

Distributed Optimization:

$$\min_{x \in \mathcal{X}} f(x) = \frac{1}{n} \sum_i f_i(x)$$

graph: $G = (V, E)$

ea. node i has their x_i and access to $f_i(x)$ and info from $N(i)$

Motivating Example:

i : separate set of training data

x : parameters

linear $f_i(x) = \|z_i - H_i x\|^2$

gen. $f_i(x) = \|h_i(z_i, x)\|^2$

Properties of $f_i(x)$:

$f_i(x)$: convex "bowl shaped"

$g_i \in \partial f_i(x)$ subgradient

$$g_i = \frac{\partial f_i}{\partial x}$$

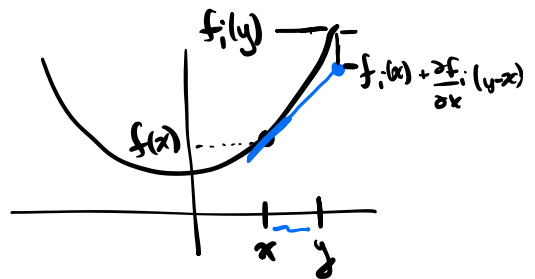


Convexity: x, y lower bound

$$f_i(y) \geq f_i(x) + \frac{\partial f_i}{\partial x}(y-x) + \epsilon \|y-x\|^2$$

convexity

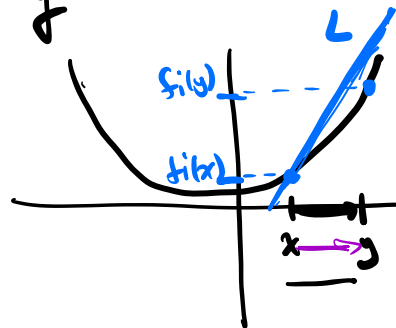
strong convexity



L-Lipschit & continuity:

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

for any $x, y \in \mathcal{X}$



for an L -Lipschitz function: f_i with $g_i = \frac{\partial f_i}{\partial x}$

$$\|g_i\|_x \leq L$$

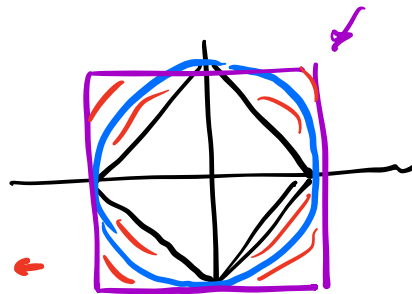
Norms:

$$\|x\|_1 = \sum |x_i|$$

$$\|x\|_2 = \left(\sum |x_i|^2\right)^{1/2}$$

$$\|x\|_p = \left(\sum |x_i|^p\right)^{1/p} \quad 1 \leq p \leq \infty$$

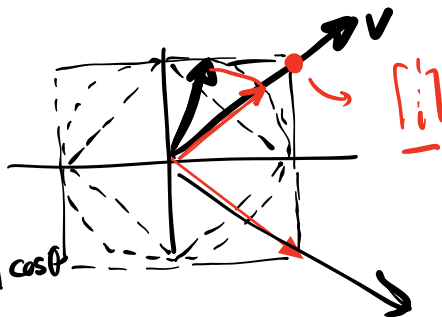
$$\|x\|_\infty = \max |x_i|$$



dual norm:

$$\|v\|_{p,*} = \sup_{\|u\|_p=1} \langle v, u \rangle$$

$\rightarrow \|v\| \|u\| \cos \theta$



$$\|v\|_{2,*} = \|v\|_2$$

$$\|v\|_{\infty,*} = \|v\|_1$$

$$\|v\|_{1,*} = \|v\|_\infty$$

Proximal Function $\Psi(x)$ used for projection $\Pi_X^\Psi(z, \alpha)$

$\Psi(x)$: strongly convex

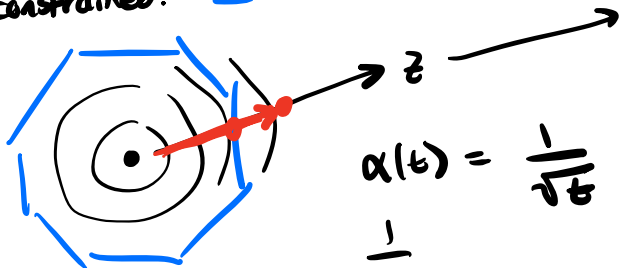
$$\Psi(x) = \frac{1}{2} \|x\|_2^2$$

$$x = \Pi_X^\Psi(z, \alpha) = \underset{x \in X}{\operatorname{argmin}} \left[-\langle z, x \rangle + \frac{1}{\alpha} \Psi(x) \right]$$

unconstrained:

$$-z^T + \frac{1}{\alpha} x^T = 0$$

$$\Rightarrow x = \alpha z$$



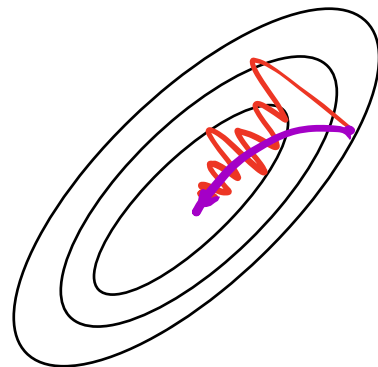
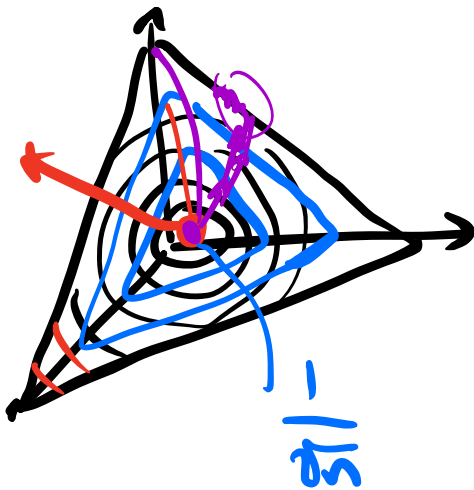
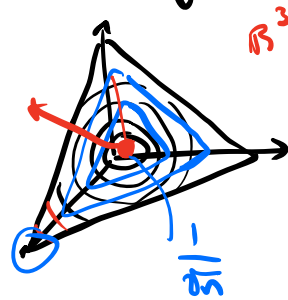
$$\alpha(t) = \frac{1}{\sqrt{t}}$$

\downarrow
 $\frac{1}{\alpha(t)}$ \uparrow

if x is on simplex "x is a discrete probability dist."

$$\Delta_n = \{x \in \mathbb{R}^n \mid \mathbf{1}^T x = 1, x \geq 0\} \quad x \in \mathbb{R}^3$$

$$\psi(x) = \sum_i x_i \log(x_i) - x_i \quad \text{and } \|\cdot\|_4$$



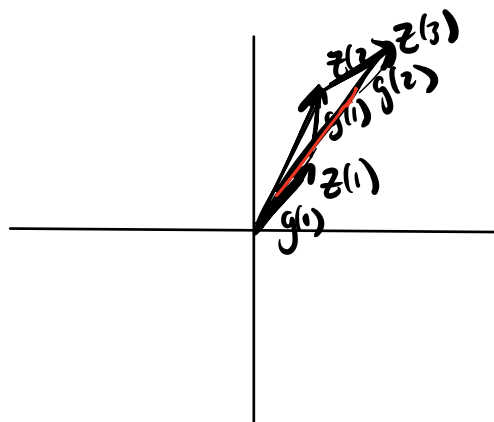
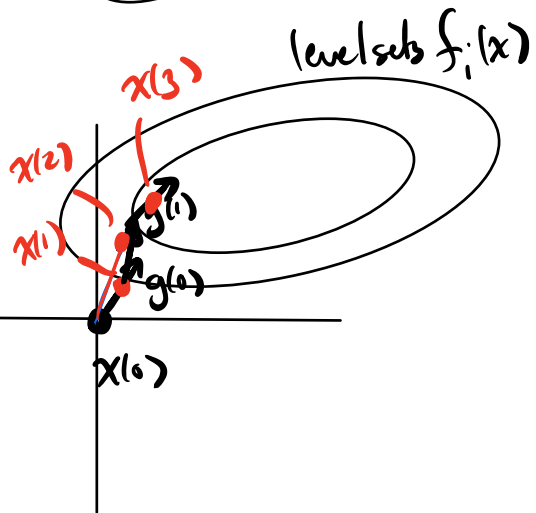
Dual Averaging: ↙

$$z(t+1) = z(t) - g(t)$$

$$g(t) \in \partial f_i(x_i(t))$$

$$g(t) = \frac{\partial f_i}{\partial x}(x_i(t))$$

$$x(t+1) = \Pi_x^\psi(-z(t+1), \alpha(t))$$



Traditional:

$$\underline{x(t+1)} = \underline{\Pi (x(t) + \alpha g(t))}$$

Linear Dynamics:

$$\underline{z(t+1)} = \underline{z(t)} - \underline{g(t)}$$

$$\underline{z(t+1)} = \underline{A z(t)} + \underline{B u(t)}$$

Diagram: A blue arrow points from $z(t)$ to $A z(t)$. A purple arrow points from $z(t)$ to A . A purple arrow points from $u(t)$ to $B u(t)$. A purple arrow points from $g(t)$ to $-g(t)$. A purple arrow points from $-I$ to A . A purple arrow points from $-I$ to B .

$$z(0) = z_0$$

$$z(1) = A z(0) + B u(0)$$

$$z(2) = \underline{A^2 z(0)} + \underline{A B u(0)} + B u(1)$$

$$z(3) = A^3 z(0) + A^2 B u(0) + A B u(1) + B u(2)$$

$$z(t) = A^t z(0) + \sum_{j=0}^{t-1} A^{t-j-1} B u(j)$$

$$z(t) = \underbrace{A^t z(0)}_{\text{drift term}} + \underbrace{\begin{bmatrix} A^{t-1} B & A^{t-2} B & \dots & A B & B \end{bmatrix}}_{\text{Reachability Matrix}} \begin{bmatrix} u(0) \\ \vdots \\ u(t-1) \end{bmatrix}$$

Distributed Scheme

ca. node i $\{x_i(t), z_i(t)\}$

compute $g_i(t) \in \partial f_i(t)$
 \hookrightarrow local f_i

receive $z_j(t) \in j \in N(i)$

Communication matrix $P \in \mathbb{R}^{n \times n}$

P is doubly stochastic, symmetric

$P_{ij} > 0$ if and only if $j \in N(i)$

$$\sum_j P_{ij} = \sum_{j \in N(i)} P_{ij} = 1 \quad P \mathbf{1} = \mathbf{1}$$

$$\sum_i P_{ij} = \sum_{i \in N(j)} P_{ij} = 1 \quad \mathbf{1}^T P = \mathbf{1}^T$$

Before: $z(t+1) = z(t) - g(t)$ $\left. \vphantom{z(t+1)} \right\} \leftarrow \leftarrow$

Now:

$$z_i(t+1) = \sum_{j \in N(i)} P_{ij} z_j(t) - g_i(t) \quad x_i(t+1) = \prod_{\alpha}^{\varphi} (z_i(t+1), \alpha(t))$$

\rightarrow vector case for z_i & x_i \star

for scalar: z_i, x_i, g_i

$$z = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad g = \begin{bmatrix} g_1 \\ \vdots \\ g_n \end{bmatrix}$$

$$z(t+1) = Pz(t) - g$$

$$z(t) = P^t z(0) - \sum_{s=0}^{t-1} P^{t-s-1} g(s) \quad *$$

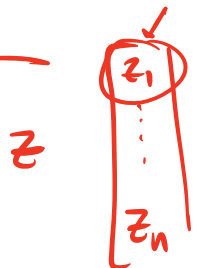
Question:

each z_i is a descent direction
 what is the average descent direction doing?

want to agree on optimal descent direction

$$\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$$

in scalar z_i case:



z

$$\frac{1}{n} \mathbf{1}^T [z(t+1) = Pz(t) - g]$$

$$\frac{1}{n} \mathbf{1}^T z(t+1) = \frac{1}{n} \mathbf{1}^T P z(t) - \frac{1}{n} \mathbf{1}^T g$$

$$\bar{z}(t+1) = \bar{z}(t) - \frac{1}{n} \sum_i g_i(x_i)$$

average descent direction

will be overall gradient when all x_i 's agree.

$$\hat{x}_i(\tau) = \frac{1}{\tau} \sum_{t=1}^{\tau} x_i(t) \rightarrow \text{time average of } x_i$$

Thm 1: Basic Convergence

$$f(\hat{x}_i(\tau)) - f(x^*) \leq \frac{1}{\tau \alpha(\tau)} \psi(x^*) + \frac{L^2}{2\tau} \sum_{t=1}^{\tau} \alpha(t-1) + \frac{3L}{\tau} \max_j \sum_{t=1}^{\tau} \alpha(t) \|\bar{z}_i(t) - z_j(t)\|$$

AP ✓
converges to 0

Thm 2: Rates (spectral gap)

$$f(\hat{x}_i(\tau)) - f(x^*) \leq 8 \frac{RL}{\sqrt{\tau}} \frac{\log(\tau \sqrt{n})}{\sqrt{1 - \sigma_2(P)}} \quad \text{for all } i \in V$$

$1 - \sigma_2(P)$ = spectral gap.

ω