

## Distributed/Multi-agent optimization

- we have a collection of  $N$  agents

- optimization problem:  $\min_{x \in \mathbb{R}^n} \sum_{i=1}^N f_i(x)$ ,

where  $f_i$  is the cost function of  $i$ -th agent

-  $f_i$  is only known to agent  $i$

- agents are unwilling to share cost function  $f_i$  with others.

- Example: Suppose we wish to solve:  $\min_{\omega} \sum_{i=1}^D (\omega^T x_i - y_i)^2$ ,

$(x_i, y_i)_{i=1}^D$  are collectively held by  $N$  agents.

agent 1 holds

$$(x_i, y_i)_{i=1}^{D_1}$$

agent 2 holds

$$(x_i, y_i)_{i=D_1+1}^{D_1+D_2}$$

:

agent  $N$  holds

$$(x_i, y_i)_{i=D_1+D_2+\dots+D_{N-1}+1}^D$$

$$\min_{\omega} \left[ \sum_{i=1}^{D_1} (\omega^T x_i - y_i)^2 + \sum_{i=D_1+1}^{D_1+D_2} (\omega^T x_i - y_i)^2 + \dots + \sum_{i=D_1+D_2+\dots+D_{N-1}+1}^D (\omega^T x_i - y_i)^2 \right]$$

$$f_1(\omega)$$

$$f_2(\omega)$$

Communication graph:  $G = (V, E)$

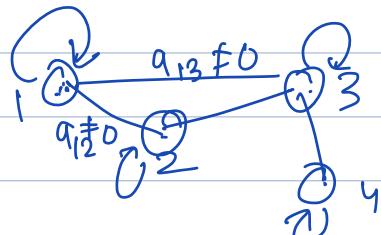
nodes  $\rightarrow$  edges

$|V| = N$ , each agent is a node

$(v_i, v_j) \in E$  if agent  $i$  &  $j$  can exchange information

the weight of this edge  $a_{ij}$

$$a_{ij} = \begin{cases} 0 & \text{if } (v_i, v_j) \notin E \\ > 0 & \text{if } (v_i, v_j) \in E \end{cases}$$



We will consider graphs that are connected & undirected.

$$a_{ij} = a_{ji}$$

for node  $i$ : we define  $N_i$  to be its set of neighbors

$$\underline{N_i} = \{v_j \in V \mid a_{ij} \neq 0\}$$

### Distributed Gradient Descent

- each agent initializes a solution

- at every time  $t = 0, 1, 2, \dots$

- gather:  $\underline{x_t^j}$  from  $j \in N_i$

$$\underline{v_{t+1}^i} = \boxed{\sum_{j \in N_i} a_{ij} x_t^j}$$

$$\underline{v_{t+1}^i} = \underline{v_{t+1}^i} - \eta_t \nabla f_i^p(v_{t+1}^i)$$



### Assumptions

i) on the network:  $a_{ii} \geq 0 \forall i$ ,  $\sum_{i=1}^N a_{ij} = 1 \forall j$ ,  $\sum_{j=1}^N a_{ij} = 1 \forall i$

$$A\mathbf{1} = \mathbf{1}$$

$$\mathbf{1}^T A = \mathbf{1}^T$$

$$A = \left[ \begin{array}{c|cc|c} a_{11} & a_{12} & \cdots & \\ \hline a_{21} & a_{22} & \cdots & \\ \hline a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right]$$

ii) step-sizes:  $\eta_t \geq 0$ ,  $\sum_{t=0}^{\infty} \eta_t = \infty$ ,  $\sum_{t=0}^{\infty} \eta_t^2 < \infty$

iii) on the cost functions:  

- each  $f_i$  is convex,  $\|\nabla f_i(x)\| \leq C_i \forall i$
- $\min_x \sum_{i=1}^N f_i(x)$  has at least one optimal solution.

Theorem: The solutions generated by the distributed GD algorithm converges to an optimal solution  $x^*$  for every agent, i.e.,

$$\lim_{t \rightarrow \infty} \|\underline{x}_t^i - x^*\| = 0 \quad \forall i \in \{1, 2, \dots, N\}.$$

Proof sketch:

i) consensus: let  $\bar{x}_t = \frac{1}{N} \sum_{i=1}^N x_t^i$

$$\lim_{t \rightarrow \infty} \|\underline{x}_t^i - \bar{x}_t\| = 0 \quad \forall i.$$

ii)  $\lim_{t \rightarrow \infty} \sum_{t=0}^T \eta_t \|\underline{x}_t^i - \bar{x}_t\| < \infty \quad \forall i$

iii)  $\lim_{t \rightarrow \infty} \|\bar{x}_t - x^*\| = 0.$

Note: we cannot use constant step-size because each agent is only computing partial gradient, rather than true gradient. This leads to slow convergence.

### Gradient Tracking Algorithm

- Initialize:  $x_0^i, y_0^i = \nabla f_i(x_0^i)$

$$x_{t+1} = x_t - \eta \underline{\nabla f(x_t)}$$

- At each time t

- gather:  $x_t^j$  from  $j \in N_i$

- updates:  $\underline{x}_{t+1}^i = \sum_{j \in N_i} a_{ij} x_t^j - \eta \underline{y}_t^i$

- gather:  $\underline{y}_t^j$  from  $j \in N_i$

- update :-  $y_{t+1}^i = \sum_{j \in N_i} a_{ij} y_t^j + \left[ \nabla_{f_i}^p(\underline{x}_{t+1}^i) - \nabla_{f_i}^p(\underline{x}_t^i) \right]$

privately held info agent i

Assumptions:

i) each  $f_i$  is  $\frac{\beta}{\gamma}$  smooth & strongly convex

ii) step-size  $\eta < \bar{\eta}$ , where  $\bar{\eta} \in (0, \frac{N}{\beta})$

Theorem :  $\|\underline{x}_{t+1}^i - \underline{x}^*\| \leq \rho \|\underline{x}_t^i - \underline{x}^*\|$  for  $\rho \in (0, 1)$ .

$$\Rightarrow \|\underline{x}_t^i - \underline{x}^*\| \leq \rho^t \|\underline{x}_0^i - \underline{x}^*\|$$

Performance of these algorithms are evaluated as follows.

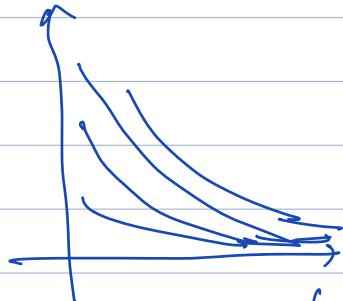
- consensus :  $\bar{e}_t^i = \|\underline{x}_t^i - \bar{x}_t\|$  vs.  $t$

- optimal value :  $|f(\bar{x}_t) - f(\underline{x}^*)|$  vs.  $t$

- optimal solution :  $\|\bar{x}_t - \underline{x}^*\|$  vs.  $t$

$$\max_i \|\underline{x}_t^i - \underline{x}^*\| \text{ vs. } t$$

- log-scale

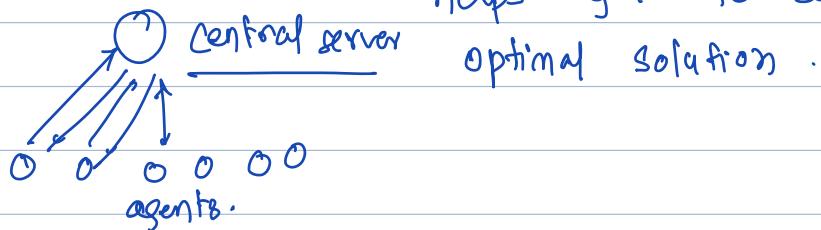


Many extensions :- time-varying graph, directed graph, asynchronous communication, malicious agents, accelerated schemes, and so on.

20th March

## Local Stochastic Gradient Descent (Local SGD)

- decentralized algorithm: there is a central server which helps agents to converge to the optimal solution.



- consider the problem:  $\min_{x \in \mathbb{R}^n} \frac{1}{N} \sum_{i=1}^N f_i(x)$ ,

where  $f_i(x) = \sum_{j=1}^{N_i} \underline{l_j^i(x, \xi_j)}$ ,

with  $(\xi_j)_{j=1}^{N_i}$  is the data held by agent  $i$

### Local SGD

- Let  $\mathcal{I}_T \subseteq \{1, 2, \dots, T\}$  be the time stamps at which communication takes place.
- every agent starts with  $x_0^i$
- at every time  $t$ :



$$x_{t+1}^i = \begin{cases} x_t^i - \eta_t \frac{\nabla l_{i_t}^i(x_t^i)}{v_t}, & \text{if } t \notin \mathcal{I}_T, \\ \frac{1}{N} \sum_{k=1}^N x_{t+1}^k - \eta_t \frac{\nabla l_{k_t}^x(x_{t+1}^k)}{v_t}, & \text{if } t \in \mathcal{I}_T \end{cases}, \quad i_t \in \{1, 2, \dots, N_i\}$$

- each agent  $i$  evaluates  $[\cdot]$
- communicates with central server
- server takes avg., & sends it back.

Theorem Let  $\eta^t = \frac{1}{\beta(a+t)}$ , where  $a$  is a constant &

$\beta$  is the smoothness parameter of the function.

Let the function be  $\alpha$ -strongly convex as well.

Then

$$\mathbb{E} \left[ f \left( \frac{1}{N} \sum_{i=1}^N x_i^t \right) \right] - f(x^*) \leq \mathcal{O} \left( \frac{1}{KT} \right),$$

where  $K = \frac{\beta}{\alpha}$  is the condition number.

① What to do when the function is not differentiable?

In those settings, we make use of subgradients.

Defn: For a function  $f$ , the set of subgradients at  $x_0$  is defined as

$$\partial f(x_0) = \left\{ g \in \mathbb{R}^n \mid f(y) \geq f(x_0) + g^T(y - x_0) \quad \forall y \in \mathbb{R}^n \right\}.$$

$\partial f(x_0)$  is called the subdifferential.

Example:

$$f(x) = |x|.$$

$$\partial f(1) = 1$$

$$\partial f(0) = \underline{\{[-1, 1]\}}$$

- In general  $\partial f(x_0)$  may be an empty-set.

- However, if the function is convex, then

$\partial f(x_0)$  is non-empty.

- If  $f$  is differentiable at  $x_0$ , then  $\partial f(x_0) = \{\nabla f(x_0)\}$ .

For such functions, we use subgradient (descent) algorithms:

$$x_{t+1} = x_t - \eta_t g_t, \text{ where } g_t \in \partial f(x_t).$$

Example:  $f(x) = \sum_{i=1}^N |x_i| = \|x\|_1, x \in \mathbb{R}^n$

find  $\partial f([0, 0, \dots, 0]) = \{g \in \mathbb{R}^n \mid g_i \in [-1, 1], \forall i\}$ .

we want to show that if  $w \in [-1, 1]^n$ , then

$$f(y) \geq f(x) + w^T(y-x) \quad \forall y \in \mathbb{R}^n \text{ at } x=0.$$

$$\Rightarrow \|y\|_1 \geq w^T y$$

from holder's inequality:  $w^T y \leq \|w\|_\infty \|y\|_1 \leq \|y\|_1$

$$\|w\|_\infty = \max_{1 \leq i \leq n} |w_i|$$

- These functions are not smooth, but may be strongly convex.
- Under the assumption that  $\|g(x)\| \leq G$ , then we can establish convergence of subgradient descent in an analogous manner to the convergence of gradient descent under bounded gradient assumption.

### Application: Optimal Detection.

Let  $\Theta = \{1, 2, \dots, m\}$  be the set of hypothesis.

$Y = \{1, 2, \dots, n\}$  be the set of observations.

eg:  $\Theta = \{\text{common cold, covid-19, covid-2, flu}\}$

$Y = \{\text{fever, sneezing, sore throat, fever+sore throat, fever+sneezing}\}$

For every hypothesis, we have a distribution over  $Y$ .

Let  $p_i$  denote the conditional distribution over  $Y$  when  $\Theta=i$

$$\underline{p_i} = \begin{bmatrix} P(Y=1 | \Theta=i) \\ P(Y=2 | \Theta=i) \\ \vdots \\ P(Y=n | \Theta=i) \end{bmatrix} \in \mathbb{R}^n, \quad p_i \geq 0, \quad \sum_i p_i = 1$$

are Known

$$P = [p_1, p_2, \dots, p_m]$$

Question: Given an observation  $Y$ , estimate (distribution) over  $\Theta$ .

For every observation, we have a different distribution over  $\Theta$ .

Let the distribution for  $Y=i$  be denoted as  $t_i \in \mathbb{R}^m$

$$\underline{t_i} = \begin{bmatrix} P(\hat{\Theta}=1 | Y=i) \\ P(\hat{\Theta}=2 | Y=i) \\ \vdots \\ P(\hat{\Theta}=m | Y=i) \end{bmatrix}, \quad t_i \geq 0, \quad \sum_i t_i = 1$$

$$T = [t_1, t_2, \dots, t_n]$$

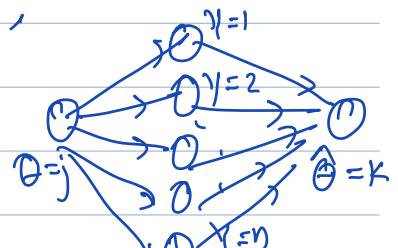
we need to determine  $(t_1, t_2, \dots, t_n)$ ,

How to evaluate the performance of  $t_i$ ?

Let the ground truth be  $\Theta=j$ .

we want  $(t_1, t_2, \dots, t_n)$  be such that

$$\begin{aligned} & P(\hat{\Theta}=j | \Theta=j) \text{ is high} \\ & P(\hat{\Theta} \neq j | \Theta=j) \text{ is low.} \end{aligned}$$



$$\begin{aligned} P(\hat{\Theta}=k | \Theta=j) &= \sum_{i=1}^N P(Y=i | \Theta=j) \cdot P(\hat{\Theta}=k | Y=i) \\ &= \sum_{i=1}^N (p_j)_i \cdot (t_{i,k})_k \end{aligned}$$

$$= \underbrace{[TP]_{kj}}_{\text{multiplication of } k^{\text{th}} \text{ row of } T \text{ with } j^{\text{th}} \text{ column of } P}$$

Thus, the matrix  $TP$  has the property that

$$\underbrace{[TP]_{ij}}_{\substack{\text{detection} \\ \text{matrix}}} = P(\hat{\theta} = i | \Theta = j)$$

we want  $\underbrace{[TP]_{ii}}_{\text{should be high}} \text{ with } i \neq j \text{ should be small.}$

Note that:  $[TP]_{ii} + \sum_{j \neq i} [TP]_{ij} = 1$

$$\Rightarrow 1 - [TP]_{ii} = \sum_{j \neq i} [TP]_{ij}$$

We can now formulate the following optimization problem.

$$\begin{aligned} & \min_{t_1, t_2, \dots, t_n} \max_{1 \leq i \leq n} (1 - [TP]_{ii}) \\ \text{s.t. } & \boxed{t_i \geq 0, \sum t_i = 1} \quad \forall i = 1, 2, \dots, n \\ & \text{min-max detector} \end{aligned}$$

Is the above problem convex??

- constraints: convex set

- cost function: it needs to be shown that  $t_i$ ,  $1 - [TP]_{ii}$  is a convex function.

$$\begin{aligned}
 f_i(t_1, t_2, \dots, t_n) &= 1 - [TP]_{ii} \\
 &= 1 - \sum_{k=1}^n t_{ik} p_{ki} \quad : \text{affine in } t_i
 \end{aligned}$$

$$\left( T_X(y) - y \right)^T (y - x) \leq 0 \quad \forall x$$