

- Algorithm:
 - proximal viewpoint of S_G
 - Non-Euclidean geometry
 - Bregman divergence
 - properties
- Convergence analysis : Convex \Leftarrow L -Lipschitz

Recall, Projected (subgradient) method

$$\left. \begin{array}{l} \text{minimize } f(\underline{x}) \\ \underline{x} \in C \end{array} \right\} \begin{array}{l} \underline{x}_{t+1} = \text{Proj}_C(\underline{x}_t - \eta_t \underline{g}_t) \\ \underline{g}_t \in \partial f(\underline{x}_t) \end{array}$$

if f is L -Lipschitz $\left[\|\underline{g}\| \leq L \text{ or } |f(\underline{x}) - f(\underline{y})| \leq L \|\underline{x} - \underline{y}\| \text{ for } \underline{x}, \underline{y} \in \text{dom}(f) \text{ and } \underline{g} \in \partial f(\underline{x}) \right]$, then $\eta_t = \eta = \frac{R}{L\sqrt{T}}$

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\underline{x}_t) - f(\underline{x}^*) \leq R L \frac{1}{\sqrt{T}}$$

- Are these dimension independent?
- Are dimension dependent const. in L ?

Recall: Proximal view of proj. subgradient method

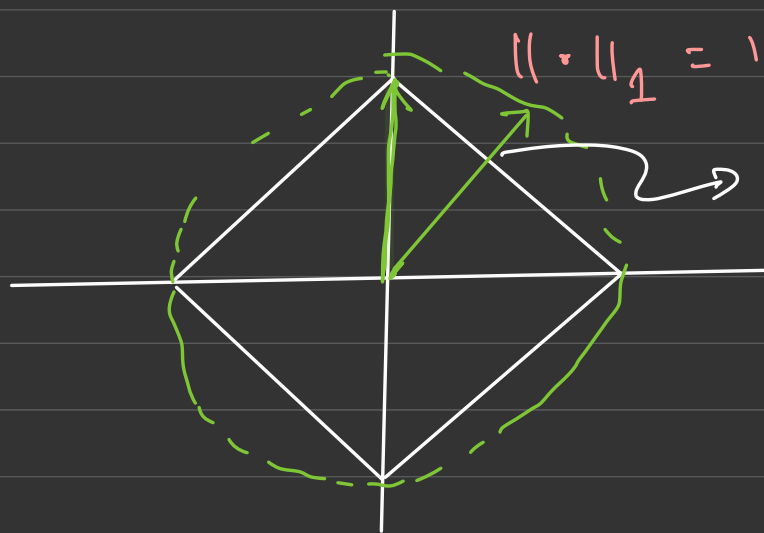
$$\underline{x}_{t+1} = \arg \min_{\underline{x}} f(\underline{x}_t) + \underline{g}_t^\top (\underline{x} - \underline{x}_t) + \frac{1}{2\eta} \|\underline{x} - \underline{x}_t\|_2^2$$

$$= \arg \min_{\underline{x}} \eta \underline{g}_t^\top \underline{x} + \frac{1}{2} \|\underline{x} - \underline{x}_t\|_2^2$$

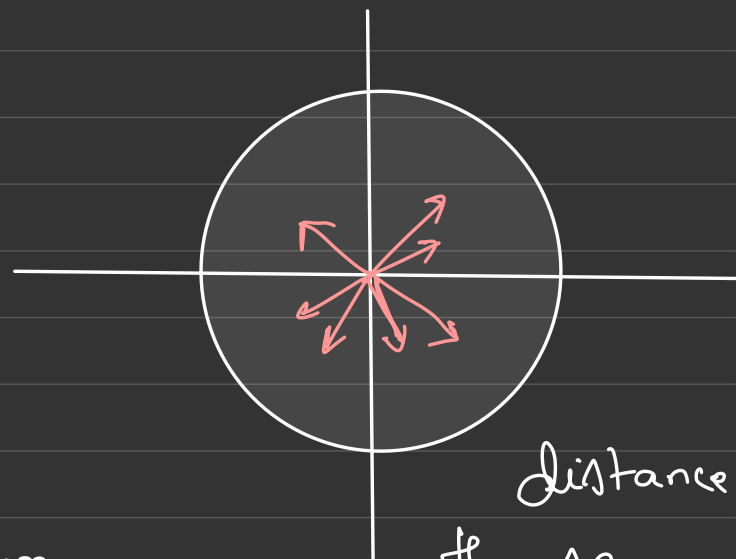
(proximal term)

Euclidean norm: Spherical Symmetry

Not true about other norms:



further away as measured by the l_1 -norm



distance is the same for all these steps.

Example of Quadratic minimization:

$$f(x_1, x_2) = x_1^2 \cdot \frac{1}{100} + x_2^2 \cdot 100$$



Suppose we are at $x_t = \begin{pmatrix} -10 \\ -0.1 \end{pmatrix}$

$$\nabla f(x_t) = \begin{pmatrix} 2x_1/100 \\ 2x_2 \cdot 100 \end{pmatrix} \Bigg|_{\begin{pmatrix} -10 \\ -0.1 \end{pmatrix}} = \begin{pmatrix} -1/5 \\ -20 \end{pmatrix}$$



$$\underline{x}_{t+1} = \underset{\underline{x}}{\operatorname{arg\,min}} \quad \eta \nabla f(\underline{x}_t)^T \underline{x} + \underbrace{\frac{1}{2} (\underline{x} - \underline{x}_t)^T \mathbf{I} (\underline{x} - \underline{x}_t)}_{\|\underline{x} - \underline{x}_t\|_g^2}$$

Suppose:

$$\underline{x}_{t+1} = \underset{\underline{x}}{\operatorname{arg\,min}} \quad \eta \nabla f(\underline{x}_t)^T \underline{x} + \frac{1}{2} (\underline{x} - \underline{x}_t)^T \mathcal{Q} (\underline{x} - \underline{x}_t)$$

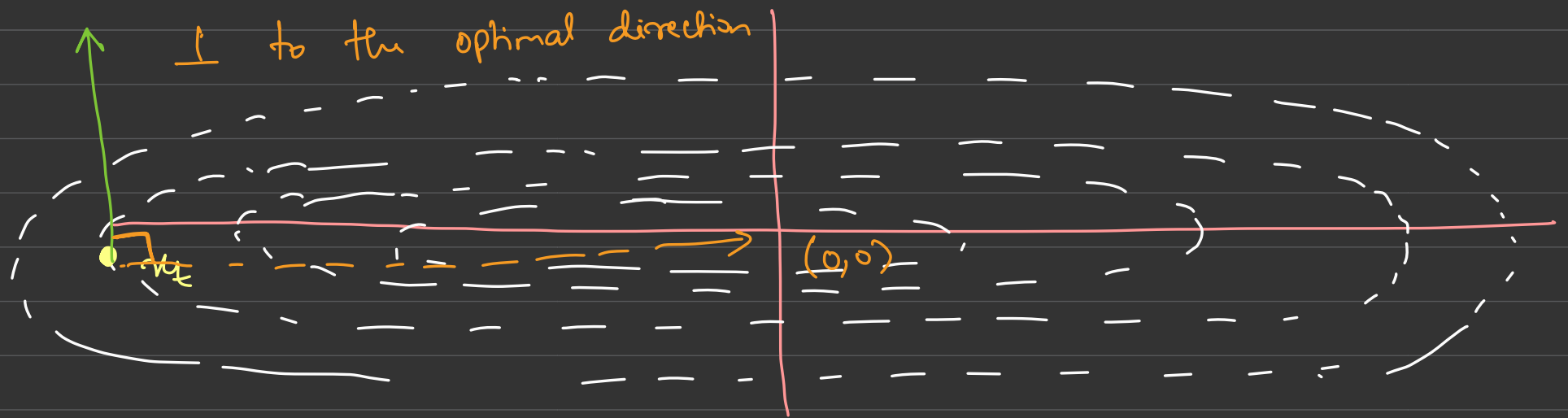
$$\eta \nabla f(\underline{x}_t)^T + \mathcal{Q} (\underline{x} - \underline{x}_t) = 0$$

$$\underline{x}_{t+1} = \underline{x}_t - \eta \mathcal{Q}^{-1} \nabla f(\underline{x}_t)^T$$

skew gradients
to fit
geometry
better

$$Q = \begin{pmatrix} \frac{1}{50} & 0 \\ 0 & 200 \end{pmatrix} \Rightarrow x_{t+1} = x_t - \eta \begin{pmatrix} 50 & 0 \\ 0 & \frac{1}{200} \end{pmatrix} \begin{pmatrix} -\frac{1}{5} \\ -20 \end{pmatrix}$$

$$= \begin{pmatrix} -10 \\ -0.1 \end{pmatrix} - \eta \begin{pmatrix} -10 \\ -0.1 \end{pmatrix}$$



- Mirror descent:
- measure distance using a different norm?
 - adjust gradient updates to fit the geometry of the problem.

This changes Lipschitzness of the function.

$$|f(\underline{x}) - f(\underline{y})| \leq L \|\underline{x} - \underline{y}\|_2$$

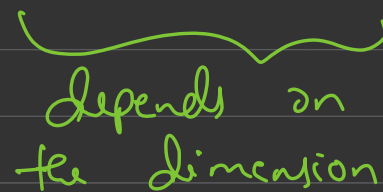
• Example: $f(\underline{x}) = \|\underline{x}\|_1$ with $\underline{x} = \underline{1}$ and $\underline{y} = (1+\epsilon)\underline{1}$

$$|\|\underline{x}\|_1 - \|\underline{y}\|_1| \leq L \|\underline{x} - \underline{y}\|_2$$

$$|\eta - (1+\epsilon)\eta| = \epsilon\eta \leq L \|\underline{x} - \underline{y}\|_2 = L \|\underline{1} - (1+\epsilon)\underline{1}\|_2$$

$$\epsilon\eta \leq L \epsilon\sqrt{\eta}$$

For the upper bound to be valid: $L = \sqrt{\eta}$

 depends on
the dimension

$$\Rightarrow \|\nabla f(\underline{x})\|_2 \leq \sqrt{\eta}$$

• Suppose f is 1-Lipschitz w.r.t. a

different norm say $\|\nabla f(x)\|_\infty \leq 1$

(we had this for $\|\cdot\|_1$) $\Rightarrow \|\nabla f(x)\|_2 < \sqrt{n}$

Mirror descent fixes this by:

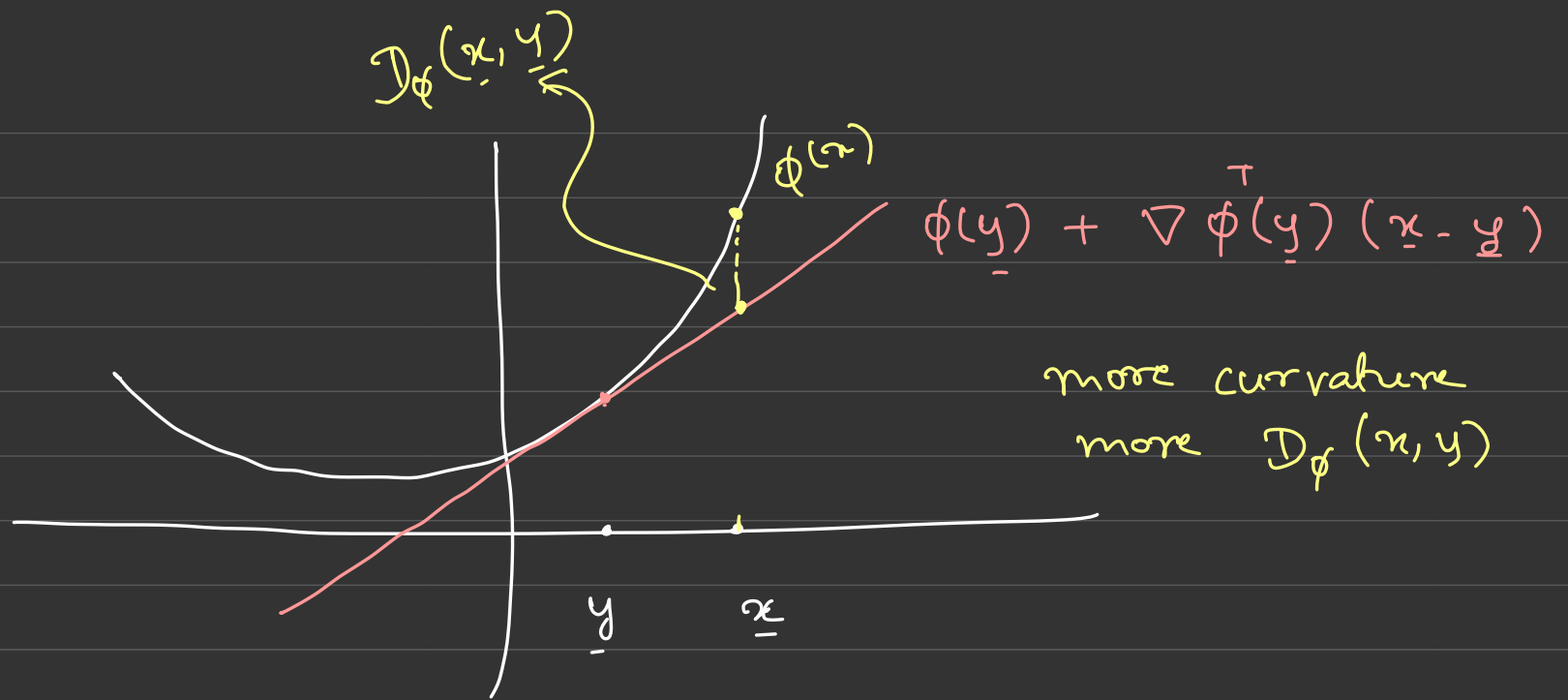
replacing the $\|x_t - x_{t-1}\|_2^2$ with

distance-like metric

$$D_\phi(x, y) = \phi(x) - \left[\phi(y) + \nabla \phi^\top(y) (x - y) \right]$$



Bregman divergence for convex and differentiable $\phi(x)$



Example:

$$\phi(\underline{x}) = \frac{1}{2} \|\underline{x}\|_2^2$$

$$\mathcal{D}_\phi(\underline{x}, \underline{y}) = \frac{1}{2} \|\underline{x}\|_2^2 - \left[\frac{1}{2} \|\underline{y}\|_2^2 + \underline{y}^\top (\underline{x} - \underline{y}) \right]$$

$$= \frac{1}{2} \left[\|\underline{x}\|_2^2 - 2 \underline{y}^\top \underline{x} - \|\underline{y}\|_2^2 + 2 \|\underline{y}\|_2^2 \right]$$

$$= \frac{1}{2} \|\underline{x} - \underline{y}\|_2^2$$

$$\textcircled{2} \quad \phi(\underline{x}) = \frac{1}{2} \underline{x}^T Q \underline{x} \quad Q \succ 0$$

$$D_\phi(\underline{x}, \underline{y}) = \frac{1}{2} (\underline{x} - \underline{y})^T Q (\underline{x} - \underline{y})$$

$$Q = \begin{bmatrix} 1/50 & \\ & 200 \end{bmatrix}$$

$D_\phi(\underline{x}, \underline{y})$
vs $\|\underline{x} - \underline{y}\|_2^2$?

Squared Mahalanobis distance

$$\textcircled{3} \quad \phi(\underline{x}) = \sum_i x_i \log x_i \quad (\text{negative entropy})$$

$$D_\phi(\underline{x}, \underline{y}) = \text{KL}(\underline{x} \parallel \underline{y}) = \sum_i x_i \log \left(\frac{x_i}{y_i} \right)$$

(check as an exercise)

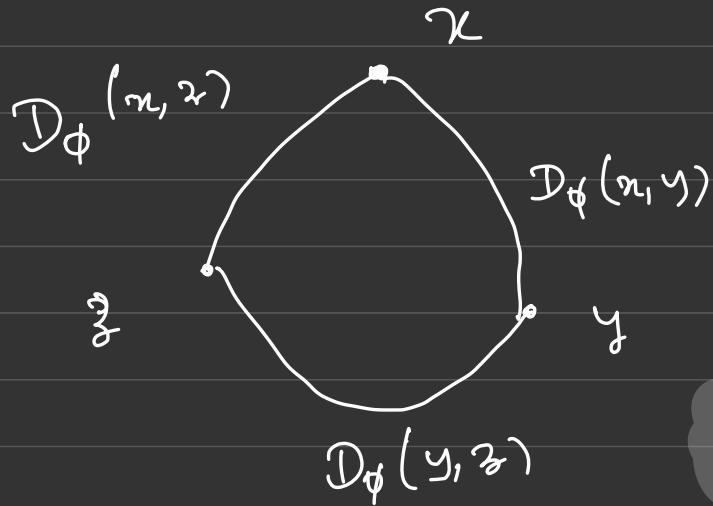
Home work: with $\Delta = \{ \underline{x} \in \mathbb{R}_+^n \mid \underline{1}^T \underline{x} = 1 \}$ probability simplex

MD results in "entropic descent".

Properties of Bregman divergence:

Three-point Lemma:

$$D_\phi(\underline{x}, \underline{z}) = D_\phi(\underline{x}, \underline{y}) + D_\phi(\underline{y}, \underline{z}) - (\nabla\phi(\underline{z}) - \nabla\phi(\underline{y}))^\top (\underline{x} - \underline{y})$$



cosine law for the Euclidean case.

$$\bullet \quad \|\underline{x} - \underline{z}\|_2^2 = \|\underline{x} - \underline{y}\|_2^2 + \|\underline{y} - \underline{z}\|_2^2 - 2(\underline{z} - \underline{y})^\top (\underline{x} - \underline{y})$$

Proof:

$$\begin{aligned} & D_\phi(\underline{x}, \underline{y}) + D_\phi(\underline{y}, \underline{z}) - D_\phi(\underline{x}, \underline{z}) \\ &= \phi(\underline{x}) - \phi(\underline{y}) - \nabla\phi(\underline{y})^\top (\underline{x} - \underline{y}) + \phi(\underline{y}) - \phi(\underline{z}) \\ &\quad - \nabla\phi(\underline{z})^\top (\underline{y} - \underline{z}) - \left[\phi(\underline{x}) - \phi(\underline{z}) - \nabla\phi(\underline{z})^\top (\underline{x} - \underline{z}) \right] \\ &= \cancel{\phi(\underline{x})} - \cancel{\phi(\underline{y})} - \nabla\phi(\underline{y})^\top (\underline{x} - \underline{y}) + \cancel{\phi(\underline{y})} - \cancel{\phi(\underline{z})} \\ &\quad - \nabla\phi(\underline{z})^\top (\underline{y} - \underline{z}) - \left[\cancel{\phi(\underline{x})} - \cancel{\phi(\underline{z})} - \nabla\phi(\underline{z})^\top (\underline{x} - \underline{z}) \right] \end{aligned}$$

$$= -\nabla\phi^T(\underline{y}) (\underline{x} - \underline{y}) - \nabla\phi^T(\underline{z}) (\underline{y} - \underline{z}) \\ + \nabla\phi^T(\underline{z}) (\underline{x} - \underline{z})$$

$$= \left[\nabla\phi(\underline{z}) - \nabla\phi(\underline{y}) \right]^T (\underline{x} - \underline{y})$$

② Convexity of $D_\phi(\underline{x}, \underline{y})$ in \underline{x}

$$D_\phi(\underline{x}, \underline{y}) = \underbrace{\phi(\underline{x}) - \phi(\underline{y}) - \nabla\phi^T(\underline{y}) (\underline{x} - \underline{y})}$$

follows from convexity of $\phi(\cdot)$

L-Lipschitz definition for arbitrary norm:

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

Defn. of subgradients : $f(y) \geq f(x) + g^T (y - x)$

$$f(x) - f(y) \leq g^T (y - x)$$

$$|f(x) - f(y)| \leq \|g\|_* \|x - y\|$$

(Generalize Cauchy-Schwarz or Holder's inequality)

This gives us another definition

of L-Lipschitz w.r.t. $\|\cdot\|$

$$\|g\|_* \leq L$$

Convex and Lipschitz Problems:

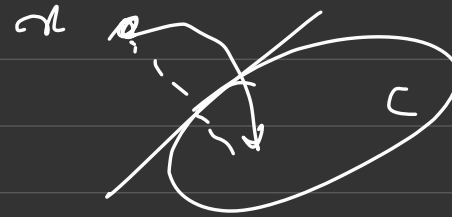
minimize $f(x)$

s. to $x \in C$

- f is convex and Lipschitz continuous
 - $\|g\|_* \leq L$ for any $g \in \partial f(x)$
- $\phi(\cdot)$ is ρ -strongly convex w.r.t. $\|\cdot\|$

$$\begin{aligned} \Rightarrow D_\phi(\underline{x}, \underline{y}) &= \phi(\underline{x}) - [\phi(\underline{y}) + \nabla \phi(\underline{y})^\top (\underline{x} - \underline{y})] \\ &\geq \frac{\rho}{2} \|\underline{x} - \underline{y}\|^2 \end{aligned}$$

Bregman projection:



Given a point \underline{x} ,

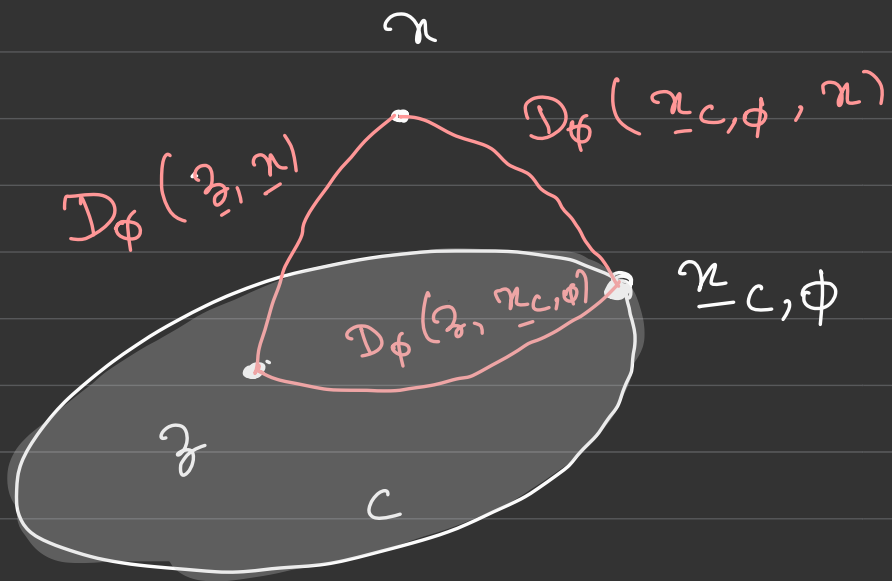
$$P_{C, \phi}(\underline{x}) = \arg \min_{\underline{z} \in C} D_{\phi}(\underline{x}, \underline{z})$$

if $\phi = \frac{1}{2} \|\underline{x}\|_2^2$, we have

$$D_{\phi}(\underline{x}, \underline{z}) = \frac{1}{2} \|\underline{x} - \underline{z}\|_2^2.$$

Then, $P_{C, \phi}(\underline{x}) = P_C(\underline{x})$ [Orthogonal Euclidean projection]

Generalized Pythagorean theorem:



$$\text{if } \underline{x}_{C,\phi} = P_{C,\phi}(\underline{x})$$

Then,

$$D_\phi(\underline{z}, \underline{x}) \geq D_\phi(\underline{z}, \underline{x}_{C,\phi})$$

$$+ D_\phi(\underline{x}_{C,\phi}, \underline{x})$$

$$\forall \underline{z} \in C$$

$$\underline{x}_{C,\phi} = \underset{\underline{z} \in C}{\text{arg min}} D_\phi(\underline{z}, \underline{x})$$

Recall optimality condn: $\underline{g}^\top (\underline{z} - \underline{x}_{C,\phi}) \geq 0 \quad \forall \underline{z} \in C$

$$\text{Let } \underline{g} = \left. \nabla_{\underline{z}} D_\phi(\underline{z}, \underline{x}) \right|_{\underline{z} = \underline{x}_{C,\phi}}$$

$$D_{\phi}(\underline{z}, \underline{x}) = \phi(\underline{z}) - \left[\phi(\underline{x}) + \nabla \phi^{\top}(\underline{x})(\underline{z} - \underline{x}) \right]$$

$$\Rightarrow \underline{g} = \nabla \phi(\underline{x}_{c,\phi}) - \nabla \phi(\underline{x})$$

$$\underline{g}^{\top}(\underline{z} - \underline{x}_{c,\phi}) \geq 0$$

$$\Rightarrow \left[\nabla \phi(\underline{x}_{c,\phi}) - \nabla \phi(\underline{x}) \right]^{\top} [\underline{z} - \underline{x}_{c,\phi}] \geq 0$$

Three-point Lemma : $D_{\phi}(\underline{x}, \underline{z}) = D_{\phi}(\underline{x}, \underline{y}) + D_{\phi}(\underline{y}, \underline{z})$
 $- (\nabla \phi(\underline{z}) - \nabla \phi(\underline{y}))^{\top} (\underline{x} - \underline{y})$

$$0 \geq D_{\phi}(\underline{z}, \underline{x}_{c,\phi}) + D_{\phi}(\underline{x}_{c,\phi}, \underline{x}) - D_{\phi}(\underline{z}, \underline{x})$$

