# Computing proximal maps and projections

**＊ Computing proximal maps:**

Let $h(z) = I_\chi(z) = \begin{cases} \infty & z \notin \chi \\ 0 & z \in \chi \end{cases}$

$\therefore \quad P^* = \begin{bmatrix} \nabla f_0(x) \\ 1 \\ x \in \chi \end{bmatrix} = \nabla f_0(x) + \underbrace{I_\chi(z)}_{h(z)}$

$\{ [\cdot], D\}$

So any general convex optimization problem can be solved by

taking $\quad x_{k+1} = \underset{s_k \nabla f_0(x_k)}{\text{prox}} (x_k - s_k \nabla f_0(x_k)) = [x_k - s_k \nabla f_0(x_k)]_\chi \quad \{\text{Euclidean projection of } x_k - s_k \nabla f_0(x_k) \text{ on } \chi \text{ will give the next point}\}$

$\text{since } \text{prox}_{\frac{s}{\chi}}(x) = \underset{z}{\text{argmin}} \left( s \left[ h(z) + \frac{1}{2} \| z - x \|_2^2 \right] \right) = \underset{z}{\text{argmin}} \begin{pmatrix} \frac{1}{2} \| z - x \|_2^2, & \text{if } z \in \chi \\ \infty, & \text{if } z \notin \chi \end{pmatrix} = \underset{z \in \chi}{\text{argmin}} \left( \frac{1}{2} \| z - x \|_2^2 \right) = \underset{z \in \chi}{\text{argmin}} \left( \| z - x \|_2 \right) = [x]_\chi \quad \{\text{Euclidean projection of } x \text{ on } \chi\}$

$\boxed{\text{prox}_{\frac{s}{\chi}}(x) = [x]_\chi} \quad \{\text{even } s \text{ is not present inside the expression}\}$

*will result to infeasibility not valid case as we are assuming the optimization problem is feasible.*

**＊ Projection onto a halfspace:**

$\chi = \{ x : a^T x \leq b \}, \, a \neq 0$

$\therefore \quad \text{prox}_{I_\chi}(x) = \underset{z \in \chi}{\text{argmin}} \| z - x \|_2^2 = [x]_\chi \quad$ now clearly if $x \in \chi$ the $[x]_\chi = x$ (কারণ $x, \chi$ মধ্যে $x$ এর থেকে কাছে projection আর কি আছে)

now if $x \notin \chi$ then $[x]_\chi$ will be the projection of $x$ on the hyperplane $\{z : a^T z = b\}$

[eq: Euclidean_proj_on_hyperplane] (করা আছে)

$[x]_{(a^T \square = b)} = x - \dfrac{a^T x - b}{\| a \|_2^2} a$

so for $x \notin \chi, \quad [x]_\chi = x - \dfrac{a^T x - b}{\| a \|_2^2} a$

$\therefore \quad [x]_{(a^T \square \leq b)} = \begin{cases} x, & \text{if } a^T x \leq b \\ x - \dfrac{a^T x - b}{\| a \|_2^2} a, & \text{if } a^T x > b \end{cases} \quad$ Ans.

**＊ Projection onto the positive orthant:**

$\chi = \mathbb{R}_+^n = \{ x : \mathbb{R}^n : x \geq 0 \}$

$\therefore \quad [x]_\chi = \underset{z \in \chi}{\text{argmin}} \| z - x \|_2 = \underset{z \geq 0}{\text{argmin}} \sum_{i=1}^n (z_i - x_i)^2 = [x]_+ \quad$ [ eq: Projection on positive orthant ]

*This can be solved by inspection*
*if $x_i \geq 0$, we take $z_i = x_i, z_i \geq 0$*
*then $x_i < 0$, we cannot take $z_i = x_i$, as it would violate $z \geq 0$*
*so the closest we can get to $z_i < 0$, is if we take $z_i = 0$*
*combining both: $z_i^* = \max(0, x_i)$ which reflects the decision above*
*in vector notation: $z^* = \max(0, x) = [x]_+$, where max is intended elementwise.*

**＊ Projections onto the standard simplex:** [Projection onto the standard simplex]

$\chi = \{ x \in \mathbb{R}^n : x \geq 0, \, \mathbf{1}^T x = 1 \}$

$[x]_\chi = \underset{z}{\text{argmin}} \left[ \frac{1}{2} \| z - x \|_2^2 \mid z \geq 0, \, \mathbf{1}^T z = 1 \right] = z^*$

$\begin{cases} \nabla \frac{1}{2} \| z - x \|_2^2 \\ 1 \\ z \geq 0 \\ \mathbf{1}^T z = 1 \end{cases} = \begin{pmatrix} \nabla \frac{1}{2} \| z - x \|_2^2 + \lambda \nabla(z \geq 0) \\ 1 \\ \mathbf{1}^T z = 1 \end{pmatrix} \quad$ partial lagrangian  $L(z, \nu) = \frac{1}{2} \| z - x \|_2^2 + \nu(\mathbf{1}^T z - 1) + \lambda_0$ // remember: for convex optimization problem the dual problem // associated with a partial lagrangian will result in strong duality // if slaters condition holds

$\therefore g(\nu) = \underset{z \geq 0}{\inf} L(z, \nu)$ // Full lagrangian তো এই $\inf_z$ has no constraints

*যেহেতু dual problemটা রান করছি, $g(\nu)$ এ $\text{argmin } z^*(\nu)$ পেয়ে যাব, যেটা this is also the solution of $P_z L(z, \nu)$*

$= \underset{z \geq 0}{\inf} \left( \frac{1}{2} \| z - x \|_2^2 + \nu(\mathbf{1}^T z - 1) \right)$

$= \underset{z \geq 0}{\inf} \left( \dfrac{\sum_{i=1}^n (z_i - x_i)^2}{2} + \nu \left( \sum_{i=1}^n z_i \right) - \nu \right)$

$= \underset{z \geq 0}{\inf} \left( \frac{1}{2} \sum_{i=1}^n (z_i - x_i)^2 + \nu \left( \sum_{i=1}^n z_i \right) - \nu \right)$

$= \underset{z \geq 0}{\inf} \left[ \sum_{i=1}^n \left( \frac{1}{2} (z_i - x_i)^2 + \nu z_i \right) - \nu \right]$

$= -\nu + \underset{z \geq 0}{\inf} \left[ \left( \sum_{i=1}^n \frac{1}{2} (z_i - x_i)^2 + \nu z_i \right) \right]$

*note that the objective is separable in each $z_i$*

$= -\nu + \sum_{i=1}^n \left( \underset{z_i \geq 0}{\inf} \left( \frac{1}{2} (z_i - x_i)^2 + \nu z_i \right) \right)$ // want to simplify this in $z_i$

$\frac{1}{2}(z_i - x_i)^2 + \nu z_i = \frac{1}{2}(z_i^2 + x_i^2 - 2 z_i x_i) + \nu z_i = \frac{x_i^2}{2} + \frac{1}{2}(z_i^2 - 2 x_i z_i + 2 \nu z_i) = \frac{x_i^2}{2} + \frac{1}{2}(z_i^2 - 2(x_i - \nu) z_i + (x_i - \nu)^2) - \frac{1}{2}(x_i - \nu)^2$

$\begin{pmatrix} \frac{1}{2}(z_i - x_i)^2 + \nu z_i \\ 1 \\ z_i \geq 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(z_i - (x_i - \nu))^2 \\ 1 \\ z_i \geq 0 \end{pmatrix}$

$= \frac{1}{2}\left( z_i - (x_i - \nu) \right)^2 + \underbrace{\frac{1}{2}(x_i^2 - (x_i - \nu)^2)}_{\text{constant w.r.t } z}$

*now if $x_i - \nu \geq 0$ then $z_i^*(\nu) = x_i - \nu \geq 0$ ( objective squared $\leq$ minimum পাবে at $z_i = x_i - \nu \geq 0$ which will make it zero)*

*if $x_i - \nu < 0$ then $z_i^*(\nu) = 0$ ( কারণ $z_i^* = x_i - \nu < 0$ হলে $z_i \geq 0$ violated হবে, in that case the square term will be minimized by setting $z_i = 0$, which will minimize the distance between $\{z_i : z_i \geq 0\}$ and $(x_i - \nu)$.)*

$\nabla$ combining both: $\forall_i \quad z_i^*(\nu) = \max\{0, x_i - \nu\}$   // max is taken elementwise

$\therefore z^*(\nu) = \left( \max\{0, x_i - \nu\} \right) = \max\{0, x - \nu \mathbf{1}\}$

$= \frac{1}{2}\left( \| z^*(\nu) - x \|_2^2 \right) + \nu \left( \mathbf{1}^T z^*(\nu) - 1 \right)$ // remember though $g(\nu)$ is a weird looking function, by construction it is concave (see boyd)

// now $R g(\nu)$ solve করে which is an unconstrained optimization problem in $\nu$

// $\nu^*$ পেলে, then $z^* = [x]_\chi = z^*(\nu^*)$ // By KKT condition $\nabla_z L(z, \nu, \lambda) = 0$ solve করে // আবার $x(\nu, \lambda) \triangleright$ optimal $\nu^*, \lambda^*$ বসালে আমরা primal // minimizer $x^*$ পেয়ে যাব।

Now, by KKT condition:

• $\nabla_z L(z, \nu) = 0 \Rightarrow z^*(\nu) = \max\{0, x - \nu\}$

• primal feasibility: $\mathbf{1}^T z(\nu) = 1$

$\quad \overset{?}{=} \mathbf{1}^T \max\{0, x - \nu\} = \sum \max\{0, x - \nu\} = 1$ // এখানে $\nu^*$ পেয়ে যাব।   then $\boxed{z^*(x) = \max\{0, x - \nu^*\}}$

given

a scalar equation

• dual feasibility: নেই

• complementary slackness: নেই

$\nabla$ Similarly when $\chi = \{ x \mid x \geq 0, \mathbf{1}^T x = b \}$ then $z^* = [x]_\chi = \max\{0, x - \nu^*\}$

$\sum_{i=1}^n \max\{0, x_i - \nu\} = b$

**涁 Projection onto the Euclidean ball:**

$\chi = \{ x \in \mathbb{R}^n : \| x \|_2 \leq 1 \}$

$[x]_\chi = \underset{z \in \chi}{\text{argmin}} \| z - x \|_2 = \underset{z \in \chi}{\text{argmin}} \frac{1}{2} \| z - x \|_2^2$

optimization problem $\| x \|_2 \leq 1$

optimization problem $\min_{z} \|z-x\|_2 \leq 1$

$\begin{cases} \forall \; \frac{1}{2}\|z-x\|_2^2 \\ 1 \end{cases} \|z\|_2 \leq 1$

$\|z\|_2^2 \leq 1$

• if $z \in X$, then $z^* = x$ with 0 objective value
  which is a norm so obviously that is $z^* = x$ is an argmin. (a norm can have 0 as a absolutely value)

• if $x \notin X \Rightarrow \|x\|_2 > 1$

$\rightarrow L(z,\lambda) = \frac{1}{2}\|z-x\|_2^2 + \lambda\left(\|z\|_2^2 - 1\right)$

• $(\nabla_z) \rightarrow (z-x) + \lambda 2z = 0 \rightarrow (1+2\lambda)z = x \rightarrow z = \frac{x}{1+2\lambda}$

KKT
• $\lambda \geq 0$
• $\lambda\left(\|z\|_2^2 - 1\right) = 0$
• $\|z\|_2^2 \leq 1$

$\frac{\|x\|_2^2}{(1+2\lambda)^2} \leq 1 \rightarrow \|x\|_2^2 \leq (1+2\lambda)^2 \rightarrow 1 < \|x\|_2^2 \leq (1+2\lambda)^2$
$\rightarrow (1+2\lambda)^2 > 1 \rightarrow \lambda > 0$

$\|z\|_2^2 - 1 = 0 \rightarrow \|z\|_2^2 = 1 \quad \left\{ z = \frac{x}{1+2\lambda}\right\}$

$\rightarrow \frac{\|x\|_2^2}{(1+2\lambda)^2} = 1 \rightarrow 1+2\lambda = \|x\|_2 \rightarrow \lambda = \frac{\|x\|_2 - 1}{2} = \lambda^\#$

$z^* = \frac{x}{1+2\lambda^\#} = \frac{x}{(1+2\cdot\frac{\|x\|_2-1}{2})\cdot} = \frac{x}{\|x\|_2}, \quad [z^*=\cdots]$

$$[x]_X = \begin{cases} x, & \text{if } \|x\|_2 \leq 1 \\ \dfrac{x}{\|x\|_2}, & \text{if } \|x\|_2 > 1 \end{cases}$$

---

**# Projection onto the $l_1$ norm ball:**

$X = \{x \in \mathbb{R}^n : \|x\|_1 \leq 1\}$

$[x]_X = \underset{\|z\|_1 \leq 1}{\text{argmin}} \; \frac{1}{2}\|z-x\|_2^2$

↳ (eq: Optimization Problem for l-1 norm projection)

$\begin{cases} \forall \; \frac{1}{2}\|z-x\|_2^2 \\ 1 \\ \|z\|_1 \leq 1 \end{cases} \leadsto L(z,\lambda) = \frac{1}{2}\|z-x\|_2^2 + \lambda\left(\|z\|_1 - 1\right)$  // $l_1$ norm এর square এর মানে নাই তাই।

$= \left(\frac{1}{2}\sum_{i=1}^n (z_i-x_i)^2 + \lambda\sum_{i=1}^n |z_i|\right) - \lambda$  // $\lambda$ is a constant w.r.t $z$

$= \left(\sum_{i=1}^n \left[\frac{1}{2}(z_i-x_i)^2 + \lambda|z_i|\right]\right) - \lambda$   

$\underset{z}{\text{argmin}}\; L(z,\lambda) = \underset{z}{\text{argmin}}\left(\frac{1}{2}\|z-x\|_2^2 + \lambda\|z\|_1 - \lambda\right) = \left[\text{shr}_\lambda(z_i)\right]_{i=1}^n$

(constant সরা গেলে argmin same থাকে)

Note that this is what we are trying to find argmin of with \lambda = s_k

$g(\lambda) = \inf_z\left(\left(\sum_{i=1}^n \left[\frac{1}{2}(z_i-x_i)^2 + \lambda|z_i|\right]\right) - \lambda\right)$

$= -\lambda + \inf_z \sum_{i=1}^n \left[\frac{1}{2}(z_i-x_i)^2 + \lambda|z_i|\right]$

note that this is separable in $z_i$

$= -\lambda + \sum_{i=1}^n \inf_{z_i}\left(\frac{1}{2}(z_i-x_i)^2 + \lambda|z_i|\right) = -\lambda + \sum_{i=1}^n \inf_{z_i} \phi(z_i;\lambda)$  $\left[\phi(y;\lambda) = \frac{1}{2}(y-x_i)^2 + \lambda|y|\right]$

$\phi(z_i;\lambda)$

$z_i^*(\lambda) = \underset{z_i}{\text{argmin}}\left(\frac{1}{2}(z_i-x_i)^2 + \lambda|z_i|\right) = \underset{z_i}{\text{argmin}}\; \phi(z_i;\lambda)$  [eq: z_i^**(\lambda) Orginal]

We will use the identity: $|z| = \underset{|\rho|\leq 1}{\max}\; \rho z$

e.g. $|-3| = \underset{-1\leq\rho\leq 1}{\max}(-3\rho) = \max[-3,3] = 3$ so it works!

$|3| = \underset{-1\leq\rho\leq 1}{\max}(3\rho) = \max[-3,3] = 3$

$\forall \; \frac{1}{2}(z_i-x_i)^2 + \lambda|z_i|) = \forall\left(\frac{1}{2}(z_i-x_i)^2 + \lambda \underset{|\rho_i|\leq 1}{\max}\rho_i z_i\right) = \forall_{z_i}\left(\frac{1}{2}(z_i-x_i)^2 + \underset{|\rho_i|\leq 1}{\max}\lambda\rho_i z_i\right) = \forall_{z_i}\underset{|\rho_i|\leq 1}{\max}\left(\frac{1}{2}(z_i-x_i)^2 + \lambda\rho_i z_i\right) = \underset{z_i}{\wedge}\underset{|\rho_i|\leq 1}{\vee}\left(\frac{1}{2}(z_i-x_i)^2 + \lambda\rho_i z_i\right)$

(constant w.r.t $\rho_i$, so সরতে পারে ভিতরে)

$\rho_i \in \mathbb{R} \mapsto \{D,[\varepsilon]\}$

**# Sion's minimax theorem**

$\begin{cases}\forall_{\rho_i} & \phi(\cdot,\rho_i) = \Box + \rho_i\cdot\Box = \Box\text{affine}\} \\ \phi(z_i;\cdot) \end{cases}\begin{cases}\forall_{\rho_i} & \phi(z_i,\cdot) = \frac{1}{2}(z_i-\Box)^2 + \Box z_i = D \;[\because\text{affine}+ D=D]\end{cases}$

$\{D,\Omega\}$

$\Rightarrow$ so $\phi$ is [concave (affine), $\cup$] in $\rho_i$ the maximizing variable

• • $[D, \cup]$  ↳ $z_i$ : minimizing "  
maximizing set is compact  

$\Rightarrow$ Sion's minimax theorem will hold.

**# Sion's minimax theorem** [function of x (y fixed)]

$\left(X\{D,[\varepsilon]\},\subseteq\mathbb{R}^n\}, Y\{D,\subseteq\mathbb{R}^m\}, \phi:\{\cdot\}: X\times Y\rightarrow\mathbb{R}, \begin{array}{c}\forall\\y\in Y\end{array} \phi(\cdot,y)\{D_x\}_\cup, \begin{array}{c}\vee\\x\in X\end{array}\phi(x,\cdot)\{\Box_y\}_\cap\right) \Rightarrow \sup_{y\in Y}\min_{x\in X}\phi(x,y) = \min_{x\in X}\sup_{y\in Y}\phi(x,y)$

minimizer variable / function of y (x fixed)

$\left(X\{D,\subseteq\mathbb{R}^n\}, Y\{D,[\varepsilon],\subseteq\mathbb{R}^m\}, \phi:\{\cdot\}: X\times Y\rightarrow\mathbb{R}, \begin{array}{c}\vee\\y\in Y\end{array}\phi(\cdot,y)\{D_x\}_\cup, \begin{array}{c}\vee\\x\in X\end{array}\phi(x,\cdot)\{\Box_y\}_\cap\right) \Rightarrow \max_{y\in Y}\inf_{x\in X}\phi(x,y) = \inf_{x\in X}\max_{y\in Y}\phi(x,y)$  [eq: Sion's mini-max theorem]

maximizer variable / minimizer variable

// mnemonic:
// minimizer variable> functional convex  
maximizer variable> functional concave হবে  
// যে: one of the set compact হবে Sion's minimax theorem ধরবে।

$= \underset{|\rho_i|\leq 1}{\wedge}\left(\underset{z_i}{\vee}\left(\frac{1}{2}(z_i^2+x_i^2-2z_i x_i) + \lambda\rho_i z_i\right)\right)$ (constant w.r.t $z_i$)

$\frac{1}{2}(z_i^2 + x_i^2 - 2z_i x_i + 2\lambda\rho_i z_i)$

$= \frac{1}{2}\left(z_i^2 - 2(x_i-\lambda\rho_i)z_i + (x_i-\lambda\rho_i)^2\right) + \frac{x_i^2}{2} - \frac{1}{2}(x_i-\lambda\rho_i)^2$

$= \frac{1}{2}\left(z_i - (x_i-\lambda\rho_i)\right)^2 + \frac{x_i^2}{2} - \frac{1}{2}(x_i-\lambda\rho_i)^2$

function to be $\vee$  (constant w.r.t $z_i$)

so $z_i$ minimum achieve করার জন্য  
$z_i^* = z_i^\#(x_i,\lambda) = x_i - \lambda\rho_i$ with 0 norm$^2$ value [eq: z^*(\lambda)]

optimum objective

$= \underset{|\rho_i|\leq 1}{\wedge}\left[\frac{x_i^2}{2} - \frac{1}{2}(x_i-\lambda\rho_i)^2\right]$

$\frac{x_i^2}{2} - \frac{1}{2}(x_i^2 + \lambda^2\rho_i^2 - 2\lambda x_i\rho_i)$

$= -\frac{1}{2}\lambda^2\rho_i^2 + \lambda x_i\rho_i = \lambda\left(x_i\rho_i - \frac{1}{2}\lambda\rho_i^2\right)$

$= \lambda\underset{|\rho_i|\leq 1}{\wedge}\left(x_i\rho_i - \frac{1}{2}\lambda\rho_i^2\right)$

$\text{# } \rho_i = -\nabla = \frac{1}{2}\rightarrow -\frac{1}{2}\lambda\left(\rho_i^2 - \frac{2x_i}{\lambda}\rho_i + \left(\frac{x_i}{\lambda}\right)^2 - \frac{x_i^2}{\lambda^2}\right) = -\frac{1}{2}\lambda\left(\rho_i - \frac{x_i}{\lambda}\right)^2 + \frac{1}{2}\cdot\frac{x_i^2}{\lambda^2}$

$= -\lambda\underset{|\rho_i|\leq 1}{\vee}\left(\frac{1}{2}\lambda\left(\rho_i - \frac{x_i}{\lambda}\right)^2 - \frac{1}{2}\frac{x_i^2}{\lambda^2}\right)$

$= -\lambda\left[-\frac{1}{2}\frac{x_i^2}{\lambda^2} + \frac{\lambda}{2}\underset{|\rho_i|\leq 1}{\vee}\left(\rho_i - \frac{x_i}{\lambda}\right)^2\right]$

$\left[\frac{|x_i|}{\lambda}\in\pi:\{\Box:|\Box|\leq 1\}\right]$

$\left[\frac{|x_i|}{\lambda}\leq 1\right]\rho_i^*(\lambda) = \frac{x_i}{\lambda}$ with 0 minvalue

$\left[\frac{x_i}{\lambda} > 1\right]\rho_i^*(\lambda) = 1 \quad \begin{array}{c}+1\\-1\end{array}\Big\}\rho_i$

$\left[-\frac{x_i}{\lambda} > 1\right]\rho_i^*(\lambda) = -1$  
$\frac{|x_i/\lambda|<1}{}$

↳ in this cases we can impose the following formula:  
$\rho_i^*(\lambda) = \text{sgn}(x_i)$ *** as when $x_i > \lambda > 0$, $\rho_i^*(\lambda) = \text{sgn}(x_i) = 1$  
" $x_i < -\lambda < 0$, $\rho_i^*(\lambda) = \text{sgn}(x_i) = -1$ }

$\begin{array}{c}x_i/\lambda\\-1\end{array}\Big\}$ min distance

$\begin{array}{c}+1\\-1\end{array}\Big\}\rho_i$

$x_i/\lambda$

combining all of these  

$\underset{|\rho_i|\leq 1}{\vee}\left[-\left(\rho_i - \frac{x_i}{\lambda}\right)^2\right] = \rho_i^*(\lambda) = \begin{cases} x_i/\lambda, & \text{if } |x_i|\leq\lambda \\ \text{sgn}(x_i), & \text{else}\end{cases}$

Then, from [eq: z^**(\lambda)], $\rho_i^*(\lambda)$ এর value বসালে পরে, $z_i^\#(\lambda) = x_i - \rho_i^*(\lambda)\lambda = x_i - \lambda\begin{cases}x_i/\lambda, & \text{if } |x_i|\leq\lambda \\ \text{sgn}(x_i), & \text{else}\end{cases}$

$= \begin{cases} 0, & \text{if } |x_i|\leq\lambda \\ x_i - \lambda\text{sgn}(x_i), & \text{else}\end{cases}$

So, now we have $z_i^\#(\lambda)$  

$z^\#(\lambda) = [z_i^\#(\lambda)]_{i=1}^n$ is $\vee$ of the Lagrangian, giving us one piece of the KKT puzzle.

$z_i^\#(\lambda) = \begin{cases} 0, & \text{if } |x_i|\leq\lambda \\ x_i - \lambda\text{sgn}(x_i), & \text{else}\end{cases} = \text{shr}_\lambda(x_i)$

$\text{shr}_\lambda(x_i) = \text{sgn}(x_i)\left[|x_i| - \lambda\right]_+$

this is called soft threshold function

**# By Sion's mini-max theorem** [eq: Sion's mini-max theorem] $\underset{z_i}{\wedge}\underset{|\rho_i|\leq 1}{\vee}\left(\frac{1}{2}(z_i-x_i)^2 + \lambda\rho_i z_i\right) = \underset{|\rho_i|\leq 1}{\vee}\underset{z_i}{\wedge}\left(\frac{1}{2}(z_i-x_i)^2 + \lambda\rho_i z_i\right)$

and by Sion's minimax theorem $z_i^\#(\lambda)$ will also be the minimizer of the Lagrangian as in [eq: z_i^**(\lambda) Orginal] thus giving us information regarding one equation (vanishing (sub) gradient of the Lagrangian) of the KKT condition

$\rho_i^* = \begin{cases}x_i/\lambda, & \text{if } |x_i|\leq\lambda \\ \text{sgn}(x_i), & \text{else}\end{cases}$

$\tilde{z}_i(\lambda) = (z_i^*(\lambda))_{i=1}^n$ is ... of the Lagrangian, giving us one piece of the KKT puzzle.

this is called soft threshold function / shrinkage operator

[eq: Optimization Problem for l-1 norm projection] (पहले वाला) ... , strong duality holds as $z=0$ is strictly feasible point & the optimal value is finite because the set is bounded.

Lets try to understand how $sthr_\lambda(x_i)$ works.
Need to finish

now:
$$\nabla_{z_i}\left[\tfrac{1}{2}(z_i-x_i)^2+\lambda|z_i|\right]=\nabla_{z_i,|r_i|\leq 1}\left[\tfrac{1}{2}(z_i-x_i)^2+\lambda r_i z_i\right]=\lambda_{|r_i|\leq 1}\left[\nabla_{z_i}\tfrac{1}{2}(z_i-x_i)^2+\lambda r_i z_i\right]$$

$$z_i^*(\lambda)=\nabla_{z_i}\left[\tfrac{1}{2}(z_i-x_i)^2+\lambda|z_i|\right]=(z_i^*(r_i^*),\lambda)=\begin{cases}0,&\text{if }|x_i|\leq\lambda\\ x_i-\lambda\,sgn(x_i),&\text{else}\end{cases}$$

[eq: z**(\lambda)] (दूसरा)

A better explanation needed

...बाकी (rest) ... optimal $\lambda$ ... from the KKT condition.

Rest of the KKT condition are:
Primal feasibility: $\|z^*(\lambda)\|_1\leq 1$
dual feasibility: $\lambda\geq 0$
complementary slackness: $\lambda\left(\|z^*(\lambda)\|_1-1\right)=0$

now: if $\lambda=0$ $\Rightarrow$ $\left(\|z^*(\lambda)\|_1\leq 1\text{ is an inactive constraint}\right)$ $z^*(\lambda)=[z_i^*(r_i^*)]_{i=1}^n[x_i]_{i=1}^n=x$

$[x]_X=z_i^*(\lambda)=x\Rightarrow\|x\|_1\leq 1$ [गर] then $[x]_X=x$. $x\in X$

$\therefore\lambda=0\to[x]_X=x$

यानि:
$[x]_X=x=z^*(\lambda)=[z_i^*(\lambda)]=[\{...\}]_{i=1}^n$ if $\lambda=0$

$\lambda=0\leftrightarrow\|x\|_1\leq 1$

$\bullet$ Now consider, then $\lambda>0\to\|z^*(\lambda)\|_1-1=0\leftrightarrow\left(\sum_{i=1}^n max\{|x_i|-\lambda,0\}=1\xrightarrow{solve}\lambda^*\right)$, then, $[x]_X=sthr_{\lambda^*}(x)$ {elementwise}

$$z_i^*(\lambda)=\|\begin{cases}0,&\tfrac{|x_i|}{\lambda}\leq 1\\ x_i-\lambda\,sgn(x_i),&else\end{cases}\| =\begin{cases}|0|=0,&\tfrac{|x_i|}{\lambda}\leq 1\\ |x_i-\lambda\,sgn(x_i)|,&\tfrac{|x_i|}{\lambda}>1\end{cases} =\begin{cases}0,&\text{if }|x_i|-\lambda\leq 0\\ |x_i|-\lambda,&\text{if }|x_i|-\lambda>0\end{cases}=max\{|x_i|-\lambda,0\}$$

(combining both $|x_i-\lambda\,sgn(x_i)|=|x_i|-\lambda$)

$[x]_X=\begin{cases}x,&\text{if }\|x\|_1\leq 1\\ sthr_{\lambda^*}(x),&\text{if }\|x\|_1>1\\ \quad\{solve_\lambda\\ \quad(\sum_{i=1}^n max\{|x_i|-\lambda,0\}=1)\end{cases}$  Ans:

---

$\bullet$ Projection onto the positive semidefinite cone.
$X=\{X\in S^n:X\succeq 0\}=S_+^n$

Given: $X\notin S_+^n$

$\square^T$: $[X]_X=\underset{Z\in X}{argmin}\|Z-X\|_F^2$ $\{\|X\|_F^2=\sum_{i=1}^n\zeta_i^2\}$

// A symmetric matrix can be diagonalized by orthogonal similarity transformation

$X=U\Lambda U^T$ {U orthogonal matrix}
$\Lambda=diag(\lambda_1,\dots,\lambda_n)$

$\|Z-X\|_F^2=\|Z-U\Lambda U^T\|_F^2=\|UU^TZU U^T-U\Lambda U^T\|_F^2=\|U(U^TZU-\Lambda)U^T\|_F^2$ {norm does not change when you multiply something by orthogonal matrix}

$=\|U^TZU-\Lambda\|_F^2=\|\tilde{Z}-\Lambda\|_F^2$ {the operator: $U^T\square U:S^n\to S^n$ is a one to one operator $(X\neq Y\to U^TXU\neq U^TYU)$}

so $[X]_X=\underset{Z\in X}{argmin}\|Z-X\|_F^2=\underset{\tilde{Z}}{argmin}\|U^TZU-\Lambda\|_F^2$

$=U\left[\underset{\tilde{Z}\in X}{argmin}\|\tilde{Z}-\Lambda\|\right]U^T$ #Say, $Z^*$ solves $\underset{Z\in X}{\nabla}\|U^TZU-\Lambda\|_F^2$

$\therefore\underset{\tilde{Z}\geq 0}{argmin}\|\tilde{Z}-\Lambda\|_F^2=[\Lambda]_+=diag([\lambda_1]_+,\dots,[\lambda_n]_+)$

|| Analogous with $\underset{Z\geq X}{argmin}\|Z-x\|_2^2=[X]_+$   [ eq: Projection on positive orthant ]

$\boxed{[X]_X=U[\Lambda]_+U^T}$  Ans:

---

② Proximal map of $\ell_1$ regularization:

Lasso problem: ($\ell_1$ regularized least square)

$\underset{x}{\nabla}\ \tfrac{1}{\gamma}\|Ax-b\|_2^2+\|x\|_1$
$\quad f_0(x)\qquad h(x)$
$\{\square_{strongly}\}\quad\{D,\text{nondifferentiable}\}$
$\uparrow$
A {fullrank}

Proximal algorithm:
$\forall_k\ x_{k+1}=\underset{s_k h}{prox}\left(x_k-s_k\nabla f_0(x_k)\right)$

$\nabla\tfrac{1}{\gamma}\|Ax_k-b\|_2^2=2\tfrac{1}{\gamma}A^T(Ax_k-b)$

$=\underset{s_k h}{prox}\left(x_k-s_k\tfrac{2}{\gamma}A^T(Ax_k-b)\right)$

$=\underset{s_k h}{prox}\left(1x_k-\tfrac{2s_k}{\gamma}A^TAx_k+\tfrac{2s_k}{\gamma}A^Tb\right)$

$=\underset{s_k h}{prox}\left((1-\tfrac{2s_k}{\gamma}A^TA)x_k+\tfrac{2s_k}{\gamma}A^Tb\right)=sthr_{s_k}\left((1-\tfrac{2s_k}{\gamma}A^TA)x_k+\tfrac{2s_k}{\gamma}A^Tb\right)$ $\to\ \therefore\forall_k\ x_{k+1}=sthr_{s_k}\left((1-\tfrac{2s_k}{\gamma}A^TA)x_k+\tfrac{2s_k}{\gamma}A^Tb\right)$

NOW! a closed form iteration :)

in terms of fixed point theory:
$$x^*=sthr_{s_k}\left((1-\tfrac{2s_k}{\gamma}A^TA)x^*+\tfrac{2s_k}{\gamma}A^Tb\right)$$ जो solve करने पर हमें optimal solution (लास प्रॉब) ① ....

{ $\underset{h}{prox}(x)=\underset{z}{argmin}\left(h(z)+\tfrac{1}{2}\|z-x\|_2^2\right)$ }

$\underset{s_k h}{prox}(x)=\underset{z}{argmin}\left(s_k h(z)+\tfrac{1}{2}\|z-x\|_2^2\right)$
$=\underset{z}{argmin}\left(s_k\|z\|_1+\tfrac{1}{2}\|z-x\|_2^2\right)=sthr_{s_k}(x)$

$\underset{z}{argmin}\left(s_k\|z\|_1+\tfrac{1}{2}\|z-x\|_2^2\right)=[sthr_{s_k}(x_i)]_{i=1}^n=sthr_{s_k}(x)$

as [ eq: Optimization Problem for l-1 norm projection ] (ऊपर बताया) (ग्राफ)

$\underset{z}{argmin}\left(\tfrac{1}{2}\|z-x\|_2^2+\lambda\|z\|_1\right)=sthr_\lambda(x)=[sthr_\lambda(x_i)]_{i=1}^n$

---

$\bullet$ ISTA (Iterative shrinkage-thresholding algorithm) for Lasso:
For lasso:
$f(x)=\tfrac{1}{\gamma}\|Ax-b\|_2^2$ ...

gradient Lipschitz constraint

For lasso:

$$f_0(x) = \frac{1}{\gamma}\|Ax-b\|_2^2$$

$$\nabla f_0(x) = \frac{1}{\gamma} 2A^T(Ax-b) = \frac{1}{\gamma}(2A^TAx - 2A^Tb)$$

$$\nabla^2 f_0(x) = \frac{2}{\gamma} A^TA \quad // \text{Matrix cookbook 96-98:} \quad f = x^T\Lambda x + b^Tx \to \begin{array}{l} \nabla_x f = \Lambda x \\ \nabla_x^2 f = \Lambda \end{array}$$

now যেহেতু $f_0(x)$ এর জন্য $L, m$ এর হিসাব করে রাখা is come up with a stopping criterion for proximal gradient algorithm.

$$\{\;[\text{eq. 12.60}]\; \text{থেকে পাই:}\quad \|g_k\|_2^2 \le 2\epsilon\frac{mL}{L-m} \Rightarrow f(x_{k+1}) - f(x^*) \le \epsilon \;\}$$

gradient Lipschitz constraint
Strong convexity constraint

✱ Strong convexity constraint for $f_0$:

✱ $f_0(x) \{\square_{\text{strongly}} \}  \Leftrightarrow \forall_{x \in \text{dom} f_0} \quad \nabla^2 f_0(x) \succeq m I ]$

$$\Leftrightarrow \forall_{x \in \text{dom} f_0}, \quad \nabla^2 f_0(x) - mI \succeq 0$$

∯ now the matrix $A^TA$ symmetric positive semidefinite, so orthogonal similarity transformation is possible with all eigenvalues non negative:

$$A^TA = Q\Lambda Q^T = \sum_{i=1}^n \lambda_i\, q_i q_i^T \quad \{\lambda_1 \ge \cdots \ge \lambda_n \ge 0\}$$

$$\therefore \nabla^2 f_0(x) - mI = \frac{2}{\gamma} Q\Lambda Q^T - m\underbrace{q q^T}_{1} = Q\left(\frac{2}{\gamma}\Lambda - mI\right)Q^T$$

eigenvalues of $\nabla^2 f_0(x) - m]$  $\nabla^2 f_0(x) - m]$ symmetric
and we have just found an orthogonal similarity transformation
of that, so the diagonal matrix will correspond to the eigenvalues

$$\to \text{all eigenvalues} \ge 0 \;\Leftrightarrow\; \forall_i \quad \frac{2}{\gamma}\lambda_i - m \ge 0$$

$$\Leftrightarrow \frac{2}{\gamma}\min(\lambda_i) - m \ge 0$$

$$\Leftrightarrow \frac{2}{\gamma}\lambda_{min}(A^TA) \ge m$$

$$\Leftrightarrow \boxed{m_{max} = \frac{2}{\gamma}\lambda_{min}(A^TA)} \quad [\text{eq: Strong Convexity Constant Lasso}]$$

this can be set as the strong convexity
constraint of $f_0(x) = \frac{1}{\gamma}\|Ax-b\|_2^2$

✱ Finding a global Lipschitz constraint:

∯ From Lemma 12·1·1:  $f\; \{f: \mathbb{R}^n \to \mathbb{R}, \text{ gradient\_Lipschitz\_continuous}, \delta_{g,2}\} \Leftrightarrow \forall_\square \quad \|\nabla^2 f(\square)\|_F = \left(\sum c_i^2\right)^{1/2} \le L\; \}$   $\lambda_i^2$ (for a symmetric PSD matrix $\lambda_i = c_i$)

$$f_0(x) = \frac{1}{\gamma}\|Ax-b\|_2^2\; \{f:\mathbb{R}^n\to\mathbb{R},\quad n\quad,\quad n\} \Leftrightarrow \forall_x \quad \|\nabla^2 f(x)\|_F = \|\frac{2}{\gamma}A^TA\|_F = \frac{2}{\gamma}\|A^TA\|_F = \frac{2}{\gamma}\left(\sum_{i=1}^n \lambda_i^2(A^TA)\right)^{1/2} \le L$$

$\{\ast: \|\alpha x\| = |\alpha|\|x\|\}$

$$\therefore L_{min} = \frac{2}{\gamma}\left(\sum_{i=1}^n \lambda_i^2(A^TA)\right)^{1/2} \quad [\text{eq: Gradient Lipshitz Constant Lasso}]$$

this we set as the Gradient Lipschitz
constant.

[ eq: Strong Convexity Constant Lasso ]    [ eq: Gradient Lipshitz Constant Lasso ]

Now, we know both $m$ and $L$, so let's give the proximal gradient algorithm for Lasso (constant stepsize):

Require: $\epsilon > 0, x_0, A$ full rank.

1. Compute $m = \frac{2}{\gamma}\lambda_{min}(A^TA),\; L = \frac{2}{\gamma}\left(\sum_{i=1}^n c_i^2(A^TA)\right)^{1/2}$

2. $k := 0,\; s = 1/L$

3. $\square\; \nabla f_0(x_k) = \frac{2}{\gamma}\left(A^T(Ax_k - b)\right)$

4. $\square\; x_{k+1} = \text{sthr}_s\left(x_k - s\nabla f_0(x_k)\right)$

5. $\square\; \|g_k\|_2 = \|x_k - x_{k+1}\|_2 / s\quad \{\;x_{k+1} = x_k - s_k g_k \to \|g_k\|_2 = \frac{\|x_k - x_{k+1}\|}{s}\;\}$

6. If $\left(\|g_k\|_2^2 \le 2\epsilon\frac{mL}{L-m}\right)$
      done!, return $x^* = x_{k+1}$

   else
        $k := k+1,$
        go to 3

✱ Fast Proximal gradient (constant step sizes)

Normal proximal gradient এর convergence rate $(\frac{1}{k})$, suitable modification যা মারফত $(\frac{1}{k^2})$ convergence rate achieve করা যায়, ঐ type যা algorithm কে Fast Proximal Gradient algorithm বলে, it has two versions

1. When $L$ is known, $s_k = s = 1/L$
2. $L$ is not known. then backtracking type যা line search করতে হয়  এবং  $s_k$ কে বের করা যায়

[ Probably need to elaborate later ]