot\|_p = (\sum (|x_i|)^p)^{1/p}$

일반적으로 $\begin{cases} p=2 & \|\beta\|_2 = (\sum |\beta_i|^2)^{1/2} \text{ (aka. Ridge Regression)} \\ p=1 & \|\beta\|_1 = (\sum |\beta_i|) \text{ (aka. LASSO Regression)} \end{cases}$ 을 사용.

· 다중공선성 & 오버피팅 으로 인한 모수 추정 오차 ($Var(\hat{\beta})$) 을 잡기 위한 알고리즘.

· 다중공선성 $\Rightarrow X^T X$ 행렬이 Singular 해지는 현상.

$$\begin{pmatrix} \Rightarrow \det(X^T X) = 0 \quad \text{하단 ~ 가 아님.} \\ \det(X^T X) = \prod_{i=1}^n \lambda_i(X^T X) \\ \Rightarrow \text{아이겐 벡터를 0 이 된 채 하는 것.} \end{pmatrix}$$

(Pseudo-singular. \Rightarrow 아이겐 벡터를 0 에 매우 근접한 값이 존재하는 것)

$$X_1 = X_2 + X_3 + \epsilon$$

4	1	3
3	2	1
2	1	1
5	4	1

2대 문제? $X^T X = T \Lambda T^T \Rightarrow (X^T X)^{-1} = T \Lambda^{-1} T^T$

$$\Lambda^{-1} = \begin{pmatrix} \lambda_1^{-1} & & \\ & \ddots & \\ & & \lambda_n^{-1} \end{pmatrix} \Rightarrow \text{Singular 행렬 } \lambda_n^{-1} = 0^{-1} \text{ 이므로. 아래 불가.}$$

Pseudo singular 행렬 λ_n^{-1} 가 매우 큰 값이 되므로 $(X^T X)^{-1}$ 의 원소들이 매우 크게 됨.
 λ_n^{-1} 는 λ_n 이 조금만 변해도 매우 크게 변하므로 $(X^T X)^{-1}$ 이 매우 불안정 해짐.

(ex) $\lambda_n \approx 0.0001 \Rightarrow \lambda_n^{-1} = 10000$
 $\lambda_{n-1} = 0.00005 \Rightarrow \lambda_{n-1}^{-1} = 20000$ ($Var(\hat{\beta})$ 가 커짐)

오버피팅 $MSE(\hat{\beta}) = E((Y - X\hat{\beta})^2) = \underbrace{E((Y - X\beta)^2)}_{\text{실제 데이터}} + \underbrace{E((X\beta - X\hat{\beta})^2)}_{\text{모델}}$

$$\approx \text{Bias}(\beta) + \text{Var}(\hat{\beta})$$



$$\underbrace{(X, Y)}_{\text{실제 데이터}} \Leftarrow \underbrace{X\beta}_{\text{모델 (가설)}} \Leftarrow \underbrace{X\hat{\beta}}_{\text{구현된 모델}}$$

Bias는 가설이 실재를 얼마나 잘 표현할 수 있는지
 Variance는 가설을 얼마나 제대로 구현할 수 있는지

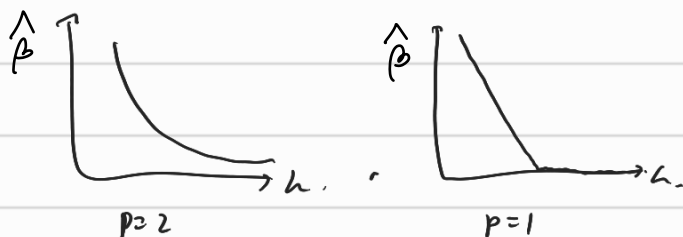
$$\min \|Y - X\beta\|_2^2 + \lambda \|\beta\|_p^p$$

$$\|Y - X\beta\|_2^2 \approx \text{minimize} \Rightarrow \beta = (X^T X)^{-1} X^T Y$$

$$\|\beta\|_p \approx \text{minimize} \Rightarrow \beta = 0$$

\approx λ 의 크기가 커짐에 따라 $\|\beta\|_p$ 의 중요도가 커지며

$$\beta = (X^T X)^{-1} X^T Y \longrightarrow 0 \text{ 으로 바뀐다.}$$



여기서 중요한 특징이.

- Ridge의 경우 필연적으로 β_i 를 0에 근접하게 만들지만 완전히 0으로 만들진 못한다.
- LASSO의 경우 필연적으로 β_i 를 이제 0으로 만들.

실제 solution을 살펴보자

Ridge: $\min (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta$

$$\begin{aligned} \nabla \mathcal{L}(\beta) &= X^T (Y - X\beta) + 2\lambda \beta \\ &= X^T Y - X^T X \beta + 2\lambda \beta \end{aligned}$$

$$= X^T Y - (X^T X + 2\lambda I) \beta = 0. \quad (F, 0.0)$$

$$X^T Y = (X^T X + 2\lambda I) \beta$$

$$\Rightarrow \beta^* = (X^T X + 2\lambda I)^{-1} X^T Y$$

(CS, 이를 Tikhonov - Regularized Inverse 라고 한다)

→ 2개의 reg의
closed-form solution

Write X_1, \dots, X_p for columns of X . Then our condition reads:

$$\begin{cases} X_i^T(y - X\beta) = \lambda \cdot \text{sign}(\beta_i) & \text{if } \beta_i \neq 0 \\ |X_i^T(y - X\beta)| \leq \lambda & \text{if } \beta_i = 0 \end{cases}$$

$$\Rightarrow |X_i^T \hat{e}| \Rightarrow X_i \perp \hat{e}$$

Note: subgradient optimality conditions don't lead to closed-form expression for a lasso solution ... however they do provide a way to **check lasso optimality**

They are also helpful in understanding the lasso estimator; e.g., if $|X_i^T(y - X\beta)| < \lambda$, then $\beta_i = 0$ (used by screening rules, later?)

X_i 가 필요없다.

$X_i \perp Y \Rightarrow$ 예측 선택하는 X_i 의 정보력이 없다.

Example: soft-thresholding

Simplified lasso problem with $X = I$: $n \times p \rightarrow n \times n$

$$\min_{\beta} \frac{1}{2} \|y - \beta\|_2^2 + \lambda \|\beta\|_1$$

This we can solve directly using subgradient optimality. Solution is $\beta = S_{\lambda}(y)$, where S_{λ} is the **soft-thresholding operator**:

$$[S_{\lambda}(y)]_i = \begin{cases} y_i - \lambda & \text{if } y_i > \lambda \\ 0 & \text{if } -\lambda \leq y_i \leq \lambda, \quad i = 1, \dots, n \\ y_i + \lambda & \text{if } y_i < -\lambda \end{cases}$$

Check: from last slide, subgradient optimality conditions are

$$\begin{cases} y_i - \beta_i = \lambda \cdot \text{sign}(\beta_i) & \text{if } \beta_i \neq 0 \\ |y_i - \beta_i| \leq \lambda & \text{if } \beta_i = 0 \end{cases}$$

문제 : y 와 3차원 가짜곡면에서 너무 크면 β 를
잡는 것. 직접적인 사용처라기 보다는 보조적으로
사용되는 알고리즘.

(만약 $\min \frac{1}{2} \|y - \beta\|_2^2$ 의 경우 $y = \beta$ 로 두는 것이 최적
즉 $\min \| \beta \|_1$ 의 경우 $\beta = 0$ 역시 최적
따라서 β 는 $0 \sim y$ 사이의 어떤 값이 될 것

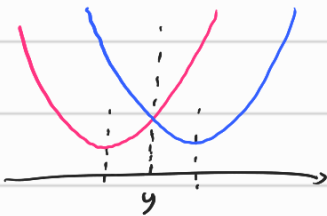
ex) $y \in \mathbb{R}^1$ 이라고 가정.

$$f(\beta) = \frac{1}{2} (y - \beta)^2 + h|\beta|$$

$$\approx (y - \beta)^2 + 2h|\beta|$$

$$= \beta^2 - 2y\beta + y^2 + 2h|\beta|$$

$$= \begin{cases} \text{Case 1. } \beta \geq 0 \Rightarrow \beta^2 - 2(y-h)\beta + y^2 \Rightarrow (\beta^* = (y-h) \text{ 일때 } y^2 - (y-h)^2) \\ \text{Case 2. } \beta < 0 \Rightarrow \beta^2 - 2(y+h)\beta + y^2 \Rightarrow (\beta^* = (y+h) \text{ 일때 } y^2 - (y+h)^2) \end{cases}$$



이 두 값과 0의 관계에 따라서 최소값이 결정.

⇒ 규칙을 따 내기에는 too 복잡 (cf, 식별-2, soft thresholding)

⇒ Subgradient Optimality Condition을 쓰자.

$$f(\beta) = \frac{1}{2} \|y - \beta\|_2^2 + h\|\beta\|_1$$

$$\partial f(\beta) \supset \partial \left(\frac{1}{2} \|y - \beta\|_2^2 + h\|\beta\|_1 \right)$$

$$= \partial \left(\frac{1}{2} \|y - \beta\|_2^2 \right) + h \partial (\|\beta\|_1) \quad (\text{by Additivity and Multiplicity})$$

$$= -(y - \beta) + h \partial (\|\beta\|_1) \quad (\partial f = \nabla f \text{ if } \nabla f \text{ exists})$$

$$\left(\begin{array}{l} \text{Let } f(v) = v^T v, \quad (\nabla f(v) = 2v), \text{ then} \\ \frac{\partial}{\partial \beta} \|y - \beta\|_2^2 = \frac{\partial}{\partial \beta} f(y - \beta) = \frac{\partial f(y - \beta)}{\partial (y - \beta)} \frac{\partial (y - \beta)}{\partial \beta} \\ = (y - \beta)^T \cdot (-1) \\ = -(y - \beta) \end{array} \right)$$

$$= \begin{cases} -(y - \beta) + h & \text{if } \beta > 0 \quad \text{Case 1.} \\ -(y - \beta) + h \cdot v \quad \{v \in [-1, 1]\} & \text{if } \beta = 0 \quad \text{Case 2.} \\ -(y - \beta) - h & \text{if } \beta < 0 \quad \text{Case 3.} \end{cases}$$

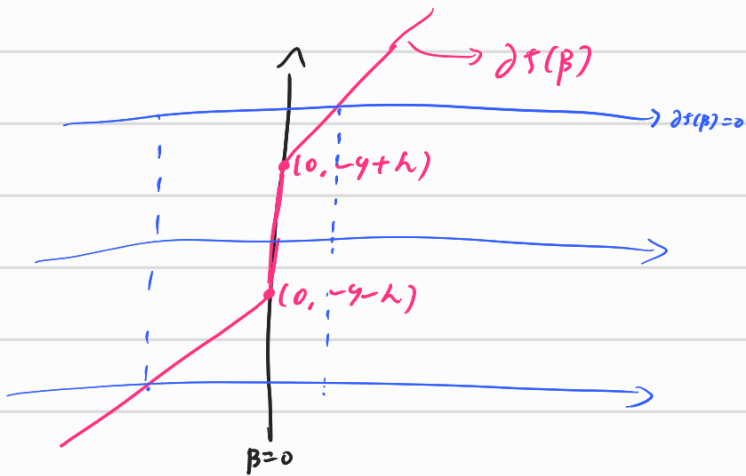
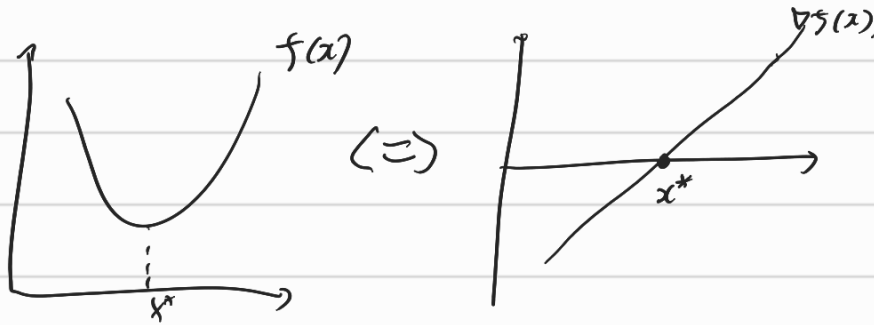
If β^* satisfies $0 \in \partial f(\beta^*)$, β^* is global minimum.

Case 1. $0 = -(y - \beta) + h \Rightarrow \beta = y - h$ if $\beta = y - h > 0$ {i.e. $y > h$ }

Case 2. $0 \in -(y - \beta) + h \cdot v \Rightarrow \beta \in h \cdot v - y = [-h - y, h - y]$ if $\beta = 0$ {i.e. $|y| < h$ }

Case 3. $0 = -(y - \beta) - h \Rightarrow \beta = y + h$ if $\beta = y + h < 0$ {i.e. $y < -h$ }

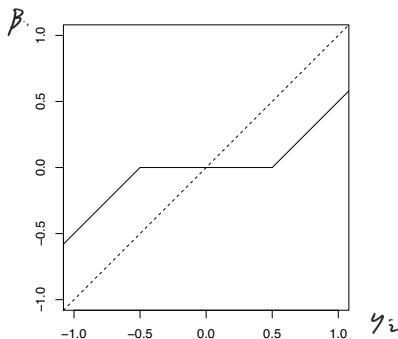
Remark that to find critical point (극점) of $f(x)$, we have to find a point $(x; \nabla f(x) = 0)$ (i.e. a point that $f'(x)$ and x -axis ($y=0$ -line) meet)



Now plug in $\beta = S_\lambda(y)$ and check these are satisfied:

- When $y_i > \lambda$, $\beta_i = y_i - \lambda > 0$, so $y_i - \beta_i = \lambda = \lambda \cdot 1$
- When $y_i < -\lambda$, argument is similar
- When $|y_i| \leq \lambda$, $\beta_i = 0$, and $|y_i - \beta_i| = |y_i| \leq \lambda$

Soft-thresholding in
one variable:



Example: distance to a convex set

Recall the **distance function** to a closed, convex set C :

$$\text{dist}(x, C) = \min_{y \in C} \|y - x\|_2$$

This is a convex function. What are its subgradients?

Write $\text{dist}(x, C) = \|x - P_C(x)\|_2$, where $P_C(x)$ is the projection of x onto C . It turns out that when $\text{dist}(x, C) > 0$,

$$\partial \text{dist}(x, C) = \left\{ \frac{x - P_C(x)}{\|x - P_C(x)\|_2} \right\}$$

Only has one element, so in fact $\text{dist}(x, C)$ is differentiable and this is its gradient

We will only show one direction, i.e., that

$$\frac{x - P_C(x)}{\|x - P_C(x)\|_2} \in \partial \text{dist}(x, C)$$

Write $u = P_C(x)$. Then by first-order optimality conditions for a projection,

$$(x - u)^T(y - u) \leq 0 \quad \text{for all } y \in C$$

Hence

$$C \subseteq H = \{y : (x - u)^T(y - u) \leq 0\}$$

Claim:

$$\text{dist}(y, C) \geq \frac{(x - u)^T(y - u)}{\|x - u\|_2} \quad \text{for all } y$$

Check: first, for $y \in H$, the right-hand side is ≤ 0

Now for $y \notin H$, we have $(x - u)^T(y - u) = \|x - u\|_2 \|y - u\|_2 \cos \theta$ where θ is the angle between $x - u$ and $y - u$. Thus

$$\frac{(x - u)^T(y - u)}{\|x - u\|_2} = \|y - u\|_2 \cos \theta = \text{dist}(y, H) \leq \text{dist}(y, C)$$

as desired

Using the claim, we have for any y

$$\begin{aligned} \text{dist}(y, C) &\geq \frac{(x - u)^T(y - x + x - u)}{\|x - u\|_2} \\ &= \|x - u\|_2 + \left(\frac{x - u}{\|x - u\|_2} \right)^T (y - x) \end{aligned}$$

Hence $g = (x - u)/\|x - u\|_2$ is a subgradient of $\text{dist}(x, C)$ at x

References and further reading

- S. Boyd, Lecture notes for EE 264B, Stanford University, Spring 2010-2011
- R. T. Rockafellar (1970), “Convex analysis”, Chapters 23–25
- L. Vandenberghe, Lecture notes for EE 236C, UCLA, Spring 2011-2012

