# Chapter 8. Gradient Descent for Constrained & non-smooth optimization.

## 1. Constrained Optimization
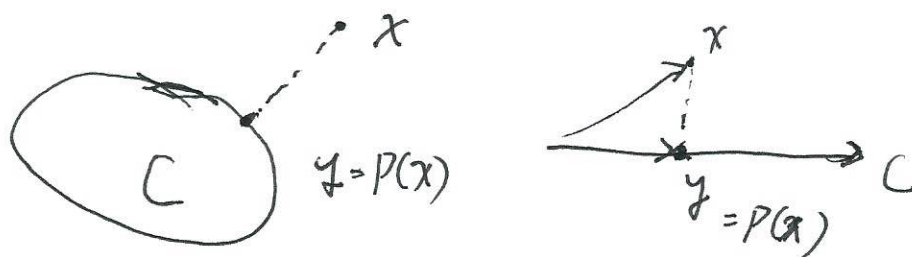
$$\min_{x} f(x) \quad st \quad x \in C.$$

$C$: Convex, set

Recall gradient descent: $x^{k+1} \leftarrow x^k - \eta^k \cdot \nabla f(x^k)$

$\uparrow$
may not $\in C$
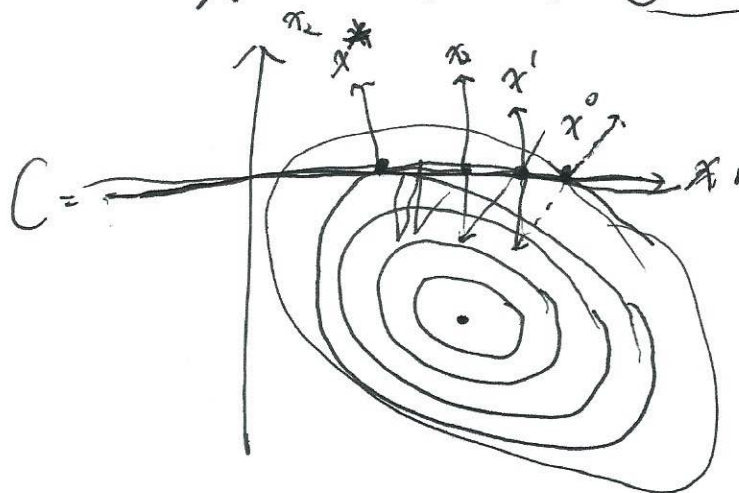
Def: Projection of $x$ onto $C$

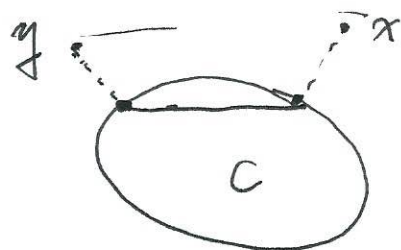$$P_C(x) = \operatorname*{argmin}_{y \in C} \| x - y \|_2$$



$y = P(x)$        $y = P(x)$

Projected Gradient Descent.

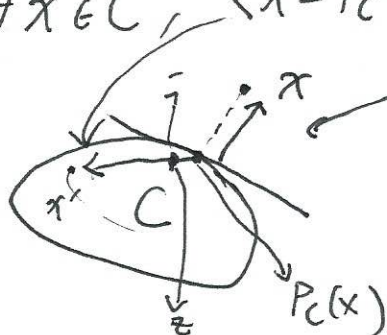$$- \quad x^{k+1} = P_C \left( x^k - \eta^k \cdot \nabla f(x^k) \right)$$

$$C: \{ x \mid x_2 = 0 \}$$

Lemma. $\|P_C(x) - P_C(y)\| \le \|x - y\| \quad \forall x, y$



pf: ① Step 1: $\forall x' \in C, \quad \langle x' - P_C(x), x - P_C(x)\rangle \le 0$



( If this not true,

$\exists x', \text{ st } \langle \qquad \rangle > \bigcirc > 0$

$z \ne P_C(x), \quad \|z - x\|^2 = \|z - P_C(x) + P_C(x) - x\|^2$

$\qquad = \|z - P_C(x)\|^2 + 2\langle z - P_C(x), P_C(x) - x\rangle$

$\qquad\qquad + \|P_C(x) - x\|^2$

take $z = P_C(x) + \varepsilon \cdot (x' - x)$

$\|z - x\|^2 = \text{const.} \; \varepsilon^2 - \text{const.} \; \varepsilon$

$\qquad\qquad + \|P_C(x) - x\|^2$

$\varepsilon \to 0 \quad, \quad \|z - x\|^2 < \|P_C(x) - x\|^2$

$\qquad\qquad\qquad \text{contradict} \quad )$

② $\langle P_C(y) - P_C(x), x - P_C(x)\rangle \leq 0$

$+ \langle P_C(x) - P_C(y), y - P_C(y)\rangle \leq 0$

$\langle P_C(y) - P_C(x), x - P_C(x) - y + P_C(y)\rangle \leq 0$

$\|P_C(x) - P_C(y)\|^2 = \langle P_C(x) - P_C(y), P_C(x) - P_C(y)\rangle \leq \langle P_C(y) - P_C(x), y - x\rangle$

$\|P_C(x) - P_C(y)\|^2 \leq \langle P_C(y) - P_C(x), y - x\rangle \leq \|P_C(y) - P_C(x)\|\|y - x\|$

$\Rightarrow \|P_C(x) - P_C(y)\| \leq \|y - x\|$ #

---

Convergce.: If $f$ is convex, bnded below,

continuosly differentable,

Then $\{x^k\}$ generated by PGD

with line search ( backtracking line search)

$\lim_{k \to \infty} x^k = x^*$

---

Example:

$$\min_{x} \frac{1}{2} x^T Q x + b^T x \quad \text{st} \quad x \geq 0$$

PGD: For $k = 0, 1, \cdots$
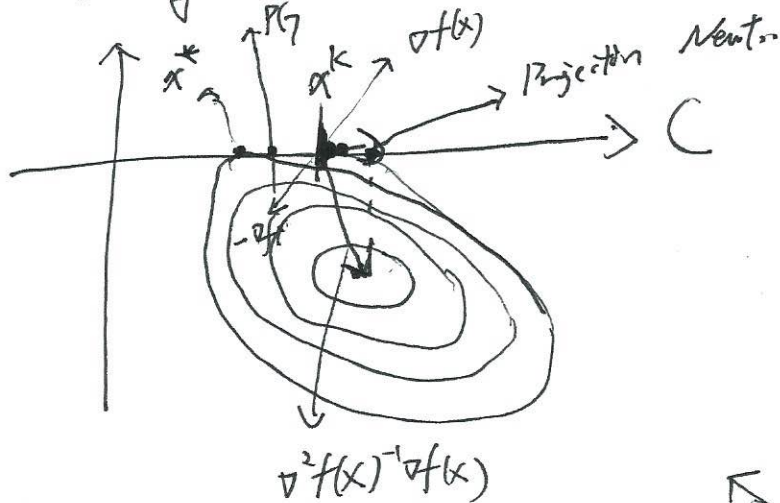
$$g = Q x^k + b$$

$$\bar{x} = x - \eta^k g$$

$$x_i^{k+1} = \begin{cases} \bar{x}_i & \text{if } \bar{x}_i \geq 0 \\ 0 & \text{if } \bar{x}_i < 0 \end{cases} = P_C(\bar{x})$$

end.

---

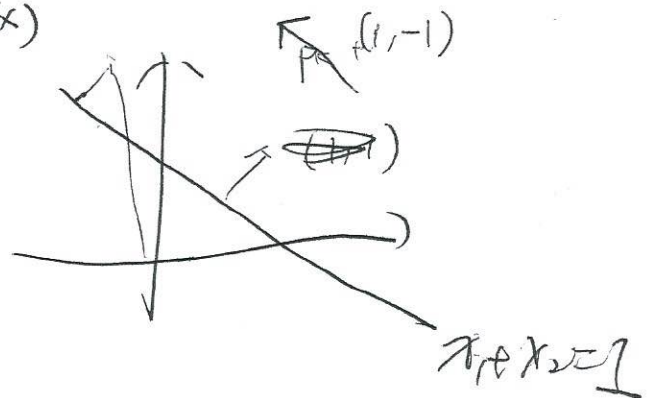Can we do projected Newton?

$$x^{k+1} = P_C \left( x^k - \eta^k \cdot \nabla^2 f(x^k)^{-1} \nabla f(x^k) \right) \quad \times$$



$$\nabla^2 f(x)^{-1} \nabla f(x)$$

$$P_A f(1,-1)$$

$$p^T x \cdot \left( \frac{p}{\|p\|^2} \right)$$
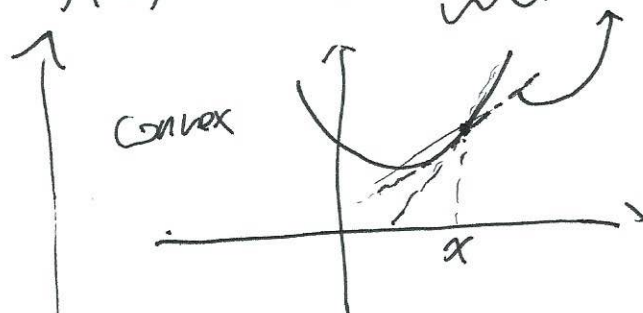
$$x_1 + x_2 = 1$$

# Gradient Descent for Non-smooth Functions.
### Non-(differentiable)

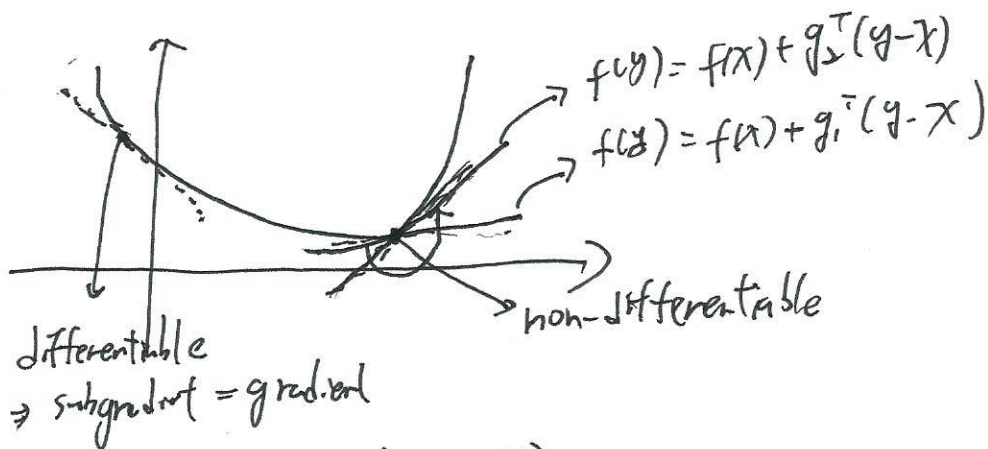— Subgradient (for non-differentiable function)

recall : $\nabla f(x)$ for a convex function

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) \quad \forall y$$

Convex

$x$

Def : $g \in R^n$ is a subgradient for $f(\cdot)$ on $X$
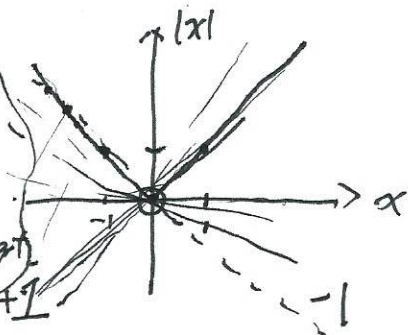
iff

$$f(y) \geq f(x) + g^T(y-x) \quad \forall y$$

$$f(y) = f(x) + g_2^T(y-x)$$
$$f(y) = f(x) + g_1^T(y-x)$$

differentiable
$\Rightarrow$ subgradient = gradient

non-differentiable

Example : Lasso : $f(x) = |x|$  $x \in R$

If $x > 0 \Rightarrow f(x) = x$
$$\nabla f(x) = 1$$

if $x < 0 \Rightarrow f(x) = -x$
$$\nabla f(x) = -1$$

if $x = 0$
$$f(y) \geq f(0) + g^T(y-0)$$
$$= g^T y$$
$$\forall g \in [-1, 1] \text{ satfy}$$

if $x \neq 0$
$\quad \nabla f(x) = \text{sign}(x)$

if $x = 0$
$\quad [-1, 1]$
$\quad$ can be subgradient

$|x|$

$+1$

$x$
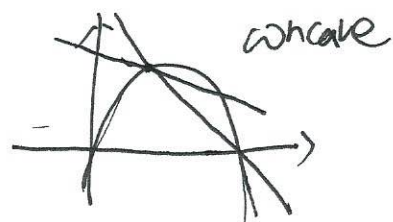
$-1$

Def: Subdifferential

$$\partial f(x) = \{ g : g \text{ is a subgradient of } f \text{ at } x \}$$

① Existence:

$\partial f(x)$ is nonempty if $f$ is convex

② If $f$ is differentiable, then

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

concave

③ If $\partial f(x) = \{ \nabla f(x) \}$, then $f$ is differentiable at $x$.

Example: If $C$ is a convex set
we can define $I_C(x)$ (indicator function)

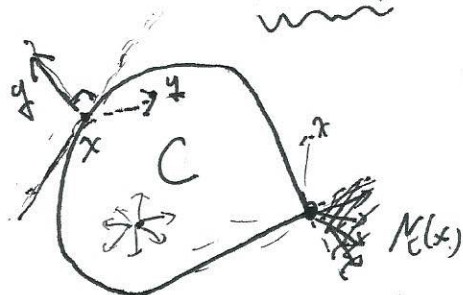$$I_C(x) = \begin{cases} 0 & \text{if } x \in C \\ \infty & \text{if } x \notin C \end{cases} \checkmark$$

$$\left( \begin{array}{l} \min_x f(x) \text{ st } x \in C \\ \underset{\text{equivalent}}{\Longleftrightarrow} \quad \min_x f(x) + I_C(x) \end{array} \right)$$

Property: $\partial I_C(x) := N_C(x)$ (normal vectors), $x \in C$

$$N_C(x) = \{ g : g^T(y - x) \leq 0, \ \forall y \in C \}$$

Why? By definition, we want.

$$I_C(y) \geq I_C(x) + g^T(y-x) \ \forall y$$

If $y \notin C \Rightarrow I_C(y) = \infty$ ok

If $y \in C \Rightarrow I_C(y) = 0$

$\rightarrow 0 \geq 0 + g^T(y-x) \Rightarrow g^T(y-x) \leq 0$

want $\rightarrow$ $\Leftrightarrow N_C(x)$

$C$ $N_C(x)$

if $x \in C$
$I_C(x) = 0$

Rec Optimality condition

recall if $f$ is differentiable

$x^*$ is optimal iff $\nabla f(x^*) = 0$

For non-differentiable function: (convex function)

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

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

For constrained optimization:
$\qquad\qquad\qquad\qquad\qquad\qquad$ $\nearrow$ differentiable

$$\min_x f(x) \quad s.t \quad x \in C \quad , \quad (f, C \text{ are convex})$$

$\boxed{x^* \text{ is optimal solution iff } \nabla f(x^*)^T(y - x) \geq 0 \quad \forall y \in C}$

$\downarrow$ equivalent

why? $\min_x \{f(x) + I_C(x)\} = F(x)$

$x^*$ is optimal $\Longrightarrow$ ~~$\nabla f(x^*)$~~ $0 \in \partial F(x^*)$
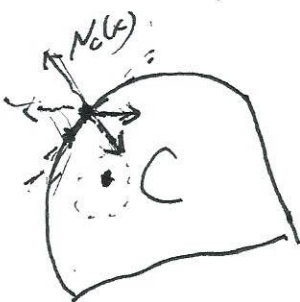
$0 \in \partial F(x^*) = \{\nabla f(x^*)\} + \partial I_C(x^*)$

$\qquad\qquad\qquad = \{\nabla f(x^*)\} + N_C(x^*)$

$N_C(x) = \{g : g^T(y - x) \leq 0 , \forall y \in C\}$

$-\nabla f(x^*)^T(y - x) \leq 0 \quad \forall y \in C$

$\Longrightarrow \nabla f(x^*)^T(y - x) \geq 0 \quad \forall y \in C$



$\#$

Example:

$$\min_{x} \left\{ g(x) + \lambda \|x\|_1 \right\} := f(x)$$
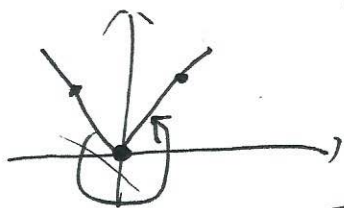
Lasso $g(x) = \frac{1}{2}\|Ax - b\|_2^2$

$$\|x\|_1 := \sum_i |x_i|$$

$\{x_i\} + \{y_1, y_2, \dots, y_n\}$
$= \{x_1 + y_1, x_1 + y_2 \dots \}$

$x^*$ Optimal solution $\Rightarrow$ $0 \in \partial f(x^*)$

$= \{\nabla g(x^*)\} + \lambda \cdot \partial(\|x^*\|_1)$

$= \nabla g(x^*) + \lambda V$

if $a \in \mathbb{R}$

$$\partial |a| = \begin{cases} 1 & \text{if } a > 0 \\ -1 & \text{if } a < 0 \\ [-1, 1] & \text{if } a = 0 \end{cases}$$

$$V_i = \begin{cases} 1 & \text{if } x_i^* > 0 \\ -1 & \text{if } x_i^* < 0 \\ [-1, 1] & \text{if } x_i^* = 0 \end{cases}$$



$\Rightarrow x^*$ is optimal iff

$$\begin{cases} \nabla_i g(x^*) = -\lambda & \text{if } x_i^* > 0 \\ \nabla_i g(x^*) = \lambda & \text{if } x_i^* < 0 \\ \nabla_i g(x^*) \in [-\lambda, \lambda] & \text{if } x_i^* = 0 \end{cases}$$

Simple Example: (Soft - thresholding)

$(Ax-b)$

$$\min_{x} \frac{1}{2}\|x - b\|_2^2 + \lambda \|x\|_1$$

$\nabla g(x^*) = x^* - b$

Optimal solution $x^*$ iff

$$\begin{cases} x_i^* - b_i = -\lambda & \text{if } x_i^* > 0 \\ x_i^* - b_i = \lambda & \text{if } x_i^* < 0 \\ x_i^* - b_i \in [-\lambda, \lambda] & \text{if } x_i^* = 0 \end{cases}$$
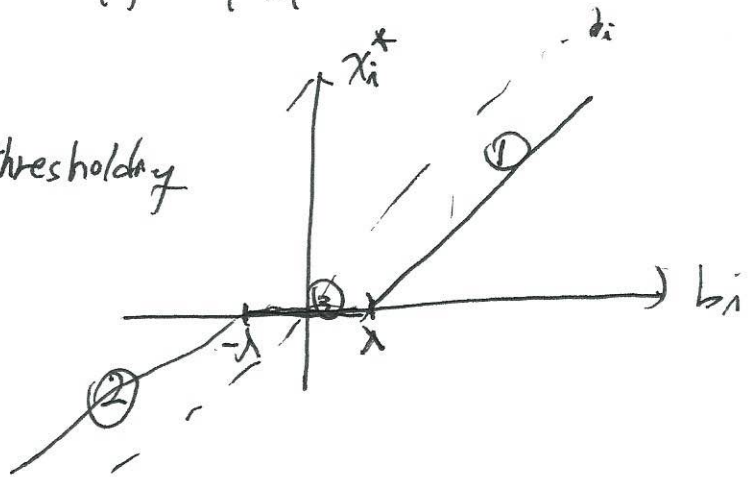
$\{ \text{if } b_i^* > 0$

$$\Rightarrow \begin{cases} \text{if } b_i - \lambda > 0, & x_i^* = b_i - \lambda \\ \text{if } b_i + \lambda < 0, & x_i^* = b_i + \lambda \\ \text{if } |b_i| < \lambda, & x_i^* = 0 \end{cases}$$

$$\Rightarrow \begin{cases} x_i^* = b_i - \lambda & \text{if } b_i > \lambda & \textcircled{1} \\ x_i^* = b_i + \lambda & \text{if } b_i < -\lambda & \textcircled{2} \\ x_i^* = 0 & \text{if } |b_i| \leq \lambda \vee \textcircled{3} \end{cases} := S_\lambda(b)$$

$\downarrow$

we call this soft-thresholding

$$x^* = S_\lambda(b)$$



---

How to minimize non-differentiable functions?

$$\min_x f(x)$$

Subgradient Descent:

$$x^{k+1} \leftarrow x^k - \eta^k \cdot g^k$$

$$g^k \in \partial f(x^k)$$

This will converge to $x^*$, but slow.

# Proximal Gradient Descent for Decomposible Functions:

- Decomposible function:

$$f(x) = g(x) + h(x).$$

$g$: is convex & differentiable

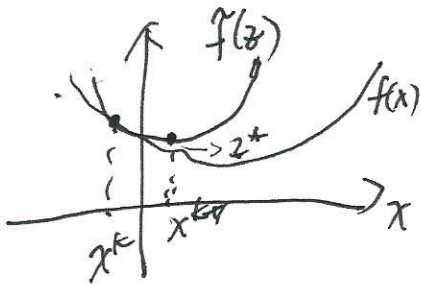$h$: is convex, but may be non-differentiable.

Example: $\frac{1}{2}\|Ax - b\|^2 + \lambda\|x\|_1$

$\underbrace{\phantom{\frac{1}{2}\|Ax-b\|^2}}_{g(x)}$  $\underbrace{\phantom{\lambda\|x\|_1}}_{h(x)}$

For ML problems: $g(\cdot)$: loss function

$h(\cdot)$: regularisation.

---

Recall gradient descent. $x^{k+1} \leftarrow x^k - \nabla f(x^k)$

Form approximate function:



$$\tilde{f}(z) = f(x^k) + \nabla f(x^k)^T (z - x^k) + \frac{1}{2\eta}\|z - x^k\|^2$$

$z^* = \arg\min_z \tilde{f}(z)$

$x^{k+1} = z^*$

$\nabla\tilde{f}(z^*) = 0$
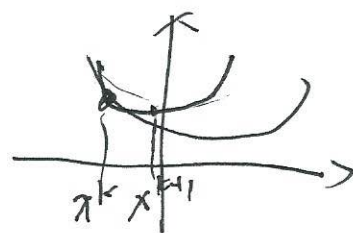
$\nabla f(x^k) + \frac{1}{\eta}(z^* - x^k) = 0$

$\Rightarrow z^* = x^k - \eta \cdot \nabla f(x^k)$

proximal gradient descent.

$$\min_{x} \{ g(x) + h(x) \} := f(x)$$

approximate function

$$\tilde{f}(z) = \tilde{g}(z) + h(z)$$

$$= \tilde{g}(x^k) + \nabla g(x^k)^T(z - x^k) + \frac{1}{2\eta} \| z - x^k \|^2 + h(z)$$

$$\frac{1}{2\eta} \| z - x^k \|^2$$
$$+ (z - x^k)^T \nabla g(x^k)$$
$$+ \eta^2 \| \nabla g(x^k) \|^2 \cdot \frac{1}{2\eta}$$

$$= \frac{1}{2\eta} \| z - x^k + \eta \cdot \nabla g(x^k) \|^2 + h(z) + \text{constants}$$

$$= \frac{1}{2\eta} \| z - (x^k - \eta \nabla g(x^k)) \|^2 + h(z) + \text{const}$$

$$z^* = \arg\min_{z} \tilde{f}(z)$$

$$x^{k+1} = z^*$$

close to $x$     to have small $h(z)$

"proximal operator"

$$\text{prox}_{\eta}(x) = \arg\min_{z} \frac{1}{2\eta} \cdot \| z - x \|^2 + h(z)$$

proximal gradient descent:

$$x^{k+1} = \text{prox}_{\eta^k}(x^k - \eta^k \cdot \nabla g(x^k))$$
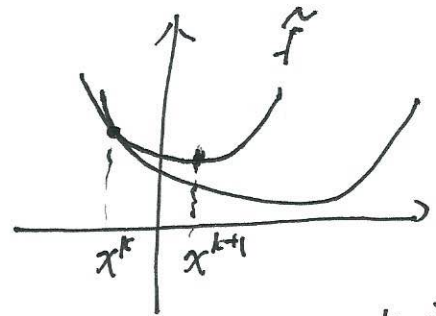
$$\left( = P_C(x^k - \eta^k \cdot \nabla g(x^k)) \right)$$

Decomposable functions

$$\min_{x} \left\{ g(x) + h(x) \right\} := f(x)$$

Smooth      non-differentiable
loss function    $(\lambda \|x\|_1)$
$(\|Ax - L\|_2^2)$    $= \lambda \sum_i |x_i|$

## Proximal Gradient Descent

At each iteration



$$\hat{f}(y) = \tilde{g}(y) + h(y)$$

$$= g(x^k) + \nabla g(x^k)^T (y - x^k) + \frac{1}{2\eta} \|y - x^k\|^2 + h(y)$$

$$x^{k+1} = \arg\min_{y} \hat{f}(y_{..})$$

$$= \arg\min_{y} \left[ \frac{1}{2\eta} \|y - x^k + \eta \cdot \nabla g(x^k)\|^2 + h(y) \right] + \text{const}$$

$$\frac{1}{2\eta} \|y - x^k\|^2 + \nabla g(x^k)^T (y - x^k) + \frac{1}{2\eta} \cdot \|\eta \nabla g(x^k)\|^2$$

$$= \arg\min_{y} \frac{1}{2} \|y - (x^k - \eta \nabla g(x^k))\|^2 + \eta h(y)$$

$$= \text{prox}_{\eta} \left( x^k - \eta \nabla g(x^k) \right)$$

$$\left( \text{prox}_{\eta}(x) = \arg\min_{y} \frac{1}{2} \|y - x\|^2 + \eta h(y) \right)$$

when $h(x) = \lambda \|x\|_1$

$\text{prox}_h(x) = \arg\min_y \frac{1}{2} \|y - x\|^2 + \lambda \eta \cdot \|y\|_1$

$$= S_{\eta\lambda}(x) = \begin{cases} x - \eta\lambda & \text{if } x > \eta\lambda \\ x + \eta\lambda & \text{if } x < -\eta\lambda \\ 0 & \text{if } |x| < \eta\lambda \end{cases}$$

proximal gradient ~~for~~ when $h(x) = \lambda \|x\|_1$

For $k = 0, 1, \cdots$

$$z = x^k - \eta^k \nabla g(x^k)$$

$$x^{k+1} = S_{\eta^k \lambda}(z)$$

end .

( Iterative Soft Thesholding Algorithm )

ISTA

---

Convergence :

if $\eta^k = 1/k$

✗ converge to solution

① will converge to stationary points if

$$0 < \underline{\eta} \leq \eta^k \leq 1/L \quad \forall k$$

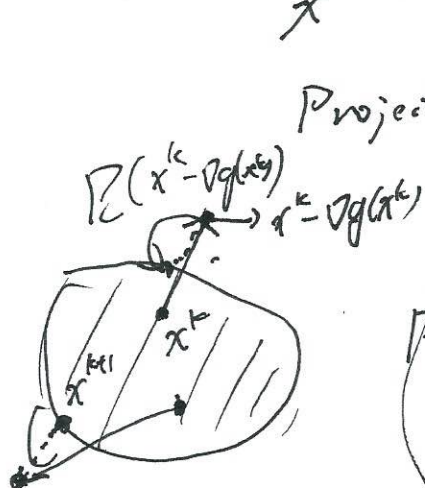$$( \nabla^2 g(x) \leq LI \quad \forall x )$$

② If $g(\cdot)$ is strongly convex

$$( mI \leq \nabla^2 g(x) \leq LI \quad \forall x )$$

$\Rightarrow$ Linear convergence .

Recall : projeted gradient for
Constrained minimization

$$\min_{x} f(x) \text{ st } x \in C, \text{ C is convex set.}$$



Projected Gradient (PG) :

$$x^{k+1} = P_C (x^k - \eta \cdot \nabla f(x^k))$$

$P_C(z)$: projection of $z$ on $C$.

$$P_C(z) := \operatorname*{argmin}_{y} \|y - z\|_2 \text{ st } y \in C$$

PG is a special case for proximal gradient.

Why?    equivalent to

$$\operatorname*{argmin}_{x} f(x) + I_C(x) \quad ---- (*)$$

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

solve $(*)$ by proximal gradient

$$x^{k+1} = \operatorname{prox}_{\eta} (x^k - \eta \nabla f(x^k))$$

$$\operatorname{prox}_{\eta}(z) := \operatorname*{argmin}_{y} \|y - z\|^2 + \eta I_C(y)$$

$$= \operatorname*{argmin}_{y} \|y - z\|^2 \text{ st } y \in C$$

$$= P_C(z)$$

# Proximal Newton:

**Newton:** $x^{k+1} \leftarrow x^k - \eta^k \cdot \nabla^2 f(x^k)^{-1} \nabla f(x^k)$

**another:** At each iteration, form approximate function          Hessian

$$\tilde{f}(y) = f(x^k) + \nabla f(x^k)^T(y - x^k) + \frac{1}{2}(y-x^k)^T \frac{\nabla^2 f(x^k)}{\eta}(y-x^k)$$

$\underset{y}{\arg\min}\ \tilde{f}(y)$

$$\nabla f(x^k) + \frac{1}{\eta}\nabla^2 f(x^k)(y - x^k) = 0$$

$$\Rightarrow y = x^k$$

$$0 = \eta \nabla f(x^k) + \nabla^2 f(x^k)(y - x^k)$$

$$\Rightarrow y = x^k - \eta \nabla^2 f(x^k)^{-1} \nabla f(x^k)$$

**proximal Newton:** for decomposable function

$$\min_{x} \left\{ g(x) + h(x) \right\} := f(x)$$

differentiable          non-differentiable

At each iteration

$$\tilde{f}(y) = \tilde{g}(y) + h(y)$$

$$= g(x^k) + \nabla g(x^k)^T(y - x^k) + \frac{1}{2}(y-x^k)^T \nabla^2 g(x^k)(y-x^k) + h(y)$$

$$x^{k+1} = \underset{y}{\arg\min}\ \tilde{f}(y)$$

$$\left(H = \sigma^2 g(x^k)\right) \qquad \operatorname*{argmin}_{y} f(y)$$

$$= \operatorname*{argmin}_{y} \frac{1}{2}\left\| y - \left(x^k - H_o^{-1}g(x^k)\right)\right\|_H^2 + h(y)$$

$$= \operatorname*{prox}^H \left(x^k - H^{-1}g(x^k)\right)$$

$$\boxed{\|x\|_H^2 = x^T H x}$$

$$\left(y - \left(x^k - H^{-1}g(x^k)\right)\right)^T H \left(y - \left(x^k - H^{-1}g(x^k)\right)\right)$$

$$= (y - x^k)^T H (y - x^k) + 2(y - x^k)^T H \cdot H^{-1} g(x^k)$$