

Постановка задачи и актуальность

Задача машинного обучения

- Сформулируем задачу машинного обучения:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{N} \sum_{i=1}^N \ell(x, z_i).$$

Задача машинного обучения

- Сформулируем задачу машинного обучения:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{N} \sum_{i=1}^N \ell(x, z_i).$$

- **Вопрос:** как решать? GD, Accelerated GD, SGD, Adam и так далее.
- Объемы данных растут, одно вычислительное устройство может долго считать даже стохастический градиент.
- **Вопрос:** что делать? Распараллеливать процесс обучения. Использовать несколько вычислительных устройств.
- **Вопрос:** как это сделать? Распределить данные между устройствами/агентами/нодами.

Задача машинного обучения

- Сформулируем задачу машинного обучения:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{N} \sum_{i=1}^N \ell(x, z_i).$$

- **Вопрос:** как решать? GD, Accelerated GD, SGD, Adam и так далее.
- Объемы данных растут, одно вычислительное устройство может долго считать даже стохастический градиент.
- **Вопрос:** что делать? Распараллеливать процесс обучения. Использовать несколько вычислительных устройств.
- **Вопрос:** как это сделать? Распределить данные между устройствами/агентами/нодами.

Задача машинного обучения

- Сформулируем задачу машинного обучения:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{N} \sum_{i=1}^N \ell(x, z_i).$$

- **Вопрос:** как решать? GD, Accelerated GD, SGD, Adam и так далее.
- Объемы данных растут, одно вычислительное устройство может долго считать даже стохастический градиент.
- **Вопрос:** что делать? Распараллеливать процесс обучения. Использовать несколько вычислительных устройств.
- **Вопрос:** как это сделать? Распределить данные между устройствами/агентами/нодами.

Федеративное обучение

- Вычислительные устройства – пользовательские устройства: ноутбуки, планшеты, телефоны. Неравномерность вычислительных мощностей.
- Данные часто сильно разнородны, как по размеру так и по природе. Дополнительные проблемы приватности.
- Желание пользователей часто отличается друг от друга и от желания владельца сервиса.

Федеративное обучение

- Вычислительные устройства – пользовательские устройства: ноутбуки, планшеты, телефоны. Неравномерность вычислительных мощностей.
- Данные часто сильно разнородны, как по размеру так и по природе. Дополнительные проблемы приватности.
- Желание пользователей часто отличается друг от друга и от желания владельца сервиса.

Коммуникации – главное узкое место

- Выигрываем в параллельности, но платим за это тратами на коммуникации. Хотелось бы плату эту уменьшить.
- Проблема присутствует во всех подходах от кластерного (с подключением по кабелю) до федеративного обучения (с неустойчивым интернетом).

Коммуникации – главное узкое место

- Выигрываем в параллельности, но платим за это тратами на коммуникации. Хотелось бы плату эту уменьшить.
- Проблема присутствует во всех подходах от кластерного (с подключением по кабелю) до федеративного обучения (с неустойчивым интернетом).

Организация коммуникаций

Типы коммуникационных архитектур: теория

- **Централизованная:** можем получить **точное** усреднение по всем устройствам (возможно, с задержками, сбоями и так далее)
- **Децентрализованная:** **точное** усреднение не предусмотрено

Типы коммуникационных архитектур: теория

- **Централизованная:** можем получить **точное** усреднение по всем устройствам (возможно, с задержками, сбоями и так далее)
- **Децентрализованная:** **точное** усреднение не предусмотрено

Типы коммуникационных архитектур: централизованная

- Посмотрим на примере, как обычный неопределенный GD становится централизованным.

Алгоритм 1 Централизованный GD

Вход: Размер шага $\gamma > 0$, стартовая точка $x_0 \in \mathbb{R}^d$, количество итераций K

```
1: for  $k = 0, 1, \dots, K - 1$  do
```

2: Отправить x_k всем рабочим

- ▷ выполняется сервером

3: **for** $i = 1, \dots, n$ параллельно **do**

4: Принять x_k от мастера

- ▷ выполняется рабочими

5: Вычислить градиент $\nabla f_m(x_k)$ в точке x_k

▷ выполняется рабочими

6: Отправить $\nabla f_m(x_k)$ мастеру

- ▷ выполняется рабочими

```
7:     end for
```

8: Принять $\nabla f_m(x_k)$ от всех рабочих

- ▷ выполняется сервером

9: Вычислить $\nabla f(x_k) = \frac{1}{M} \sum_{m=1}^M \nabla f_m(x_k)$

- ▷ выполняется сервером

10: $x_{k+1} = x_k - \gamma \nabla f(x_k)$

- ▷ выполняется сервером

```
11: end for
```

Выход: x^K

Типы коммуникационных архитектур: практика

- **Централизованная:** есть некоторая машина-сервер (мастер), с которой соединены все остальные устройства (рабочие).
- **Архитектура с AllReduce процедурой** (в теории это централизованная): задан некоторый граф связей/коммуникаций, обмен сообщениями происходит согласно этому графу, в том числе можно организовать усреднение (в теории это также централизованная архитектура).
- **Децентрализованная:** задан некоторый граф связей/коммуникаций, обмен сообщениями происходит согласно этому графу, точное усреднение не используется.

Типы коммуникационных архитектур: практика

- Централизованная: есть некоторая машина-сервер (мастер), с которой соединены все остальные устройства (рабочие).
- Архитектура с AllReduce процедурой (в теории это централизованная): задан некоторый граф связей/коммуникаций, обмен сообщениями происходит согласно этому графу, в том числе можно организовать усреднение (в теории это также централизованная архитектура).
- Децентрализованная: задан некоторый граф связей/коммуникаций, обмен сообщениями происходит согласно этому графу, точное усреднение не используется.

Ring AllReduce: второй шаг суммирования

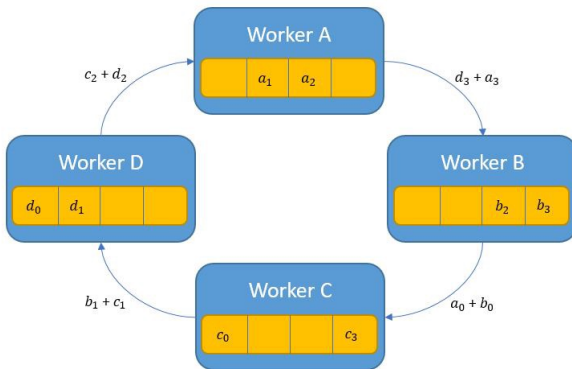


Figure: Картинка отсюда

Сжатие: несмещенные и смещенные операторы

Несмещённая компрессия: примеры

Трёхуровневая ℓ_2 -квантизация

Рассмотрим следующий оператор: $[Q(x)]_i = \|x\|_2 \text{sign}(x_i) \xi_i$, $i = 1, \dots, d$, где $[Q(x)]_i$ — i -я компонента вектора $Q(x)$ и ξ_i — случайная величина, имеющая распределение Бернулли с параметром $\frac{|x_i|}{\|x\|_2}$, то есть

$$\xi_i = \begin{cases} 1 & \text{с вероятностью } \frac{|x_i|}{\|x\|_2}, \\ 0 & \text{с вероятностью } 1 - \frac{|x_i|}{\|x\|_2}. \end{cases}$$

Таким образом, если мы хотим передать вектор $Q(x)$, то нам нужно передать вектор, состоящий из нулей и ± 1 , и вещественное число $\|x\|_2$, причём вероятность обнуления компоненты тем больше, чем компонента меньше по модулю по сравнению с остальными компонентами. Можно показать, что данный оператор является несмещённой компрессией с константой $\omega = \sqrt{d}$.



Alistarh D. et al. QSGD: Communication-Efficient SGD via Gradient Quantization and Encoding

Несмещенная компрессия: примеры

- Еще примеры несмещенных компрессией:



Beznosikov A. et al. On Biased Compression for Distributed Learning



Horváth S. et al. Natural compression for distributed deep learning



Szlendak R. et al. Permutation Compressors for Provably Faster Distributed Nonconvex Optimization

Несмещенная компрессия: идея

Несмещенная компрессия: идея

- Самая простая идея, которая приходит в голову, состоит в том, чтобы использовать параллельный GD, но к градиентам, пересылаемым от рабочих на сервер, применять несмещённую компрессию.

Квантизированный GD (QGD)

Алгоритм 1 QGD

Вход: размер шага $\gamma > 0$, стартовая точка $x_0 \in \mathbb{R}^d$, количество итераций K

```

1: for  $k = 0, 1, \dots, K - 1$  do
2:   Отправить  $x_k$  всем рабочим ▷ выполняется сервером
3:   for  $i = 1, \dots, n$  параллельно do
4:     Принять  $x_k$  от мастера ▷ выполняется рабочими
5:     Вычислить градиент  $\nabla f_m(x_k)$  в точке  $x_k$  ▷ выполняется рабочими
6:     Независимо сгенерировать  $g_{k,m} = \mathcal{Q}(\nabla f_m(x_k))$  ▷ выполняется рабочими
7:     Отправить  $g_{k,m}$  мастеру ▷ выполняется рабочими
8:   end for
9:   Принять  $g_{k,m}$  от всех рабочих ▷ выполняется сервером
10:  Вычислить  $g_k = \frac{1}{M} \sum_{m=1}^M g_{k,m}$  ▷ выполняется сервером
11:   $x_{k+1} = x_k - \gamma g_k$  ▷ выполняется сервером
12: end for
Выход:  $x^K$ 

```


Несмещенная компрессия: доказательство

- Работаем с $\mathbb{E}[g_k \mid x_k]$:

$$\begin{aligned}\mathbb{E}[g_k \mid x_k] &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[g_{k,m} \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\mathbb{E}[\mathcal{Q}(\nabla f_m(x_k)) \mid \nabla f_m(x_k)] \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\nabla f_m(x_k) \mid x_k] = \frac{1}{M} \sum_{m=1}^M \nabla f_m(x_k) = \nabla f(x_k).\end{aligned}$$

- Работаем с $\mathbb{E}[\|g_k\|^2 \mid x_k]$:

$$\mathbb{E}[\|g_k\|^2 \mid x_k] = \mathbb{E}\left[\left\|\frac{1}{M} \sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right] = \frac{1}{M^2} \mathbb{E}\left[\left\|\sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right].$$

Несмещенная компрессия: доказательство

- Работаем с $\mathbb{E}[g_k \mid x_k]$:

$$\begin{aligned}\mathbb{E}[g_k \mid x_k] &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[g_{k,m} \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\mathbb{E}[\mathcal{Q}(\nabla f_m(x_k)) \mid \nabla f_m(x_k)] \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\nabla f_m(x_k) \mid x_k] = \frac{1}{M} \sum_{m=1}^M \nabla f_m(x_k) = \nabla f(x_k).\end{aligned}$$

- Работаем с $\mathbb{E}[\|g_k\|^2 \mid x_k]$:

$$\mathbb{E}[\|g_k\|^2 \mid x_k] = \mathbb{E}\left[\left\|\frac{1}{M} \sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right] = \frac{1}{M^2} \mathbb{E}\left[\left\|\sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right].$$

Несмещенная компрессия: доказательство

- Работаем с $\mathbb{E}[g_k \mid x_k]$:

$$\begin{aligned}\mathbb{E}[g_k \mid x_k] &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[g_{k,m} \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\mathbb{E}[\mathcal{Q}(\nabla f_m(x_k)) \mid \nabla f_m(x_k)] \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\nabla f_m(x_k) \mid x_k] = \frac{1}{M} \sum_{m=1}^M \nabla f_m(x_k) = \nabla f(x_k).\end{aligned}$$

- Работаем с $\mathbb{E}[\|g_k\|^2 \mid x_k]$:

$$\mathbb{E}[\|g_k\|^2 \mid x_k] = \mathbb{E}\left[\left\|\frac{1}{M} \sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right] = \frac{1}{M^2} \mathbb{E}\left[\left\|\sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right].$$

Несмещенная компрессия: доказательство

- Работаем с $\mathbb{E}[g_k \mid x_k]$:

$$\begin{aligned}\mathbb{E}[g_k \mid x_k] &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[g_{k,m} \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\mathbb{E}[\mathcal{Q}(\nabla f_m(x_k)) \mid \nabla f_m(x_k)] \mid x_k] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}[\nabla f_m(x_k) \mid x_k] = \frac{1}{M} \sum_{m=1}^M \nabla f_m(x_k) = \nabla f(x_k).\end{aligned}$$

- Работаем с $\mathbb{E}[\|g_k\|^2 \mid x_k]$:

$$\mathbb{E}[\|g_k\|^2 \mid x_k] = \mathbb{E}\left[\left\|\frac{1}{M} \sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right] = \frac{1}{M^2} \mathbb{E}\left[\left\|\sum_{m=1}^M g_{k,m}\right\|^2 \mid x_k\right].$$

Несмещенная компрессия: доказательство

- Продолжаем и применяем первое свойство (несмещенность) в определении компрессии:

$$\begin{aligned}\mathbb{E} [\|g_k\|^2 \mid x_k] &= \frac{1}{M^2} \mathbb{E} \left[\left\| \sum_{m=1}^M g_{k,m} \right\|^2 \mid x^k \right] \\&= \frac{1}{M^2} \sum_{m=1}^M \mathbb{E} [\|g_{k,m}\|^2 \mid x^k] \\&\quad + \frac{2}{M^2} \sum_{m \neq l} \mathbb{E} [\langle g_{k,m}, g_{k,l} \rangle \mid x^k] \\&= \frac{1}{M^2} \sum_{m=1}^M \mathbb{E} [\|g_{k,m}\|^2 \mid x^k] \\&\quad + \frac{1}{M^2} \sum_{m \neq l} \mathbb{E} [\langle \mathbb{E} [g_{k,m} \mid \nabla f_m(x_k)], \mathbb{E} [g_{k,l} \mid \nabla f_l(x_k)] \rangle \mid x^k] .\end{aligned}$$

Несмещенная компрессия: доказательство

- Продолжаем и применяем первое свойство (несмещенность) в определении компрессии:

$$\begin{aligned}\mathbb{E} [\|g_k\|^2 \mid x_k] &= \frac{1}{M^2} \mathbb{E} \left[\left\| \sum_{m=1}^M g_{k,m} \right\|^2 \mid x^k \right] \\&= \frac{1}{M^2} \sum_{m=1}^M \mathbb{E} [\|g_{k,m}\|^2 \mid x^k] \\&\quad + \frac{2}{M^2} \sum_{m \neq l} \mathbb{E} [\langle g_{k,m}, g_{k,l} \rangle \mid x^k] \\&= \frac{1}{M^2} \sum_{m=1}^M \mathbb{E} [\|g_{k,m}\|^2 \mid x^k] \\&\quad + \frac{1}{M^2} \sum_{m \neq l} \mathbb{E} [\langle \mathbb{E} [g_{k,m} \mid \nabla f_m(x_k)], \mathbb{E} [g_{k,l} \mid \nabla f_l(x_k)] \rangle \mid x^k].\end{aligned}$$

Несмещенная компрессия: доказательство

- Продолжаем и применяем второе свойство в определении компрессии:

$$\begin{aligned}\mathbb{E} [\|g_k\|^2 \mid x_k] &= \frac{1}{M^2} \sum_{m=1}^M \mathbb{E} [\|Q(\nabla f_m(x_k))\|^2 \mid x^k] \\ &\quad + \frac{1}{M^2} \sum_{m \neq l} \langle \nabla f_m(x_k), \nabla f_l(x_k) \rangle \\ &\leq \frac{\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x_k)\|^2 \\ &\quad + \|\nabla f(x_k)\|^2 \\ &\leq \frac{2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x_k) - \nabla f_m(x^*)\|^2 \\ &\quad + \frac{2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 + \|\nabla f(x_k) - \nabla f(x^*)\|^2.\end{aligned}$$

Несмещенная компрессия: доказательство

- Продолжаем и применяем второе свойство в определении компрессии:

$$\begin{aligned}\mathbb{E} [\|g_k\|^2 \mid x_k] &\leq \frac{4\omega L}{M^2} \sum_{m=1}^M (f_m(x_k) - f_m(x^*) - \langle \nabla f_m(x^*), x_k - x^* \rangle) \\ &\quad + \frac{2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 + 2L(f(x_k) - f(x^*)) \\ &= \frac{4\omega L}{M} (f(x_k) - f(x^*)) \\ &\quad + \frac{2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 + 2L(f(x_k) - f(x^*)).\end{aligned}$$

Несмещенная компрессия: доказательство

- Все, что получили:

$$\mathbb{E} [\|x_{k+1} - x^*\|^2 \mid x_k] = \|x_k - x^*\|^2 - 2\gamma \langle \mathbb{E}[g_k \mid x_k], x_k - x^* \rangle + \gamma^2 \mathbb{E} [\|g_k\|^2 \mid x_k].$$

$$\mathbb{E}[g_k \mid x_k] = \nabla f(x_k).$$

$$\begin{aligned} \mathbb{E} [\|g_k\|^2 \mid x_k] &\leq \frac{4\omega L}{M} (f(x_k) - f(x^*)) \\ &+ \frac{2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 + 2L(f(x_k) - f(x^*)). \end{aligned}$$

Несмещенная компрессия: доказательство

- Объединяем:

$$\begin{aligned}\mathbb{E} [\|x_{k+1} - x^*\|^2 \mid x_k] &\leq \|x_k - x^*\|^2 - 2\gamma \langle \nabla f(x_k), x_k - x^* \rangle \\ &\quad + 2\gamma^2 L \left(\frac{2\omega}{M} + 1 \right) (f(x_k) - f(x^*)) \\ &\quad + \frac{2\gamma^2 \omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2.\end{aligned}$$

- Пользуемся сильной выпуклостью:

$$\begin{aligned}\mathbb{E} [\|x_{k+1} - x^*\|^2 \mid x_k] &\leq \|x_k - x^*\|^2 - 2\gamma \left(\frac{\mu}{2} \|x_k - x^*\|^2 + f(x_k) - f(x^*) \right) \\ &\quad + 2\gamma^2 L \left(\frac{2\omega}{M} + 1 \right) (f(x_k) - f(x^*)) \\ &\quad + \frac{2\gamma^2 \omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2.\end{aligned}$$

Несмещенная компрессия: доказательство

- Если взять полное математическое ожидание

$$\begin{aligned}\mathbb{E} [\|x_{k+1} - x^*\|^2] &\leq (1 - \gamma\mu) \mathbb{E} [\|x_k - x^*\|^2] \\ &\quad - 2\gamma \left[1 - \gamma L \left(\frac{2\omega}{M} + 1 \right) \right] \mathbb{E} [(f(x_k) - f(x^*))] \\ &\quad + \frac{2\gamma^2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2.\end{aligned}$$

- Если $\gamma \leq L^{-1} \left(\frac{2\omega}{M} + 1 \right)^{-1}$, то

$$\begin{aligned}\mathbb{E} [\|x_{k+1} - x^*\|^2] &\leq (1 - \gamma\mu) \mathbb{E} [\|x_k - x^*\|^2] \\ &\quad + \frac{2\gamma^2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2.\end{aligned}$$

QGD: СХОДИМОСТЬ

Теорема (QGD)

Пусть все локальные функции f_m являются μ -сильно выпуклыми и имеют L -Липшицев градиент, тогда если $\eta \leq L^{-1} \left(\frac{2\omega}{M} + 1 \right)^{-1}$, то

$$\mathbb{E} [\|x_K - x^*\|^2] = \mathcal{O} \left((1 - \gamma\mu)^K \|x_0 - x^*\|^2 + \frac{2\omega}{\mu M^2 K} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 \right).$$

При получении данного результата так же использовался подбор γ из работы:



Stich S. Unified Optimal Analysis of the (Stochastic) Gradient Method

- **Вопрос:** какие проблемы есть в этой оценке? Сублинейная сходимость (зависит от гетерогенности данных).

QGD: СХОДИМОСТЬ

Теорема (QGD)

Пусть все локальные функции f_m являются μ -сильно выпуклыми и имеют L -Липшицев градиент, тогда если $\eta \leq L^{-1} \left(\frac{2\omega}{M} + 1 \right)^{-1}$, то

$$\mathbb{E} [\|x_K - x^*\|^2] = \mathcal{O} \left((1 - \gamma\mu)^K \|x_0 - x^*\|^2 + \frac{2\omega}{\mu M^2 K} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 \right).$$

При получении данного результата так же использовался подбор γ из работы:



Stich S. Unified Optimal Analysis of the (Stochastic) Gradient Method

- **Вопрос:** какие проблемы есть в этой оценке? Сублинейная сходимость (зависит от гетерогенности данных).

QGD: СХОДИМОСТЬ

Теорема (QGD)

Пусть все локальные функции f_m являются μ -сильно выпуклыми и имеют L -Липшицев градиент, тогда если $\eta \leq L^{-1} \left(\frac{2\omega}{M} + 1 \right)^{-1}$, то

$$\mathbb{E} [\|x_K - x^*\|^2] = \mathcal{O} \left((1 - \gamma\mu)^K \|x_0 - x^*\|^2 + \frac{2\omega}{\mu M^2 K} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 \right).$$



При получении данного результата так же использовался подбор γ из работы:





Stich S. Unified Optimal Analysis of the (Stochastic) Gradient Method

- **Вопрос:** какие проблемы есть в этой оценке? Сублинейная сходимость (зависит от гетерогенности данных).

Несмещенная компрессия: больше

- Решена проблема с гетерогенностью за счет "памяти":
 Mishchenko K. et al. Distributed Learning with Compressed Gradient Differences
- Ускоренная версия:
 Li Z. et al. Acceleration for compressed gradient descent in distributed and federated optimization

Несмещенная компрессия: больше

- Решена проблема с гетерогенностью за счет "памяти":
 Mishchenko K. et al. Distributed Learning with Compressed Gradient Differences
- Ускоренная версия:
 Li Z. et al. Acceleration for compressed gradient descent in distributed and federated optimization

Смещенная компрессия

Смещённая компрессия

Будем называть (стохастический) оператор (x) оператором смещённой компрессии, если для любого $x \in \mathbb{R}^d$ выполняется:

$$\mathbb{E}[\|C(x) - x\|_2^2] \leq \left(1 - \frac{1}{\delta}\right) \|x\|_2^2,$$

где $\delta \geq 0$.

Смещенная компрессия: примеры

"Жадная" спарсификация (выбор наибольших по модулю компонент)

Рассмотрим стохастический оператор

$$\text{Top}_k(x) = \sum_{i=d-k+1}^d x_{(i)} e_{(i)},$$

где k — некоторое фиксированное число из множества $\{1, \dots, d\}$ (количество компонент вектора x , которые мы передаём; например, можно выбрать $k = 1$), при этом координаты отсортированы по модулю: $|x_{(1)}| \leq |x_{(2)}| \leq \dots \leq |x_{(d)}|$, (e_1, \dots, e_d) — стандартный базис в \mathbb{R}^d . Можно показать, что данный оператор является смещённой компрессией с константой $\delta = \frac{d}{k}$.



Alistarh D. et al. The convergence of sparsified gradient methods

Смещенная компрессия: примеры

- Еще примеры смещенных компрессией:



Beznosikov A. et al. On Biased Compression for Distributed Learning



Vogels T. et al. PowerSGD: Practical Low-Rank Gradient Compression for Distributed Optimization

Смещенная компрессия: идея и доказательство в случае 1 ноды

- Использовать тот же подход, что и в несмещенном случае (QGD).
- Докажем в случае одной ноды:

$$x_{k+1} = x_k - \gamma C(\nabla f(x_k)).$$

Пусть f имеет L -Липшицев градиент и является μ -сильно выпуклой.

- Начнем с того, что воспользуемся Липшицевостью градиента:

$$\begin{aligned} f(x_{k+1}) &= f(x_k - \gamma C(\nabla f(x_k))) \\ &\leq f(x_k) + \langle \nabla f(x_k), -\gamma C(\nabla f(x_k)) \rangle + \frac{L}{2} \| -\gamma C(\nabla f(x_k)) \|^2 \\ &= f(x_k) - \gamma \langle C(\nabla f(x_k)), \nabla f(x_k) \rangle + \frac{\gamma^2 L}{2} \| C(\nabla f(x_k)) \|^2. \end{aligned}$$

Смещенная компрессия: идея и доказательство в случае 1 ноды

- Использовать тот же подход, что и в несмещенном случае (QGD).
- Докажем в случае одной ноды:

$$x_{k+1} = x_k - \gamma C(\nabla f(x_k)).$$

Пусть f имеет L -Липшицев градиент и является μ -сильно выпуклой.

- Начнем с того, что воспользуемся Липшицевостью градиента:

$$\begin{aligned} f(x_{k+1}) &= f(x_k - \gamma C(\nabla f(x_k))) \\ &\leq f(x_k) + \langle \nabla f(x_k), -\gamma C(\nabla f(x_k)) \rangle + \frac{L}{2} \| -\gamma C(\nabla f(x_k)) \|^2 \\ &= f(x_k) - \gamma \langle C(\nabla f(x_k)), \nabla f(x_k) \rangle + \frac{\gamma^2 L}{2} \| C(\nabla f(x_k)) \|^2. \end{aligned}$$

Смещенная компрессия: идея и доказательство в случае 1 ноды

- Использовать тот же подход, что и в несмещенном случае (QGD).
- Докажем в случае одной ноды:

$$x_{k+1} = x_k - \gamma C(\nabla f(x_k)).$$

Пусть f имеет L -Липшицев градиент и является μ -сильно выпуклой.

- Начнем с того, что воспользуемся Липшицевостью градиента:

$$\begin{aligned} f(x_{k+1}) &= f(x_k - \gamma C(\nabla f(x_k))) \\ &\leq f(x_k) + \langle \nabla f(x_k), -\gamma C(\nabla f(x_k)) \rangle + \frac{L}{2} \| -\gamma C(\nabla f(x_k)) \|^2 \\ &= f(x_k) - \gamma \langle C(\nabla f(x_k)), \nabla f(x_k) \rangle + \frac{\gamma^2 L}{2} \| C(\nabla f(x_k)) \|^2. \end{aligned}$$

Смещенная компрессия: доказательство в случае 1 ноды

- Определение компрессора:

$$\begin{aligned} \|\nabla f(x_k)\|^2 - 2\mathbb{E}_C [\langle C(\nabla f(x_k)), \nabla f(x_k) \rangle] + \mathbb{E}_C [\|C(\nabla f(x_k))\|^2] \\ = \mathbb{E}_C [\|C(\nabla f(x_k)) - \nabla f(x_k)\|^2] \leq \left(1 - \frac{1}{\delta}\right) \|\nabla f(x_k)\|^2. \end{aligned}$$

- Откуда:

$$-\gamma \mathbb{E}_C [\langle C(\nabla f(x_k)), \nabla f(x_k) \rangle] + \frac{\gamma}{2} \mathbb{E}_C [\|C(\nabla f(x_k))\|^2] \leq -\frac{\gamma}{2\delta} \|\nabla f(x_k)\|^2.$$

Смещенная компрессия: доказательство в случае 1 ноды

- Определение компрессора:

$$\begin{aligned} \|\nabla f(x_k)\|^2 - 2\mathbb{E}_C [\langle C(\nabla f(x_k)), \nabla f(x_k) \rangle] + \mathbb{E}_C [\|C(\nabla f(x_k))\|^2] \\ = \mathbb{E}_C [\|C(\nabla f(x_k)) - \nabla f(x_k)\|^2] \leq \left(1 - \frac{1}{\delta}\right) \|\nabla f(x_k)\|^2. \end{aligned}$$

- Откуда:

$$-\gamma \mathbb{E}_C [\langle C(\nabla f(x_k)), \nabla f(x_k) \rangle] + \frac{\gamma}{2} \mathbb{E}_C [\|C(\nabla f(x_k))\|^2] \leq -\frac{\gamma}{2\delta} \|\nabla f(x_k)\|^2.$$

Смещенная компрессия: доказательство в случае 1 ноды

- С двух предыдущих слайдов:

$$f(x_{k+1}) - \leq f(x_k) - \gamma \langle C(\nabla f(x_k)), \nabla f(x_k) \rangle + \frac{\gamma^2 L}{2} \|C(\nabla f(x_k))\|^2.$$

$$-\gamma \mathbb{E}_C [\langle C(\nabla f(x_k)), \nabla f(x_k) \rangle] + \frac{\gamma}{2} \mathbb{E}_C [\|C(\nabla f(x_k))\|^2] \leq -\frac{\gamma}{2\delta} \|\nabla f(x_k)\|^2.$$

- Сложим, вычтем из обеих частей $f(x^*)$ и возьмем полное мат. ожидание:

$$\begin{aligned} \mathbb{E}[f(x_{k+1}) - f(x^*)] &\leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2} (1 - \gamma L) \mathbb{E}[\|C(\nabla f(x_k))\|^2] \\ &\quad - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2]. \end{aligned}$$

- Возьмем $\gamma \leq \frac{1}{L}$:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2].$$

Смещенная компрессия: доказательство в случае 1 ноды

- С двух предыдущих слайдов:

$$f(x_{k+1}) - \leq f(x_k) - \gamma \langle C(\nabla f(x_k)), \nabla f(x_k) \rangle + \frac{\gamma^2 L}{2} \|C(\nabla f(x_k))\|^2.$$

$$-\gamma \mathbb{E}_C [\langle C(\nabla f(x_k)), \nabla f(x_k) \rangle] + \frac{\gamma}{2} \mathbb{E}_C [\|C(\nabla f(x_k))\|^2] \leq -\frac{\gamma}{2\delta} \|\nabla f(x_k)\|^2.$$

- Сложим, вычтем из обеих частей $f(x^*)$ и возьмем полное мат. ожидание:

$$\begin{aligned} \mathbb{E}[f(x_{k+1}) - f(x^*)] &\leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2} (1 - \gamma L) \mathbb{E}[\|C(\nabla f(x_k))\|^2] \\ &\quad - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2]. \end{aligned}$$

- Возьмем $\gamma \leq \frac{1}{L}$:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2].$$

Смещенная компрессия: доказательство в случае 1 ноды

- С двух предыдущих слайдов:

$$f(x_{k+1}) - \leq f(x_k) - \gamma \langle C(\nabla f(x_k)), \nabla f(x_k) \rangle + \frac{\gamma^2 L}{2} \|C(\nabla f(x_k))\|^2.$$

$$-\gamma \mathbb{E}_C [\langle C(\nabla f(x_k)), \nabla f(x_k) \rangle] + \frac{\gamma}{2} \mathbb{E}_C [\|C(\nabla f(x_k))\|^2] \leq -\frac{\gamma}{2\delta} \|\nabla f(x_k)\|^2.$$

- Сложим, вычтем из обеих частей $f(x^*)$ и возьмем полное мат. ожидание:

$$\begin{aligned} \mathbb{E}[f(x_{k+1}) - f(x^*)] &\leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2} (1 - \gamma L) \mathbb{E}[\|C(\nabla f(x_k))\|^2] \\ &\quad - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2]. \end{aligned}$$

- Возьмем $\gamma \leq \frac{1}{L}$:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2].$$

Смещенная компрессия: доказательство в случае 1 ноды

- С предыдущего слайда:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2].$$

- Сильная выпуклость (или даже более слабое условие PL):

$$2\mu(f(x_k) - f(x^*)) \leq \|\nabla f(x_k)\|^2.$$

- Соединим два предыдущих:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \left(1 - \frac{\gamma\mu}{\delta}\right) \mathbb{E}[f(x_k) - f(x^*)].$$

Смещенная компрессия: доказательство в случае 1 ноды

- С предыдущего слайда:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \mathbb{E}[f(x_k) - f(x^*)] - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(x_k)\|^2].$$

- Сильная выпуклость (или даже более слабое условие PL):

$$2\mu(f(x_k) - f(x^*)) \leq \|\nabla f(x_k)\|^2.$$

- Соединим два предыдущих:

$$\mathbb{E}[f(x_{k+1}) - f(x^*)] \leq \left(1 - \frac{\gamma\mu}{\delta}\right) \mathbb{E}[f(x_k) - f(x^*)].$$

Смещенная компрессия: теорема в случае 1 ноды

Теорема (сходимость QGD со смещенной компрессией в случае 1 ноды)

Пусть f μ -сильно выпуклая (или PL) и имеет L -Липшицев градиент, тогда QGD для одной ноды с шагом $\gamma \leq 1/L$ и со смещенным компрессором с параметром δ сходится и выполнено:

$$f(x_K) - f(x^*) \leq \left(1 - \frac{\gamma\mu}{\delta}\right)^K (f(x_0) - f(x^*)).$$

Смещенная компрессия: не так все просто

- Рассмотрим следующую распределенную задачу с $M = 3$, $d = 3$ и локальными функциями:

$$f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_2(x) = \langle b, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_3(x) = \langle c, x \rangle^2 + \frac{1}{4} \|x\|^2,$$

где $a = (-3, 2, 2)$, $b = (2, -3, 2)$ и $c = (2, 2, -3)$.

- Вопрос:** где у нее оптимум? $(0, 0, 0)$.
- Пусть стартовая точка $x_0 = (t, t, t)$ для какого-то $t > 0$. Тогда локальные градиенты:

$$\nabla f_1(x_0) = \frac{t}{2}(-11, 9, 9), \quad \nabla f_2(x_0) = \frac{t}{2}(9, -11, 9), \quad \nabla f_3(x_0) = \frac{t}{2}(9, 9, -11).$$

- Вопрос:** как будет выглядеть шаг QGD (градиентного спуска с сжатиями), если мы будем использовать Top_1 компрессию?

$$x_1 = (t, t, t) + \eta \cdot \frac{11}{6} (t, t, t) = \left(1 + \frac{11\eta}{6}\right) x_0.$$

- Мы удаляемся от решения геометрически для любого $\eta > 0$.

Смещенная компрессия: не так все просто

- Рассмотрим следующую распределенную задачу с $M = 3$, $d = 3$ и локальными функциями:

$$f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_2(x) = \langle b, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_3(x) = \langle c, x \rangle^2 + \frac{1}{4} \|x\|^2,$$

где $a = (-3, 2, 2)$, $b = (2, -3, 2)$ и $c = (2, 2, -3)$.

- Вопрос:** где у нее оптимум? $(0, 0, 0)$.
- Пусть стартовая точка $x_0 = (t, t, t)$ для какого-то $t > 0$. Тогда локальные градиенты:

$$\nabla f_1(x_0) = \frac{t}{2}(-11, 9, 9), \quad \nabla f_2(x_0) = \frac{t}{2}(9, -11, 9), \quad \nabla f_3(x_0) = \frac{t}{2}(9, 9, -11).$$

- Вопрос:** как будет выглядеть шаг QGD (градиентного спуска с сжатиями), если мы будем использовать Top_1 компрессию?

$$x_1 = (t, t, t) + \eta \cdot \frac{11}{6} (t, t, t) = \left(1 + \frac{11\eta}{6}\right) x_0.$$

- Мы удаляемся от решения геометрически для любого $\eta > 0$.

Смещенная компрессия: не так все просто

- Рассмотрим следующую распределенную задачу с $M = 3$, $d = 3$ и локальными функциями:

$$f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_2(x) = \langle b, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_3(x) = \langle c, x \rangle^2 + \frac{1}{4} \|x\|^2,$$

где $a = (-3, 2, 2)$, $b = (2, -3, 2)$ и $c = (2, 2, -3)$.

- Вопрос:** где у нее оптимум? $(0, 0, 0)$.
- Пусть стартовая точка $x_0 = (t, t, t)$ для какого-то $t > 0$. Тогда локальные градиенты:

$$\nabla f_1(x_0) = \frac{t}{2}(-11, 9, 9), \quad \nabla f_2(x_0) = \frac{t}{2}(9, -11, 9), \quad \nabla f_3(x_0) = \frac{t}{2}(9, 9, -11).$$

- Вопрос:** как будет выглядеть шаг QGD (градиентного спуска с сжатиями), если мы будем использовать Top_1 компрессию?

$$x_1 = (t, t, t) + \eta \cdot \frac{11}{6} (t, t, t) = \left(1 + \frac{11\eta}{6}\right) x_0.$$

- Мы удаляемся от решения геометрически для любого $\eta > 0$.

Смещенная компрессия: не так все просто

- Рассмотрим следующую распределенную задачу с $M = 3$, $d = 3$ и локальными функциями:

$$f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_2(x) = \langle b, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_3(x) = \langle c, x \rangle^2 + \frac{1}{4} \|x\|^2,$$

где $a = (-3, 2, 2)$, $b = (2, -3, 2)$ и $c = (2, 2, -3)$.

- Вопрос:** где у нее оптимум? $(0, 0, 0)$.
- Пусть стартовая точка $x_0 = (t, t, t)$ для какого-то $t > 0$. Тогда локальные градиенты:

$$\nabla f_1(x_0) = \frac{t}{2}(-11, 9, 9), \quad \nabla f_2(x_0) = \frac{t}{2}(9, -11, 9), \quad \nabla f_3(x_0) = \frac{t}{2}(9, 9, -11).$$

- Вопрос:** как будет выглядеть шаг QGD (градиентного спуска с сжатиями), если мы будем использовать Top_1 компрессию?

$$x_1 = (t, t, t) + \eta \cdot \frac{11}{6} (t, t, t) = \left(1 + \frac{11\eta}{6}\right) x_0.$$

- Мы удаляемся от решения геометрически для любого $\eta > 0$.

Смещенная компрессия: не так все просто

- Рассмотрим следующую распределенную задачу с $M = 3$, $d = 3$ и локальными функциями:

$$f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_2(x) = \langle b, x \rangle^2 + \frac{1}{4} \|x\|^2, \quad f_3(x) = \langle c, x \rangle^2 + \frac{1}{4} \|x\|^2,$$

где $a = (-3, 2, 2)$, $b = (2, -3, 2)$ и $c = (2, 2, -3)$.

- Вопрос:** где у нее оптимум? $(0, 0, 0)$.
- Пусть стартовая точка $x_0 = (t, t, t)$ для какого-то $t > 0$. Тогда локальные градиенты:

$$\nabla f_1(x_0) = \frac{t}{2}(-11, 9, 9), \quad \nabla f_2(x_0) = \frac{t}{2}(9, -11, 9), \quad \nabla f_3(x_0) = \frac{t}{2}(9, 9, -11).$$

- Вопрос:** как будет выглядеть шаг QGD (градиентного спуска с сжатиями), если мы будем использовать Top_1 компрессию?

$$x_1 = (t, t, t) + \eta \cdot \frac{11}{6}(t, t, t) = \left(1 + \frac{11\eta}{6}\right) x_0.$$

- Мы удаляемся от решения геометрически для любого $\eta > 0$.

Смещенная компрессия: компенсация ошибки

- Попробуем запоминать то, что не передали в процессе общения:

$$e_{1,m} = 0 + \gamma \nabla f_m(x_0) - C(0 + \gamma \nabla f_m(x_0)).$$

- И добавлять это в будущие посылки:

$$C(e_{1,m} + \gamma \nabla f_m(x_1))$$

- На произвольной итерации это записывается так:

$$\text{Посылка: } C(e_{k,m} + \gamma \nabla f_m(x_k)),$$

$$e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(x_k) - C(e_{k,m} + \gamma \nabla f_m(x_k))$$

- Это техника называется компенсация ошибка (error feedback).



Stich S. et al. Sparsified SGD with memory

Смещенная компрессия: компенсация ошибки

- Попробуем запоминать то, что не передали в процессе общения:

$$e_{1,m} = 0 + \gamma \nabla f_m(x_0) - C(0 + \gamma \nabla f_m(x_0)).$$

- И добавлять это в будущие посылки:

$$C(e_{1,m} + \gamma \nabla f_m(x_1))$$

- На произвольной итерации это записывается так:

$$\text{Посылка: } C(e_{k,m} + \gamma \nabla f_m(x_k)),$$

$$e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(x_k) - C(e_{k,m} + \gamma \nabla f_m(x_k))$$

- Это техника называется компенсация ошибка (error feedback).



Stich S. et al. Sparsified SGD with memory

Смещенная компрессия: компенсация ошибки

- Попробуем запоминать то, что не передали в процессе общения:

$$e_{1,m} = 0 + \gamma \nabla f_m(x_0) - C(0 + \gamma \nabla f_m(x_0)).$$

- И добавлять это в будущие посылки:

$$C(e_{1,m} + \gamma \nabla f_m(x_1))$$

- На произвольной итерации это записывается так:

$$\text{Посылка: } C(e_{k,m} + \gamma \nabla f_m(x_k)),$$

$$e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(x_k) - C(e_{k,m} + \gamma \nabla f_m(x_k))$$

- Это техника называется компенсация ошибка (error feedback).



Stich S. et al. Sparsified SGD with memory

Смещенная компрессия: компенсация ошибки

- Попробуем запоминать то, что не передали в процессе общения:

$$e_{1,m} = 0 + \gamma \nabla f_m(x_0) - C(0 + \gamma \nabla f_m(x_0)).$$

- И добавлять это в будущие посылки:

$$C(e_{1,m} + \gamma \nabla f_m(x_1))$$

- На произвольной итерации это записывается так:

$$\text{Посылка: } C(e_{k,m} + \gamma \nabla f_m(x_k)),$$

$$e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(x_k) - C(e_{k,m} + \gamma \nabla f_m(x_k))$$

- Это техника называется компенсация ошибка (error feedback).



Stich S. et al. Sparsified SGD with memory

GD с error feedback

Алгоритм 1 GD с error feedback

Вход: Размер шага $\gamma > 0$, стартовая точка $x_0 \in \mathbb{R}^d$, стартовые ошибки $e_{0,m} = 0$ для всех m от 1 до M , количество итераций K

- 1: **for** $k = 0, 1, \dots, K - 1$ **do**
- 2: Отправить x_k всем рабочим ▷ выполняется сервером
- 3: **for** $m = 1, \dots, M$ параллельно **do**
- 4: Принять x_k от мастера ▷ выполняется рабочими
- 5: Вычислить градиент $\nabla f(x_k)$ в точке x_k ▷ выполняется рабочими
- 6: Сгенерировать $g_{k,m} = C(e_{k,m} + \gamma \nabla f(x_k))$ ▷ выполняется рабочими
- 7: Вычислить $e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(x_k) - g_{k,m}$ ▷ выполняется рабочими
- 8: Отправить $g_{k,m}$ мастеру ▷ выполняется рабочими
- 9: **end for**
- 10: Принять g_k от всех рабочих ▷ выполняется сервером
- 11: Вычислить $g_k = \frac{1}{M} \sum_{m=1}^M g_{k,m}$ ▷ выполняется сервером
- 12: $x_{k+1} = x_k - g_k$ ▷ выполняется сервером
- 13: **end for**

Выход: x_K

GD с error feedback: сходимость

Teorema GD c error feedback

Пусть все локальные функции f_m являются μ -сильно выпуклыми и имеют L -Липшицев градиент, тогда если $\eta \leq \frac{1}{28\delta L}$, то

$$\mathbb{E} [f(\tilde{x}_K) - f(x^*)] \leq \mathcal{O} \left(\delta L \|x_0 - x^*\|^2 \exp \left(-\frac{\gamma \mu K}{2} \right) + \frac{\delta}{\mu K} \cdot \frac{1}{M} \sum_{m=1}^M \|\nabla f_m(x^*)\|^2 \right).$$





Stich S. and Karimireddy S. The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication





Beznosikov A. et al. On Biased Compression for Distributed Learning

Смещенная компрессия: больше

- Решена проблема с гетерогенностью за счет "памяти":
 Richtarik P. et al. EF21: A New, Simpler, Theoretically Better, and Practically Faster Error Feedback
- Ускоренная версия:
 Qian X. Error Compensated Distributed SGD Can Be Accelerated

Смещенная компрессия: больше

- Решена проблема с гетерогенностью за счет "памяти":
 Richtarik P. et al. EF21: A New, Simpler, Theoretically Better, and Practically Faster Error Feedback
- Ускоренная версия:
 Qian X. Error Compensated Distributed SGD Can Be Accelerated

Несмещенная против смещенной

- Лучшая оценка на число коммуникаций для неускоренного метода с несмещенной компрессией (DIANA):

$$\mathcal{O}\left(\left[1 + \frac{\omega}{M}\right] \frac{L}{\mu} \log \frac{1}{\varepsilon}\right).$$

- Лучшая оценка на число коммуникаций для неускоренного метода со смещенной компрессией (EF-21):

$$\mathcal{O}\left([1+\delta]\frac{L}{\mu}\log\frac{1}{\varepsilon}\right).$$

Несмещенная против смещенной

- Лучшая оценка на число коммуникаций для неускоренного метода с несмещенной компрессией (DIANA):

$$\mathcal{O}\left(\left[1 + \frac{\omega}{M}\right] \frac{L}{\mu} \log \frac{1}{\varepsilon}\right).$$

- Лучшая оценка на число коммуникаций для неускоренного метода со смещенной компрессией (EF-21):

$$\mathcal{O}\left([1+\delta]\frac{L}{\mu}\log\frac{1}{\varepsilon}\right).$$

- **Вопрос:** что можно о них сказать? Как они соотносятся с несжатými методами? Они хуже. Но важно не число коммуникаций, а количество передаваемой информации.

Несмещенная против смещенной

- Компрессоры сжимают информацию в β раз и типично, что $\beta \geq \omega$ и $\beta \geq \delta$.
- Лучшая оценка на число информации для неускоренного метода с несмещенной компрессией (DIANA):

$$\mathcal{O}\left(\left[\frac{1}{\beta} + \frac{1}{M}\right] \frac{L}{\mu} \log \frac{1}{\varepsilon}\right).$$

Несмещенная против смещенной

- Компрессоры сжимают информацию в β раз и типично, что $\beta \geq \omega$ и $\beta \geq \delta$.
- Лучшая оценка на число информации для неускоренного метода с несмещенной компрессией (DIANA):

$$\mathcal{O}\left(\left[\frac{1}{\beta} + \frac{1}{M}\right] \frac{L}{\mu} \log \frac{1}{\varepsilon}\right).$$

Несмещенная против смещенной

- Компрессоры сжимают информацию в β раз и типично, что $\beta \geq \omega$ и $\beta \geq \delta$.
- Лучшая оценка на число информации для неускоренного метода с несмещенной компрессией (DIANA):

$$\mathcal{O}\left(\left[\frac{1}{\beta} + \frac{1}{M}\right] \frac{L}{\mu} \log \frac{1}{\varepsilon}\right).$$

- Несмещенный компрессор доказуемо улучшает число передаваемой информации, фактор улучшения: $\left[\frac{1}{\beta} + \frac{1}{M} \right]$.
- Смещенный компрессор не улучшает число передаваемой информации в общем случае.

Несмещенная против смещенной

- Компрессоры сжимают информацию в β раз и типично, что $\beta \geq \omega$ и $\beta \geq \delta$.
- Лучшая оценка на число информации для неускоренного метода с несмещенной компрессией (DIANA):

$$\mathcal{O}\left(\left[\frac{1}{\beta} + \frac{1}{M}\right] \frac{L}{\mu} \log \frac{1}{\varepsilon}\right).$$

- Несмещенный компрессор доказуемо улучшает число передаваемой информации, фактор улучшения: $\left[\frac{1}{\beta} + \frac{1}{M} \right]$.
- Смещенный компрессор не улучшает число передаваемой информации в общем случае.

Data similarity

Data similarity

- И снова распределенная задача обучения:

$$f(x) = \frac{1}{M} \sum_{m=1}^M f_m(x) = \frac{1}{M} \sum_{m=1}^M \left[\frac{1}{N} \sum_{i=1}^N \ell(x, z_i) \right],$$

где z_i – элемент выборки, ℓ – функция потерь.

- Предположим, что мы можем разбить обучающую выборку равномерно между устройствами (например, если используются кластерные или коллаборативные вычисления на открытых данных).
- Что это может дать? Похожесть локальных функций потерь.
- Утверждается, что для любого x

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\| \leq \delta.$$

- Вопрос:** если ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция, то в общем случае (если не предполагать равномерность распределения данных), что можно сказать о δ ? $\delta \sim L$.

Data similarity

- И снова распределенная задача обучения:

$$f(x) = \frac{1}{M} \sum_{m=1}^M f_m(x) = \frac{1}{M} \sum_{m=1}^M \left[\frac{1}{N} \sum_{i=1}^N \ell(x, z_i) \right],$$

где z_i – элемент выборки, ℓ – функция потерь.

- Предположим, что мы можем разбить обучающую выборку равномерно между устройствами (например, если используются кластерные или коллаборативные вычисления на открытых данных).
- Что это может дать? Похожесть локальных функций потерь.
- Утверждается, что для любого x

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\| \leq \delta.$$

- Вопрос:** если ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция, то в общем случае (если не предполагать равномерность распределения данных), что можно сказать о δ ? $\delta \sim L$.

Data similarity

- И снова распределенная задача обучения:

$$f(x) = \frac{1}{M} \sum_{m=1}^M f_m(x) = \frac{1}{M} \sum_{m=1}^M \left[\frac{1}{N} \sum_{i=1}^N \ell(x, z_i) \right],$$

где z_i – элемент выборки, ℓ – функция потерь.

- Предположим, что мы можем разбить обучающую выборку равномерно между устройствами (например, если используются кластерные или коллаборативные вычисления на открытых данных).
- Что это может дать? Похожесть локальных функций потерь.
- Утверждается, что для любого x

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\| \leq \delta.$$

- Вопрос:** если ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция, то в общем случае (если не предполагать равномерность распределения данных), что можно сказать о δ ? $\delta \sim L$.

Data similarity

- И снова распределенная задача обучения:

$$f(x) = \frac{1}{M} \sum_{m=1}^M f_m(x) = \frac{1}{M} \sum_{m=1}^M \left[\frac{1}{N} \sum_{i=1}^N \ell(x, z_i) \right],$$

где z_i – элемент выборки, ℓ – функция потерь.

- Предположим, что мы можем разбить обучающую выборку равномерно между устройствами (например, если используются кластерные или коллаборативные вычисления на открытых данных).
- Что это может дать? Похожесть локальных функций потерь.
- Утверждается, что для любого x

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\| \leq \delta.$$

- Вопрос:** если ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция, то в общем случае (если не предполагать равномерность распределения данных), что можно сказать о δ ? $\delta \sim L$

Data similarity

- И снова распределенная задача обучения:

$$f(x) = \frac{1}{M} \sum_{m=1}^M f_m(x) = \frac{1}{M} \sum_{m=1}^M \left[\frac{1}{N} \sum_{i=1}^N \ell(x, z_i) \right],$$

где z_i – элемент выборки, ℓ – функция потерь.

- Предположим, что мы можем разбить обучающую выборку равномерно между устройствами (например, если используются кластерные или коллаборативные вычисления на открытых данных).
- Что это может дать? Похожесть локальных функций потерь.
- Утверждается, что для любого x

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\| \leq \delta.$$

- Вопрос:** если ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция, то в общем случае (если не предполагать равномерность распределения данных), что можно сказать о δ ? $\delta \sim L$.

Матричное неравенство Хёфдинга

Теорема (Матричное неравенство Хёфдинга)

Рассмотрим конечную последовательность случайных квадратных матриц $\{X_i\}_{i=1}^N$. Пусть в этой последовательности матрицы независимы, эрмитовы и имеют размерность d . Предположим так же, что $\mathbb{E}[X_i] = 0$, и $X_i^2 \preceq A^2$ почти наверное, где A – неслучайная эрмитова матрица. Тогда с вероятностью $1 - p$ выполнено, что

$$\left\| \sum_{i=1}^N X_i \right\| \leq \sqrt{8N \|A^2\| \cdot \ln(d/p)}.$$



Tropp J. An introduction to matrix concentration inequalities



Tropp J. User-friendly tail bounds for sums of random matrices

Параметр схожести

- Локальная функция потерь:

$$f_m(x) = \frac{1}{N} \sum_{i=1}^N l(x, z_i).$$

- ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция (например, квадратичная или логрегрессия). Тогда имеем $\nabla^2 \ell(x, z_i) \preceq L I$ для любого x и z_i (здесь I – единичная матрица.).
- Распределим все данные равномерно по всем нодам.
- **Вопрос:** что нужно взять в качестве X_i в неравенстве Хёфдинга?
 $X_i = \frac{1}{N} [\nabla \ell(x, z_i) - \nabla f(x)]$. Легко проверить, что все условия матричного неравенства Хёфдинга для нее выполнены, в частности, $A^2 = \frac{4L^2}{N^2} I$.

Параметр схожести

- Локальная функция потерь:

$$f_m(x) = \frac{1}{N} \sum_{i=1}^N l(x, z_i).$$

- ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция (например, квадратичная или логрегрессия). Тогда имеем $\nabla^2 l(x, z_i) \preceq LI$ для любого x и z_i (здесь I – единичная матрица.).
- Распределим все данные равномерно по всем нодам.
- Вопрос:** что нужно взять в качестве X_i в неравенстве Хёфдинга?
 $X_i = \frac{1}{N} [\nabla l(x, z_i) - \nabla f(x)]$. Легко проверить, что все условия матричного неравенства Хёфдинга для нее выполнены, в частности, $A^2 = \frac{4L^2}{N^2} I$.

Параметр схожести

- Локальная функция потерь:

$$f_m(x) = \frac{1}{N} \sum_{i=1}^N l(x, z_i).$$

- ℓ – L -гладкая (L -Липшицев градиент), выпуклая, дважды дифференцируемая функция (например, квадратичная или логрегрессия). Тогда имеем $\nabla^2 l(x, z_i) \preceq LI$ для любого x и z_i (здесь I – единичная матрица.).
- Распределим все данные равномерно по всем нодам.
- Вопрос:** что нужно взять в качестве X_i в неравенстве Хёфдинга?
 $X_i = \frac{1}{N} [\nabla l(x, z_i) - \nabla f(x)]$. Легко проверить, что все условия матричного неравенства Хёфдинга для нее выполнены, в частности, $A^2 = \frac{4L^2}{N^2} I$.

Параметр схожести: итог

- В итоге имеем

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\| \leq \delta \sim \frac{L}{\sqrt{N}}.$$

- Для квадратичных задач можно получить оценку вида:

$$\|\nabla^2 f_m(x) - \nabla^2 f(x)\| \leq \delta \sim \frac{L}{N}.$$



Hendrikx H. et al. Statistically Preconditioned Accelerated Gradient Method for Distributed Optimization

- В любом случае следует вывод: чем больше размер локальной выборки, тем меньше параметр схожести (похожи между собой гессианы).

Метод в общем виде

- Рассмотрим зеркальный спуск:

$$x_{k+1} = \arg \min_{x \in \mathbb{R}^d} (\gamma \langle \nabla f(x_k), x \rangle + V(x, x_k)),$$

где $V(x, y)$ – дивергенция Брегмана, порожденная функцией строго-выпуклой функцией $\varphi(x)$:

$$V(x, y) = \varphi(x) - \varphi(y) - \langle \nabla \varphi(y); x - y \rangle.$$

Сходимость в общем виде: доказательство

- С предыдущего слайда:

$$f(x_{k+1}) - f(x_k) \leq L_\varphi V(x^*, x_k) - L_\varphi V(x^*, x_{k+1}) - \langle \nabla f(x_k), x_k - x^* \rangle.$$

- Относительная сильная выпуклость:

$$\mu_\varphi V(x^*, x_k) \leq f(x^*) - f(x_k) - \langle \nabla f(x_k); x^* - x_k \rangle$$

- Сложим два предыдущих и немного поперемещаем:

$$f(x_{k+1}) - f(x^*) \leq (L_\varphi - \mu_\varphi)V(x^*, x_k) - L_\varphi V(x^*, x_{k+1}).$$

- В силу того, что x^* – оптимум:

$$V(x^*, x_{k+1}) \leq \left(1 - \frac{\mu_\varphi}{L_\varphi}\right) V(x^*, x_k).$$

Сходимость для задачи data similarity: доказательство

- Найдем μ_φ . Из сильно выпуклости функции f :

$$\mu l \preceq \nabla^2 f(x) \Rightarrow \delta l \preceq \frac{2\delta}{\mu} \nabla^2 f(x) - \delta l.$$

- Из $\|\nabla^2 f_1(x) - \nabla^2 f(x)\| \leq \delta$ имеем:

$$\nabla^2 f_1(x) - \nabla^2 f(x) \preceq \delta l.$$

- Объединяем два предыдущих пункта:

$$\nabla^2 f_1(x) - \nabla^2 f(x) \preceq \frac{2\delta}{\mu} \nabla^2 f(x) - \delta l.$$

- Откуда:

$$\nabla^2 f_1(x) + \delta l \preceq \frac{2\delta + \mu}{\mu} \nabla^2 f(x) \Rightarrow \mu_\varphi = \frac{\mu}{2\delta + \mu}.$$

Сходимость для задачи data similarity: теорема

Теорема (сходимость для задачи data similarity)

Пусть f сильно выпуклая, f_i выпуклые, а ℓ - гладкие, а $\varphi(x) = f_1(x) + \delta\|x\|^2$, тогда зеркальный спуск с шагом $\gamma = 1$ сходится и выполнено:

$$V(x^*, x_K) \leq \left(1 - \frac{\mu}{\mu + 2\delta}\right)^K V(x^*, x_0).$$

Лучше?

- Оценка на число коммуникаций в условиях data similarity:

$$K = \mathcal{O} \left(\left[1 + \frac{\delta}{\mu} \right] \log \frac{1}{\varepsilon} \right).$$

- Оценка на число коммуникаций для обычного распределенного градиентного спуска:

$$K = \mathcal{O} \left(\frac{L}{\mu} \log \frac{1}{\varepsilon} \right).$$

- Напомним, что $\delta \sim \frac{L}{\sqrt{N}}$, т.е. может быть значительное улучшение.

Лучше?

- Но есть ведь и ускоренный градиентный метод, который дает оценки:

$$K = \mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon}\right).$$

Более того, эти оценки улучшаемые в общем случае.



Scaman K. et al. Optimal Convergence Rates for Convex Distributed Optimization in Networks

$$K = \Omega \left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\varepsilon} \right),$$



Лучше?

- Но есть ведь и ускоренный градиентный метод, который дает оценки:

$$K = \mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon}\right).$$

Более того, эти оценки улучшаемые в общем случае.



Scaman K. et al. Optimal Convergence Rates for Convex Distributed Optimization in Networks

Вопрос: как в общем случае доказывается какого-то алгоритма (необязательно оптимизационного) для класса задач?

- Для задачи data similarity так же имеются нижние оценки (в 2015 году):

$$K = \Omega \left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\varepsilon} \right),$$

т.е. предполагается возможно ускорение.



Arjevani Y. and Shamir O. Communication complexity of distributed convex learning and optimization

Оптимальный алгоритм

- У данной проблемы довольно большая история:

Reference	Communication complexity	Local gradient complexity	Order	Limitations
DANE [42]	$\mathcal{O}\left(\frac{\delta^2}{\mu^2} \log \frac{1}{\epsilon}\right)$	— ⁽²⁾	1st	quadratic
DiSCD [51]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \left(\log \frac{1}{\epsilon} + C^2 \Delta F_0\right) \log \frac{L}{\mu}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \left(\log \frac{1}{\epsilon} + C^2 \Delta F_0\right) \log \frac{L}{\mu}\right)$	2nd	C - self-concordant ⁽¹⁾
AIDE [40]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon} \log \frac{L}{\delta}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon} \log \frac{L}{\delta}\right)$ ⁽⁴⁾	1st	quadratic
DANE-LS [50]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \frac{\delta^{3/2}}{\mu^{3/2}} \log \frac{1}{\epsilon}\right)$ ⁽⁵⁾	1st/2nd	quadratic ⁽⁶⁾
DANE-HB [50]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \frac{\delta}{\mu} \log \frac{1}{\epsilon}\right)$ ⁽⁵⁾	1st/2nd	quadratic ⁽⁶⁾
SONATA [45]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon}\right)$	— ⁽²⁾	1st	decentralized
SPAG [21]	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}\right)$ ⁽¹⁾	— ⁽²⁾	1st	M - Lipschitz hessian
DiRegINA [12]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon} + \sqrt{\frac{M^2 R \mu}{\delta}}\right)$	— ⁽²⁾	2nd	M - Lipschitz hessian
ACN [1]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon} + \sqrt{\frac{M^2 R \mu}{\delta}}\right)$	— ⁽²⁾	2nd	M - Lipschitz hessian
Acc SONATA [46]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon} \log \frac{L}{\mu}\right)$	— ⁽²⁾	1st	decentralized
This paper	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}\right)$	1st	

Figure: Таблица из статъи Kovalev D. et al. Optimal Gradient Sliding and its Application to Distributed Optimization Under Similarity

В частности, подход зеркального спуска с необычной дивергенцией называется DANE.



Все! Спасибо за внимание!