

**Обобщенная теория децентрализованного
стохастического градиентного спуска с
изменяющейся топологией и локальными шагами**

Общероссийский семинар по оптимизации

1 июля 2020

Machine Learning and Optimization Lab

EPFL



<https://www.epfl.ch/labs/mlo/>

Обобщенная теория децентрализованного стохастического градиентного спуска с изменяющейся топологией и локальными шагами

A Unified Theory of Decentralized SGD with Changing Topology and Local Updates

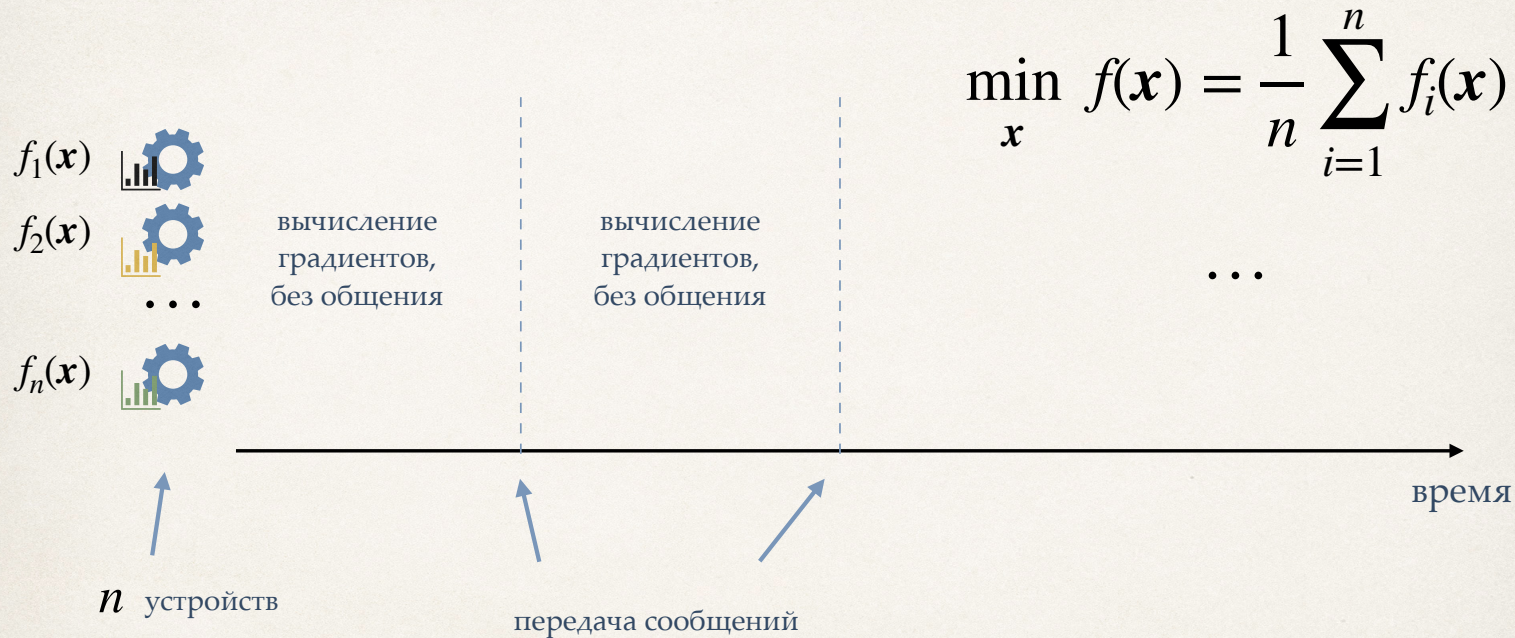
Anastasia Koloskova ^{*1} Nicolas Loizou ² Sadra Boreiri ¹ Martin Jaggi ¹ Sebastian U. Stich ^{*1}

Abstract

Decentralized stochastic optimization methods have gained a lot of attention recently, mainly because of their cheap per iteration cost, data locality, and their communication efficiency. In this

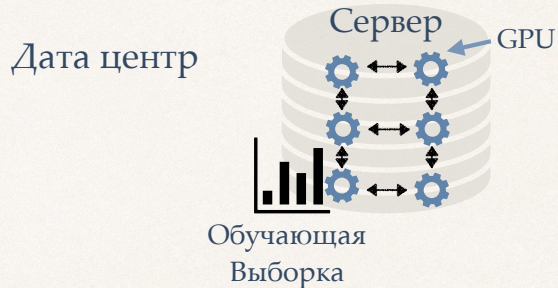
recently—though yet still at a smaller scale than federated learning (Lian et al., 2017; Assran et al., 2019; Koloskova et al., 2020). However, the community has identified a host of challenges that come along with decentralized training: notably, high communication cost (Tang et al., 2018a; Wang

Параллельные методы оптимизации

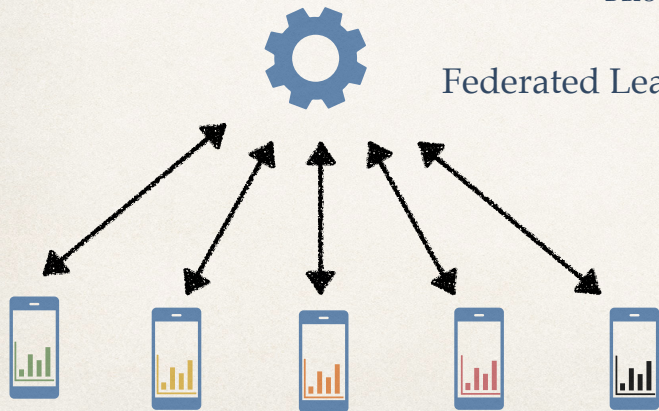


Централизованные

Децентрализованные



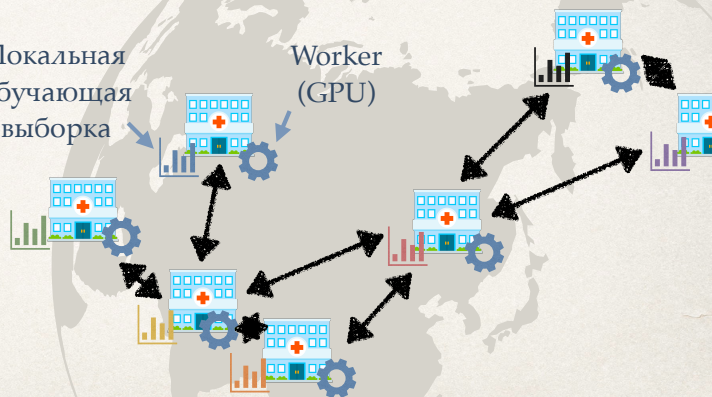
Federated Learning



Нет центрального устройства
Локальная передача данных
Распределенная выборка
(конфиденциальность&память)

Локальная
обучающая
выборка

Worker
(GPU)



Вариации

Разные статьи используют разную технику доказательства,
не все поддерживают non-iid распределенную выборку

Local SGD

(Li et al. 2020), (Khaled et al. 2020), (Wang&Joshi 2018),
(Yu et al. 2019), (Basu et al. 2019), (Patel&Dieuleveut 2019),
(Stich & Karimireddy 2019), (Li et al. 2019)

Decentralized Local SGD

(Wang&Joshi 2018), (Li et al. 2019)

Decentralized SGD

(Lian et al. 2017), (Wang&Joshi 2018),
(Olshevsky et al. 2019), (Koloskova et al. 2019),
(Li et al. 2019)

Decentralized SGD с изменяющейся топологией

(Nedic&Olshevsky 2014), (Wang et al. 2019)

and others

Основные результаты

- ✧ Мы анализируем все частные случаи в одном фреймворке
 - ✧ Покрываем non-iid data
- ✧ Скорость сходимости либо совпадает, либо улучшает предыдущие результаты
 - ✧ Предоставляем нижнюю оценку скорости сходимости
- ✧ Может помочь понять недостатки алгоритмов и для дизайна новых методов

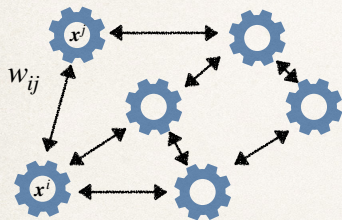
Decentralized SGD

устройство i



$$\mathbf{x}_i^{(t+\frac{1}{2})} := \mathbf{x}_i^{(t)} - \gamma_t \nabla F_i(\mathbf{x}_i^{(t)}, \xi_{j_t}) \quad \leftarrow \text{шаг SGD}$$

стохастический градиент



$$\mathbf{x}_i^{(t+1)} := \sum_j w_{ij} \mathbf{x}_j^{(t+\frac{1}{2})} \quad \leftarrow \text{шаг усреднения}$$

Например, для централизованной системы:

$$\mathbf{x}^{(t+1)} := \frac{1}{n} \sum_j \mathbf{x}_j^{(t+\frac{1}{2})} = \mathbf{x}^{(t)} - \gamma_t \frac{1}{n} \sum_j \nabla F_i(\mathbf{x}^{(t)}, \xi_{j_t})$$

Decentralized SGD

Mixing matrix:

$$W = \{w_{ij}\}$$

$$X^{(t)} = [\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_n^{(t)}]$$

$$\mathbf{x}_i^{(t+\frac{1}{2})} := \mathbf{x}_i^{(t)} - \gamma_t \nabla F_i(\mathbf{x}_i^{(t)}, \xi_{j_t})$$

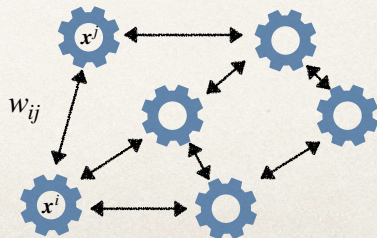
$$\mathbf{x}_i^{(t+1)} := \sum_j w_{ij} \mathbf{x}_j^{(t+\frac{1}{2})}$$

Шаг усреднения также можно записать как:

$$X^{(t+1)} := X^{(t+\frac{1}{2})} W$$

$$W = \{w_{ij}\}$$

- * Дважды стохастическая
- * Симметричная



Decentralized SGD

❖ Шаг усреднения действительно усредняет

$$\mathbb{E} \|XW - \bar{X}\|^2 \leq (1 - \underline{p}) \|X - \bar{X}\|^2$$
$$p \in [0, 1]$$

$$1 = |\lambda_1(W)| > |\lambda_2(W)| \geq \dots$$

$$\delta = 1 - |\lambda_2(W)| \text{ — spectral gap}$$

$$\mathbf{x}_i^{(t+\frac{1}{2})} := \mathbf{x}_i^{(t)} - \gamma_t \nabla F_i(\mathbf{x}_i^{(t)}, \xi_{j_t})$$

$$\mathbf{x}_i^{(t+1)} := \sum_j w_{ij} \mathbf{x}_j^{(t+\frac{1}{2})}$$

$$X^{(t)} = [\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_n^{(t)}]$$

$$\bar{X}^{(t)} = [\bar{\mathbf{x}}^{(t)}, \dots, \bar{\mathbf{x}}^{(t)}]$$

Decentralized SGD

$$\mathbb{E} \|XW - \bar{X}\|^2 \leq (1-p) \|X - \bar{X}\|^2$$

$$\|XW - \bar{X}\|_2^2 = \left\| \underbrace{(X - \bar{X})}_{\bar{X}} \underbrace{W - \frac{11^T}{n}}_0 \right\|_2^2 \quad \boxed{X^{(t)}} = [x_1^{(t)}, \dots, x_n^{(t)}]$$

$$\bar{X} = X \left(\frac{11^T}{n} \right)$$

$$\bar{X}^{(t)} = [\bar{x}^{(t)}, \dots, \bar{x}^{(t)}]$$

$$\bar{X} = \bar{X} \cdot W \quad = \left\| (X - \bar{X}) \left(W - \frac{11^T}{n} \right) \right\|_2^2 =$$

$$\leq \|X - \bar{X}\|_2^2 \underbrace{\left\| W - \frac{11^T}{n} \right\|_2^2}_{\lambda_2(W)} = \lambda_2(W) \|X - \bar{X}\|^2$$

$$(1 - \underline{\delta})^2 = (1 - p)$$

$$\boxed{p} = 1 - (1 - \underline{\delta})^2 = 1 - 1 + 2\underline{\delta} - \underline{\delta}^2 = 2\underline{\delta} - \underline{\delta}^2 \underset{\underline{\delta} \rightarrow 0}{\approx} 2\underline{\delta} = \Theta(\underline{\delta})$$

Decentralized SGD

На практике часто выбирают:

- ❖ Если все степени вершин одинаковые,

$$w_{ij} = \frac{1}{deg_i + 1}$$

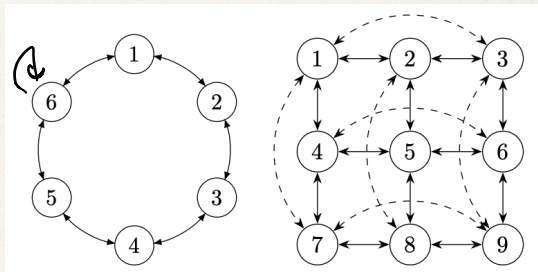
- ❖ Если разные, Metropolis-Hastings правило

$$w_{ij} = \frac{1}{\max(deg_i, deg_j) + 1}$$

$$w_{ii} = 1 - \sum w_{ij}$$

$$\mathbf{x}_i^{(t+\frac{1}{2})} := \mathbf{x}_i^{(t)} - \gamma_t \nabla F_i(\mathbf{x}_i^{(t)}, \xi_{j_t})$$

$$\mathbf{x}_i^{(t+1)} := \sum_j w_{ij} \mathbf{x}_j^{(t+\frac{1}{2})}$$



graph/topology	δ^{-1}	node degree
ring	$\mathcal{O}(n^2)$	2
2d-torus	$\mathcal{O}(n)$	4
fully connected	$\mathcal{O}(1)$	$n - 1$

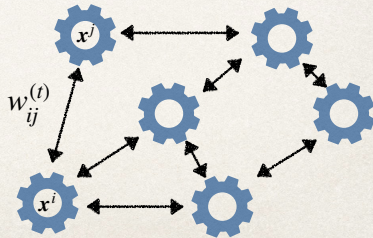
Обобщенный decentralized SGD

$$\mathbf{x}_{t+\frac{1}{2}}^i := \mathbf{x}_t^i - \gamma_t \nabla F_i(\mathbf{x}_t^i, \xi_{j_t}) \quad \leftarrow \text{шаг SGD}$$

стохастический градиент

$$\mathbf{x}_{t+1}^i := \sum_j w_{ij}^{(t)} \mathbf{x}_{t+\frac{1}{2}}^j \quad \leftarrow \text{шаг усреднения с весами } \{w_{ij}^{(t)}\} \sim \mathcal{W}^{(t)}$$

- ❖ Возможно менять топологию сети на каждой итерации
- ❖ Разрешено случайным образом выбирать веса из какого-то распределения



Предположение

$$W^{(t)} = \{w_{ij}^{(t)}\} \quad \clubsuit \text{ Дважды стохастическая}$$

$$\mathbf{x}_{t+\frac{1}{2}}^i := \mathbf{x}_t^i - \gamma_t \nabla F_i(\mathbf{x}_t^i, \xi_{j_t}^i)$$

$$\mathbf{x}_{t+1}^i := \sum_j w_{ij}^{(t)} \mathbf{x}_{t+\frac{1}{2}}^j$$

$$\underbrace{W^{(l+1)\tau-1} \dots W^{l\tau}}_{W_{l,\tau}} \dots \underbrace{W^{(2\tau-1)} \dots W^{(\tau)}}_{W_{1,\tau}} \underbrace{W^{(\tau-1)} \dots W^{(0)}}_{W_{0,\tau}}$$

\clubsuit Нужно только что комбинация каждых τ шагов усредняет, а не каждый отдельный шаг

$$\mathbb{E} \|XW_{l,\tau} - \bar{X}\|^2 \leq (1-p) \|X - \bar{X}\|^2 \quad \forall l$$

Пример: Decentralized Local SGD

$$\underbrace{WI \dots I}_{W_{l,\tau}} \cdot \dots \cdot \underbrace{WI \dots I}_{W_{1,\tau}} \underbrace{WI \dots I}_{W_{0,\tau}}$$

I — единичная матрица,
без передачи данных

Сходимость

$$f_i(x) = \mathbb{E}_{\xi} F_i(x, \xi)$$

$$f_i(x) - f_i(y) + \frac{\mu}{2} \|x - y\|_2^2 \leq \langle \nabla f_i(x), x - y \rangle$$

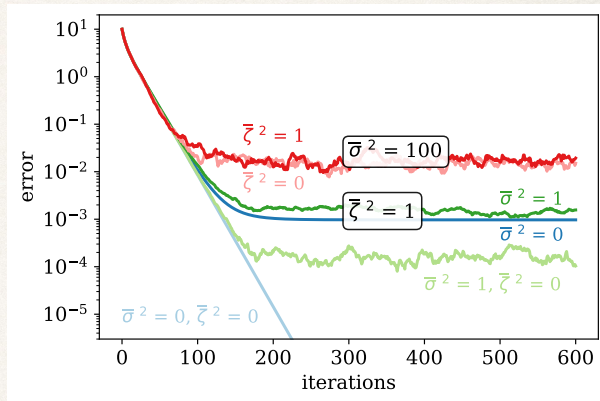
$$\|\nabla F_i(y, \xi) - \nabla F_i(x, \xi)\| \leq L\|y - x\|$$

Сильно выпуклый случай

$$\tilde{\mathcal{O}} \left(\underbrace{\frac{\bar{\sigma}^2}{nT}}_{\text{стохастические члены}} + \underbrace{\frac{\tau^2 \bar{\zeta}^2 + \tau p \bar{\sigma}^2}{p^2 T^2}}_{\text{оптимизационные члены}} + \underbrace{\frac{\tau}{p} \exp \left[-\frac{Tp}{\tau} \right]}_{\text{оптимизационные члены}} \right)$$

стохастические члены

оптимизационные члены



$$\bar{\sigma}^2 = \frac{1}{n} \sum_i \sigma_i^2, \quad \sigma_i^2 := \mathbb{E}_{\xi_i} \|\nabla F_i(\mathbf{x}^*, \xi_i) - \nabla f_i(\mathbf{x}^*)\|_2^2$$

дисперсия градиентов

$$\bar{\zeta}^2 = \frac{1}{n} \sum_i \zeta_i^2,$$

$$\zeta_i^2 := \|\nabla f_i(\mathbf{x}^*)\|_2^2$$

разница функций

Сходимость: идея доказательства

$$\eta_t = \frac{1}{t}$$

Лемма 1

$$\mathbb{E} \|\bar{\mathbf{x}}^{(t+1)} - \mathbf{x}^{\star}\|^2 \leq \left(1 - \frac{\mu\eta_t}{2}\right) \mathbb{E} \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{\star}\|^2 - \eta_t(f(\bar{\mathbf{x}}^{(t)}) - f(\bar{\mathbf{x}}^{\star})) + \frac{\bar{\sigma}^2}{n}\eta_t^2 + \eta_t 3L\Xi_t$$

$$\Xi_t = \frac{1}{n} \mathbb{E}_t \sum_{i=1}^n \|\mathbf{x}_i^{(t)} - \bar{\mathbf{x}}^{(t)}\|^2$$

$$\tilde{\mathbf{W}} \sim \begin{cases} \mathbf{W} & \text{w.p. } \frac{1}{\epsilon} \\ \mathbf{I} & \text{w.p. } 1 - \frac{1}{\epsilon}. \end{cases}$$

Лемма 2

$$\Xi_t \leq \left(1 - \frac{p}{2}\right) \underbrace{\Xi_{m\tau}}_{\text{circled}} + \frac{p}{16\tau} \sum_{j=m\tau}^{t-1} \Xi_j + 36L \frac{\tau}{p} \sum_{j=m\tau}^{t-1} \eta_j^2 (f(\bar{\mathbf{x}}^{(t)}) - f(\bar{\mathbf{x}}^{\star})) + \left(\bar{\sigma}^2 + \frac{18\tau}{p} \bar{\zeta}^2\right) \sum_{j=m\tau}^{t-1} \eta_j^2$$

Нижняя оценка

$$\tilde{\mathcal{O}} \left(\frac{\bar{\sigma}^2}{nT} + \frac{\tau^2 \bar{\zeta}^2 + \tau p \bar{\sigma}^2}{p^2 T^2} + \frac{1}{p} \exp \left[-\frac{Tp}{\tau} \right] \right)$$

Сильно выпуклый случай

$$\bar{\sigma}^2 = 0$$

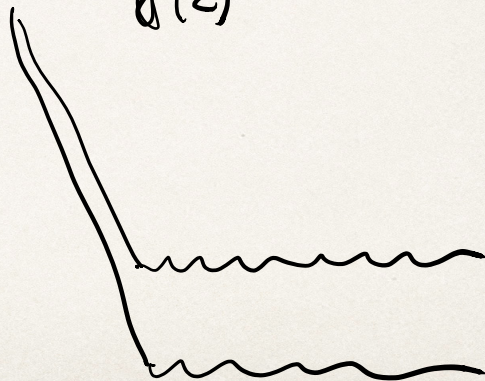
$$\tau = 1$$

$$\tilde{\Omega} \left(\frac{\bar{\zeta}^2 (1-p)}{p^2 T^2} \right)$$

Сублинейный член необходим

ε

$\chi(\varepsilon)$



Нижняя оценка: идея

Decentralized SGD с матрицей W

$$f_i(x) = \frac{1}{2}(x - y_i)^2 \quad x \in \mathbb{R}^1$$

$\mathbf{x}^{(0)} = (x_1^{(0)}, \dots, x_n^{(0)})^\top$ — собственный вектор W , отвечающий второму
собственному числу ρ

$$\mathbf{y} = \mathbf{x}^{(0)} + \mathbf{1}$$

Итерация равносильна:

$$\mathbf{x}^{(t+1)} = W\mathbf{x}^{(t+\frac{1}{2})} = W((1 - \gamma)\mathbf{x}^{(t)} + \gamma\mathbf{y}) =$$

$$\underline{W}^t (1 - \gamma)^t \underline{x}^{(0)} + \gamma \sum_{\tau=0}^{t-1} \underline{(1 - \gamma)^{t-\tau}} \underline{W}^{t-\tau} \mathbf{y}$$

$$\mathbf{x}_{t+\frac{1}{2}}^i := \mathbf{x}_t^i - \gamma_t \nabla F_i(\mathbf{x}_t^i, \xi_{j_t}^i)$$

$$\mathbf{x}_{t+1}^i := \sum_j w_{ij}^{(t)} \mathbf{x}_{t+\frac{1}{2}}^j$$

$$\bar{\sigma}^2 = 0$$

$$\tau = 1$$

$$\tilde{\Omega} \left(\frac{\bar{\zeta}^2 (1 - p)}{p^2 T^2} \right)$$

Частный случай: Local SGD



вычисление
градиентов,
без общения

$$W^{(t)} = I$$

$$W^{(t)} = I$$

...

$$W^{(t)} = \frac{11^T}{n}$$

(полный граф)

централизованная all-to-all агрегация

Наилучшая скорость
сходимости ранее:

$$\mathcal{O}\left(\frac{\bar{\sigma}^2}{nT} + \frac{\tau^2 \bar{\zeta}^2}{T}\right) \quad (\text{Li et al. 2020})$$

Эта работа:

$$\tilde{\mathcal{O}}\left(\frac{\bar{\sigma}^2}{nT} + \frac{\tau^2 \bar{\zeta}^2 + \tau \bar{\sigma}^2}{T^2} + \tau \exp\left[-\frac{T}{\tau}\right]\right)$$

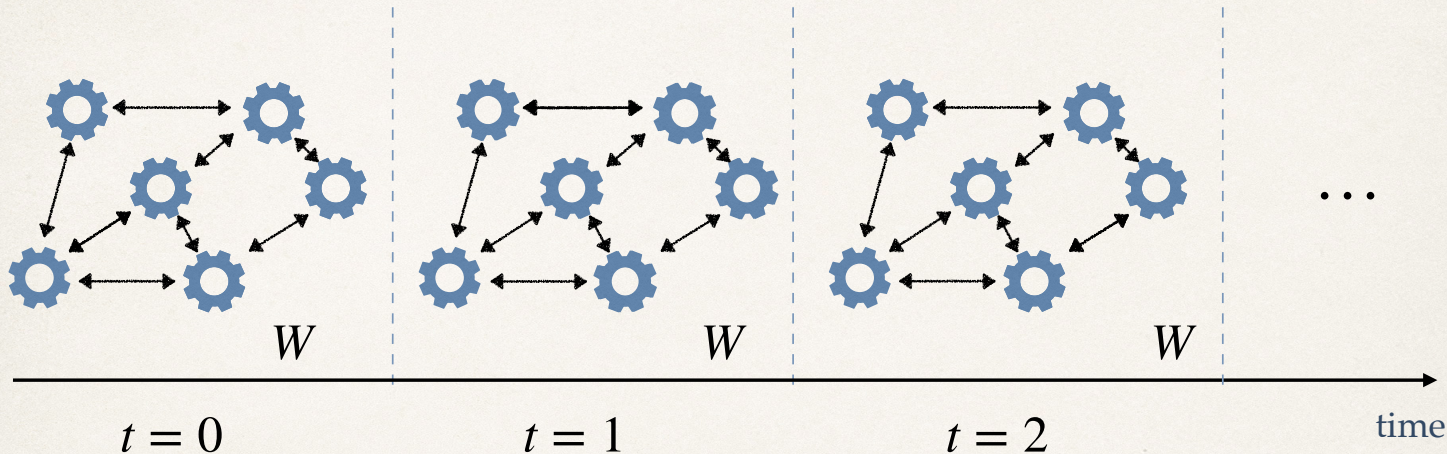
Скорость сходимости совпадает с
ранее известным результатом для
невыпуклых функций

time

Улучшает ранее известный результат
для выпуклых и сильно выпуклых
функций

Также в литературе называют FedAvg,
популярен в Federated Learning

Частный случай: Decentralized SGD



Наилучшая скорость
сходимости ранее:

$$\mathcal{O}\left(\frac{\bar{\sigma}^2}{nT} + \frac{G^2}{p^2T^2} + \frac{G^2}{p^3T^3}\right)$$

(Koloskova et al., 2019)

(Olshevsky et al., 2019)

Эта работа:

$$\tilde{\mathcal{O}}\left(\frac{\bar{\sigma}^2}{nT} + \frac{\bar{\zeta}^2 + p\bar{\sigma}^2}{p^2T^2} + \frac{1}{p} \exp[-Tp]\right)$$

$$G^2 \geq \bar{\zeta}^2 + \bar{\sigma}^2$$

Совпадает с ранее известным для
невыпуклых функций

Улучшает ранее известный результат
для выпуклых и сильно выпуклых
функций

Частный случай: Decentralized Local SGD



вычисление
градиентов,
без общения

$$W^{(t)} = I$$

$$W^{(t)} = I$$

...

Эта работа:

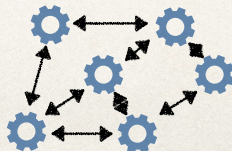
$$\tilde{\mathcal{O}} \left(\frac{\bar{\sigma}^2}{nT} + \frac{\tau^2 \bar{\zeta}^2 + \tau p \bar{\sigma}^2}{p^2 T^2} + \frac{\tau}{p} \exp \left[-\frac{Tp}{\tau} \right] \right)$$

τ

2τ

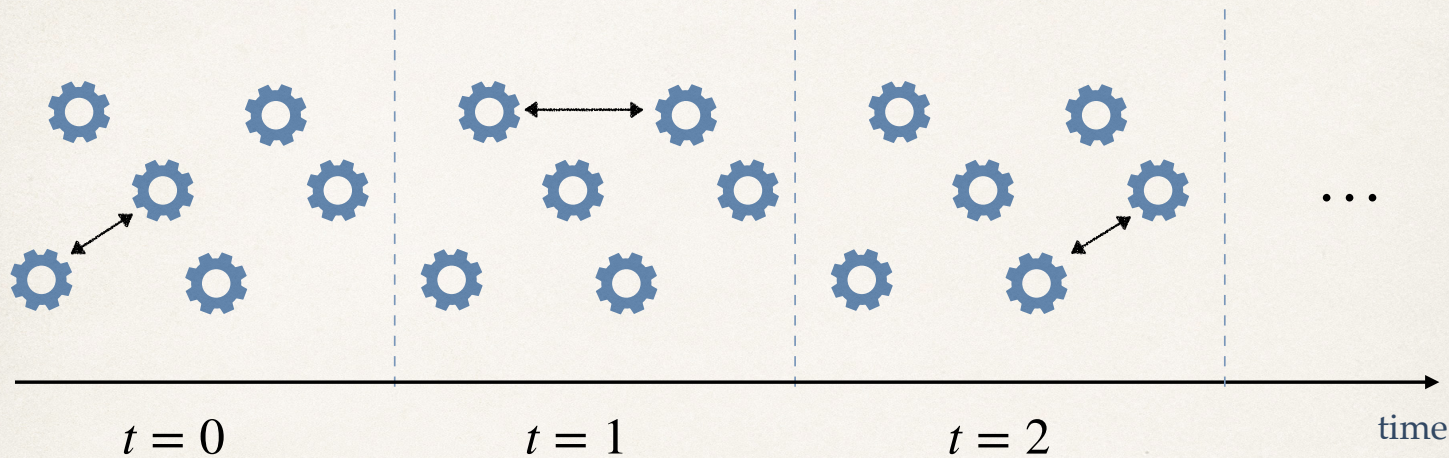
time

$$W^{(t)} = W$$



Улучшает известные результаты во
всех случаях

Частный случай: Pairwise randomised gossip



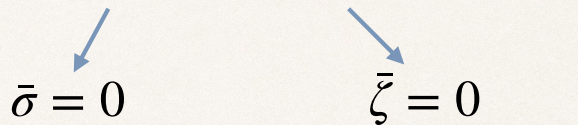
$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(\bar{x}^t) \right\|^2 \right]$$

Эта работа:

$$\frac{1}{T} \sum_{t=1}^T \underbrace{\mathbb{E} \left[f(\bar{x}^t) - f(x^*) \right]}_{\text{error}} \leq \tilde{\mathcal{O}} \left(\frac{\bar{\sigma}^2}{nT} + \frac{\bar{\zeta}^2 + p\bar{\sigma}^2}{p^2 T^2} + \frac{1}{p} \exp[-Tp] \right)$$

Частный случай: Overparametrized regime

В оптимуме, функция потерь равна нулю для каждого из обучающего примера

$$\nabla F_i(\mathbf{x}^\star, \xi) = 0 \quad \forall i, \xi$$

$$\bar{\sigma} = 0 \qquad \bar{\zeta} = 0$$

Линейная сходимость

$$\tilde{\mathcal{O}} \left(\frac{L\tau R_0^2}{p} \exp \left[-\frac{\mu T p}{\tau L} \right] \right)$$

Заключение

- ⚙ Предложен обобщенный фреймворк для параллельной оптимизации
 - ⚙ Вывели теоретическую скорость сходимости с non-iid data
- ⚙ Скорость сходимости совпадает и улучшает предыдущие результаты
 - ⚙ Вывели нижнюю оценку для decentralized SGD