

Towards better understanding of decentralized optimization using row and column stochastic matrices

Liyuan Liang, Gan Luo

1 Related Works

To be filled.

1.1 Push-Pull Algorithm

$$\mathbf{x}^{(k+1)} = A\mathbf{x}^{(k)} - \alpha\mathbf{y}^{(k)} \quad (1)$$

$$\mathbf{y}^{(k+1)} = B\mathbf{y}^{(k)} + \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)} \quad (2)$$

2 New Approach

Consider using $\bar{x}^{(k)} := \frac{1}{n}\mathbf{1}_n^T \mathbf{x}^{(k)}$ and $\bar{\mathbf{x}}^{(k)} := \frac{1}{n}\mathbf{1}_n \mathbf{1}_n^T \mathbf{x}^{(k)}$ as the true parameter.

2.1 Descent Lemma

Lemma 1.

$$f(\bar{x}^{(k+1)}) \leq f(\bar{x}^{(k)}) - \frac{\alpha}{4}\|\bar{y}^{(k)}\|^2 - \frac{\alpha}{4}\|\nabla f(\bar{x}^{(k)})\|^2 + \left(\frac{\alpha L^2}{2} + \frac{2}{\alpha n^2}\right)\|\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)}\|^2 \quad (3)$$

Proof. Notice that

$$\bar{x}^{(k+1)} = \bar{x}^{(k)} + \frac{1}{n}\mathbf{1}_n^T A(\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)}) - \alpha\bar{y}^{(k)} \quad (4)$$

By L -smooth inequality, we have

$$f(y) \leq f(x) + \langle y - x, \nabla f(x) \rangle + \frac{L}{2}\|y - x\|^2 \quad (5)$$

Thus, using the smoothness of $f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x)$, we have

$$\begin{aligned}
f(\bar{x}^{(k+1)}) &\leq f(\bar{x}^{(k)}) - \alpha \langle \bar{y}^{(k)}, \nabla f(\bar{x}^{(k)}) \rangle + \langle \frac{1}{n} \mathbf{1}_n^T A(\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)}), \nabla f(\bar{x}^{(k)}) \rangle + \frac{\alpha^2 L}{2} \|\bar{y}^{(k)}\|^2 \\
&\stackrel{A-M}{\leq} f(\bar{x}^{(k)}) - \frac{\alpha - \alpha^2 L}{2} \|\bar{y}^{(k)}\|^2 - \frac{\alpha}{2} \|\nabla f(\bar{x}^{(k)})\|^2 + \frac{\alpha}{2} \|\bar{y}^{(k)} - \nabla f(\bar{x}^{(k)})\|^2 \\
&\quad + \frac{\alpha}{4} \|\nabla f(\bar{x}^{(k)})\|^2 + \frac{2}{\alpha} \|\frac{1}{n} \mathbf{1}_n^T A(\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)})\|^2 \\
&\stackrel{\alpha \leq \frac{1}{2L}}{\leq} f(\bar{x}^{(k)}) - \frac{\alpha}{4} \|\bar{y}^{(k)}\|^2 - \frac{\alpha}{4} \|\nabla f(\bar{x}^{(k)})\|^2 + \frac{\alpha}{2} \|\bar{y}^{(k)} - \nabla f(\bar{x}^{(k)})\|^2 \\
&\quad + \frac{2}{\alpha n^2} \|\mathbf{1}_n^T A(\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)})\|^2 \\
&= f(\bar{x}^{(k)}) - \frac{\alpha}{4} \|\bar{y}^{(k)}\|^2 - \frac{\alpha}{4} \|\nabla f(\bar{x}^{(k)})\|^2 + \frac{\alpha}{2n} \left\| \sum_{i=1}^n (\nabla f_i(x_i^{(k)}) - \nabla f_i(\bar{x}^{(k)})) \right\|^2 \\
&\quad + \frac{2}{\alpha n^2} \|\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)}\|^2 \\
&\leq f(\bar{x}^{(k)}) - \frac{\alpha}{4} \|\bar{y}^{(k)}\|^2 - \frac{\alpha}{4} \|\nabla f(\bar{x}^{(k)})\|^2 + (\frac{\alpha L^2}{2} + \frac{2}{\alpha n^2}) \|\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)}\|^2
\end{aligned} \tag{6}$$

This suffices to controlling the size of consensus error. \square

2.2 Consensus Error

Lemma 2.

$$\sum_{k=0}^T \|(I - R)\mathbf{x}^{(k+1)}\|^2 \leq \frac{\alpha^2 \kappa_P^2}{1 - \beta} \sum_{k=0}^T \beta^{2k} \sum_{i=0}^k \frac{1}{\beta^{3i}} \|(I - R)\mathbf{y}^{(i)}\|_F^2 \tag{7}$$

Proof. The consensus error can be expressed by $(I - R)\mathbf{x}^{(k)}$ where $R := \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$.

$$(I - R)\mathbf{x}^{(k+1)} = (A - RA)(I - R)\mathbf{x}^{(k)} - \alpha(I - R)\mathbf{y}^{(k)} \tag{8}$$

Notice that $(A - RA)^k(I - R) = (I - R)(A - A_\infty)^k(I - R)$, its size decays exponentially fast. The same for the consensus error.

Since $A_\infty \cdot A = A \cdot A_\infty$, we can diagonalize A and A_∞ at the same time.

$$A = P \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & \beta & 0 & \dots & 0 \\ 0 & 0 & \lambda_3 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \dots & \lambda_n \end{bmatrix} P^{-1}, A_\infty = P \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} P^{-1}$$

Then we have $\|(A - \mathbf{1}_n \pi_l^T)^{k-i}\|_2^2 \leq \kappa_P^2 \beta^{2(k-i)}$.

So we have:

$$\begin{aligned}
& \|(I - R)\mathbf{x}^{(k+1)}\|_F^2 \\
&= \alpha^2 \left\| \sum_{i=0}^k (I - R)(A - A_\infty)^{k-i} (I - R)\mathbf{y}^{(i)} \right\|_F^2 \\
&= \alpha^2 \left\| \sum_{i=0}^k \frac{(1 - \beta)\beta^i}{(1 - \beta)\beta^i} (I - R)(A - A_\infty)^{k-i} (I - R)\mathbf{y}^{(i)} \right\|_F^2 \\
&\stackrel{\text{Jensen}}{\leq} \alpha^2 \sum_{i=0}^k (1 - \beta)\beta^i \left\| \frac{1}{(1 - \beta)\beta^i} (I - R)(A - A_\infty)^{k-i} (I - R)\mathbf{y}^{(i)} \right\|_F^2 \\
&= \alpha^2 \sum_{i=0}^k \frac{1}{(1 - \beta)\beta^i} \|(I - R)(A - A_\infty)^{k-i} (I - R)\mathbf{y}^{(i)}\|_F^2 \\
&\leq \alpha^2 \sum_{i=0}^k \frac{1}{(1 - \beta)\beta^i} \|I - R\|_2^2 \cdot \|A - A_\infty\|_2^{2(k-i)} \cdot \|(I - R)\mathbf{y}^{(i)}\|_F^2 \\
&\leq \alpha^2 \sum_{i=0}^k \frac{1}{(1 - \beta)\beta^i} \cdot \|A - A_\infty\|_2^{2(k-i)} \cdot \|(I - R)\mathbf{y}^{(i)}\|_F^2 \\
&\leq \alpha^2 \sum_{i=0}^k \frac{1}{(1 - \beta)\beta^i} \cdot \kappa_P^2 \beta^{2(k-i)} \cdot \|(I - R)\mathbf{y}^{(i)}\|_F^2 \\
&= \frac{\alpha^2 \kappa_P^2 \beta^{2k}}{1 - \beta} \sum_{i=0}^k \frac{1}{\beta^{3i}} \|(I - R)\mathbf{y}^{(i)}\|_F^2
\end{aligned} \tag{9}$$

Then we have,

$$\sum_{k=0}^T \|(I - R)\mathbf{x}^{(k+1)}\|_F^2 \tag{10}$$

$$\leq \sum_{k=0}^T \frac{\alpha^2 \kappa_P^2 \beta^{2k}}{1 - \beta} \sum_{i=0}^k \frac{1}{\beta^{3i}} \|(I - R)\mathbf{y}^{(i)}\|_F^2 \tag{11}$$

$$= \frac{\alpha^2 \kappa_P^2}{1 - \beta} \sum_{k=0}^T \beta^{2k} \sum_{i=0}^k \frac{1}{\beta^{3i}} \|(I - R)\mathbf{y}^{(i)}\|_F^2 \tag{12}$$

Here $\kappa_P = \|P\|_2^2 \|P^{-1}\|_2^2$ and β is the second largest eigenvalue. \square

2.3 Gradient Consensus Error

Lemma 3.

$$\sum_{k=0}^T \|(I - R)\mathbf{y}^{(k+1)}\|^2 \leq \frac{3\kappa_B^2 L^2 (2\|A - RA\|_2^2 + 1 + 2n)}{(1 - \beta)^2} \sum_{k=0}^T \|(I - R)\mathbf{x}^{(k)}\|^2 + \frac{6n^2 \alpha^2 \kappa_B^2 L^2}{1 - \beta} \sum_{k=0}^T \|\bar{\mathbf{y}}^{(k)}\|^2 \tag{13}$$

[LLY: hi, Gan Luo, my proof here is correct but its technique may be not the best when we consider stochastic gradient. In stochastic case, using this proof, the order of $1 - \beta$ may be higher than 2. Can you follow the proof in my paper (Lemma B.2-Lemma B.5) to provide a new proof in stochastic case? (and compare which one provides a better constant)]

Proof. By Cauchy inequality, (8) indicates that

$$\|(I - R)\mathbf{x}^{(k+1)}\|^2 \leq 2\|A - RA\|_2^2 \|(I - R)\mathbf{x}^{(k)}\|^2 + 2\alpha^2 \|(I - R)\mathbf{y}^{(k)}\|^2 \quad (14)$$

By Cauchy inequality, (4) indicates that

$$\|R(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})\|^2 \leq \frac{2}{n} \|\mathbf{1}_n^T A(I - R)\mathbf{x}^{(k)}\|^2 + 2n\alpha^2 \|\bar{y}^{(k)}\|^2 \leq 2n \|(I - R)\mathbf{x}^{(k)}\|^2 + 2n\alpha^2 \|\bar{y}^{(k)}\|^2 \quad (15)$$

[LLY: Please use this to do inequalities.]

$$\begin{aligned} (I - R)\mathbf{y}^{(k+1)} &= (I - B_\infty)\mathbf{y}^{(k+1)} + (n\pi - \mathbf{1}_n)\bar{y}^{(k+1)} \\ &= \sum_{i=0}^k (B - B_\infty)^{k-i} (I_n - B_\infty)(\mathbf{g}^{(i+1)} - \mathbf{g}^{(i)}) + (n\pi - \mathbf{1}_n)\bar{y}^{(k+1)} \end{aligned} \quad (16)$$

$$\begin{aligned} &\|(I - R)\mathbf{y}^{(k+1)}\|^2 \\ &= \|(B - R)\mathbf{y}^{(k)} + (I - R)(\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)})\|^2 \\ &= \left\| \sum_{i=0}^k (B - B_\infty)^{k-i} (I - R)(\mathbf{g}^{(i+1)} - \mathbf{g}^{(i)}) \right\|^2 \\ &\stackrel{\text{Jensen}}{\leq} \frac{\kappa_B^2}{1 - \beta} \sum_{i=0}^k \beta^{k-i} \|\mathbf{g}^{(i+1)} - \mathbf{g}^{(i)}\|^2 \\ &\leq \frac{\kappa_B^2 L^2}{1 - \beta} \sum_{i=0}^k \beta^{k-i} \|\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)}\|^2 \\ &\leq \frac{3\kappa_B^2 L^2}{1 - \beta} \sum_{i=0}^k \beta^{k-i} (\|(I - R)\mathbf{x}^{(i+1)}\|^2 + \|R(\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)})\|^2 + \|(I - R)\mathbf{x}^{(i)}\|^2) \\ &\stackrel{(14), (15)}{\leq} \frac{3\kappa_B^2 L^2 (2\|A - RA\|_2^2 + 1 + 2n)}{1 - \beta} \sum_{i=0}^k \beta^{k-i} \|(I - R)\mathbf{x}^{(i)}\|^2 + \frac{6n^2 \alpha^2 \kappa_B^2 L^2}{1 - \beta} \sum_{i=0}^k \beta^{k-i} \|\bar{y}^{(i)}\|^2 \end{aligned} \quad (17)$$

□

2.4 Main Theorem

For non-stochastic case When $\alpha = \mathcal{O}(\frac{1}{L})$ is small enough, we have

Theorem 1.

$$\frac{1}{T+1} \sum_{k=0}^T \|\nabla f(\bar{x}^{(k)})\|^2 \leq \frac{4(f_0 - f^*)}{\alpha(T+1)} \quad (18)$$

[LLY: Please notice that we have not added stochastic noise. When there is some noise, we can show linear speedup easily.]

Proof.

□

For stochastic gradient case, we when α satisfies some condition, we have

Theorem 2.

$$\frac{1}{T+1} \sum_{k=0}^T \|\nabla f(\bar{x}^{(k)})\|^2 \leq \left(\frac{4\sigma^2(f_0 - f^*)}{n(T+1)} \right)^{\frac{1}{2}} + \text{network influence} \quad (19)$$

References