

Применения операторов сжатия в задачах распределенной оптимизации

Касюк В.А.

Московский физико-технический институт

Москва,
2024 г.

Целью исследования является создание новых, более совершенных операторов сжатия для задач распределенной оптимизации, а также анализ их особенностей.

Проблематика данной задачи заключается в том, что предыдущие исследования показывают эмпирически лучшую сходимость смещенных операторов сжатия, однако теория мало развита и построено мало методов.

О смещенном сжатии для распределенного обучения

Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, Mher Safaryan On Biased Compression for Distributed Learning

Унифицированная теория методов SGD

Eduard Gorbunov, Filip Hanzely, Peter Richtárik A Unified Theory of SGD: Variance Reduction, Sampling, Quantization and Coordinate Descent

Дополнительные ссылки

Е.А. Воронцова, Р. Хильдебранд, А.В. Гасников, Ф.С. Стонякин Выпуклая оптимизация

Решается задача оптимизации :

$$\min_{x \in R^d} f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x), \quad (0.1)$$

где $x \in R^d$ - пространство признаков размерности d , n - количество устройств/узлов, $f_i(x) : R^d \rightarrow R$ - потери, которые получает модель на устройстве под номером i . Часто функция потерь имеем вид :

$$f_i(x) = E_{\xi \sim P_i}[f_{\xi}(x)],$$

где P_i - это распределение тренировочных данных на устройстве i . То, как данные распределяются между работниками оказывает большое влияние на процесс обучения.

Предыдущие результаты: методы Top-k, Rand-K

Базовым решением задачи (0.1) является обычный градиентный спуск (GD), который имеет вид :

$$x^{k+1} = x^k - \frac{\eta^k}{n} \sum_{i=1}^n \nabla f_i(x^k) \quad (0.2)$$

Теперь поставим задачу оптимизации, которую будем решать с помощью градиентного спуска с сжатием (CGD - compressed gradient descent). Пусть решается задача:

$$\min_{x \in R^d} f(x) \quad (0.3)$$

, где $f : R^d \rightarrow R$, L - гладкая и μ -сильно-выпуклая, методом :

$$x^{k+1} = x^k - \eta C(\nabla f(x)) \quad (0.4)$$

, где $C : R^d \rightarrow R^d$ - оператор сжатия и $\eta > 0$ - темп обучения.

Под операторами сжатия мы подразумеваем (возможно рандомизированное) отображение $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ с некоторыми требованиями. Как правило, в литературе рассматриваются *несмещенные* операторы сжатия \mathcal{C} с ограниченным вторым моментом, т.е. :

Definition

Пусть $\zeta \geq 1$. Тогда $\mathcal{C} \in \mathbb{U}(\zeta)$ если \mathcal{C} является несмещенным (то есть, $\mathbb{E} [\mathcal{C}(x)] = x$ для любых x) и второй момент ограничен^a

$$\mathbb{E} \left[\|\mathcal{C}(x)\|_2^2 \right] \leq \zeta \|x\|_2^2, \quad \forall x \in \mathbb{R}^d. \quad (0.5)$$

^a(0.5) Что также может быть переписанное в виде $\mathbb{E} [\|\mathcal{C}(x) - x\|_2^2] \leq (\zeta - 1) \|x\|_2^2$.

Предыдущие результаты: методы Top-k, Rand-K

Definition

Будем говорить, что $\mathcal{C} \in \mathbb{B}^1(\alpha, \beta)$ для некоторых $\alpha, \beta > 0$ если

$$\alpha \|x\|_2^2 \leq \mathbb{E} \left[\|\mathcal{C}(x)\|_2^2 \right] \leq \beta \langle \mathbb{E} [\mathcal{C}(x)], x \rangle, \quad \forall x \in \mathbb{R}^d. \quad (0.6)$$

Definition

Будем говорить, что $\mathcal{C} \in \mathbb{B}^2(\gamma, \beta)$ для некоторых $\gamma, \beta > 0$ если

$$\max \left\{ \gamma \|x\|_2^2, \frac{1}{\beta} \mathbb{E} \left[\|\mathcal{C}(x)\|_2^2 \right] \right\} \leq \langle \mathbb{E} [\mathcal{C}(x)], x \rangle \quad \forall x \in \mathbb{R}^d. \quad (0.7)$$

Definition

Будем говорить, что $\mathcal{C} \in \mathbb{B}^3(\delta)$ для некоторого $\delta > 0$ если

$$\mathbb{E} \left[\|\mathcal{C}(x) - x\|_2^2 \right] \leq \left(1 - \frac{1}{\delta} \right) \|x\|_2^2, \quad \forall x \in \mathbb{R}^d. \quad (0.8)$$

Для $k \in [d] \stackrel{\text{def}}{=} \{1, \dots, d\}$, несмещенный случайный (aka $\text{Rand-}k$) сжимающий определяется как :

$$\mathcal{C}(x) \stackrel{\text{def}}{=} \frac{d}{k} \sum_{i \in S} x_i e_i, \quad (0.9)$$

где $S \subseteq [d]$ подмножество индексов размера k выбранных равномерно, и e_1, \dots, e_d стандартные базисные вектора \mathbb{R}^d .

Lemma

$\text{Rand-}k$ оператор (0.9) принадлежит $U(\frac{d}{k})$.

Жадный (aka Top- k) сжимающий оператор определяется, как :

$$\mathcal{C}(x) \stackrel{\text{def}}{=} \sum_{i=d-k+1}^d x_{(i)} e_{(i)}, \quad (0.10)$$

где координаты упорядочены по их абсолютным значениям $|x_{(1)}| \leq |x_{(2)}| \leq \dots \leq |x_{(d)}|$.

Lemma

Top- k сжимающий оператор (0.10) принадлежит $\mathbb{B}^1(\frac{k}{d}, 1)$, $\mathbb{B}^2(\frac{k}{d}, 1)$, и $\mathbb{B}^3(\frac{d}{k})$.

Предлагается рассмотреть 3 новый оператора сжатия вида :

$$C(\nabla f(x)) = \text{Top}_k(\nabla f(x)|w) \quad (0.11)$$

, где w имеет смысл вектора "важности", а $\text{Top}_k(\circ|w)$ - выбор самых важных координат(обладающих самыми большими весами "важности" - координатами w_i). В зависимости от того, как определяется вектор w , меняется и сам оператор.

- (a) $w_i := f(x^k) - f(x^k - \eta [\nabla f(x^k)]_i).$
- (b) $w := \arg \min_{w \in \Delta} f(x^k - \eta \sum_{i=1}^d w_i [\nabla f(x^k)]_i).$
- (c) $w := \arg \min_{w \in [0,1]^d} f(x^k - \eta \sum_{i=1}^d w_i [\nabla f(x^k)]_i)$

Рассмотрим работу новых операторов на примере задачи логистической регрессии классическим методом и с использованием нейросетей. Следует обратить внимание, что на каждом шаге алгоритмов (b) и (c) необходимо решать свою задачу оптимизации. Конкретно для случая (b) воспользуемся методом зеркального спуска.

На каждой итерации спуска нам нужно решать задачу:

$$\min_{w \in \Delta^d} g_k(w),$$

где

$$g_k(w) := f\left(x^k - \eta w^T \cdot \mathcal{P}\right),$$

$$\mathcal{P} := \mathcal{I} \odot \nabla f\left(x^k\right) = ([\nabla f\left(x^k\right)]_1, [\nabla f\left(x^k\right)]_2, \dots, [\nabla f\left(x^k\right)]_d)$$

Решать ее будем методом зеркального спуска с KL-дивергенцией.

Решение:

$$w_i^{k+1} := \frac{w_i^k e^{-\eta [\nabla g(w^k)]_i}}{\sum_{j=1}^d w_j^k e^{-\eta [\nabla g(w^k)]_j}}$$

Здесь :

$$\nabla g(w) = -\eta \mathcal{P} \cdot \nabla f \left(x^k - \eta w^T \cdot \mathcal{P} \right)$$

Классическая логистическая регрессия :

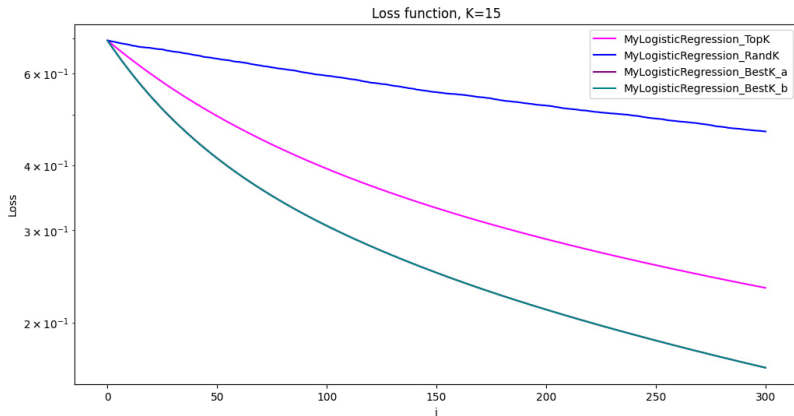


Рис.: Зависимость функции потерь от итерации

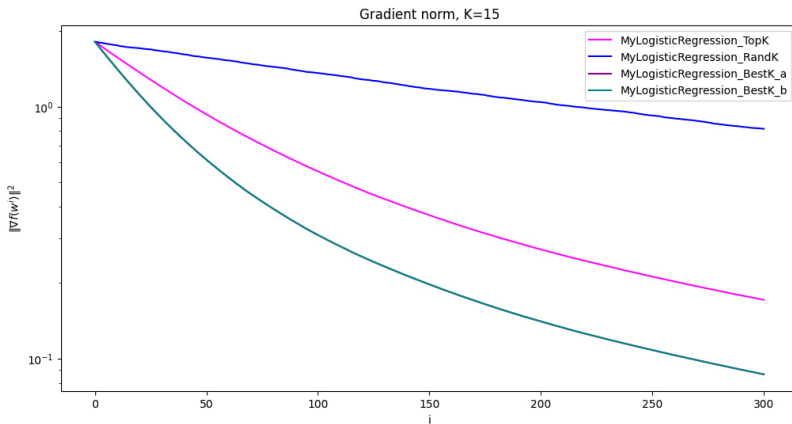


Рис.: Зависимость нормы градиента от итерации

Численные эксперименты

Теперь будем обучать нейросеть для датасета mushrooms.
Получаем следующие результаты :

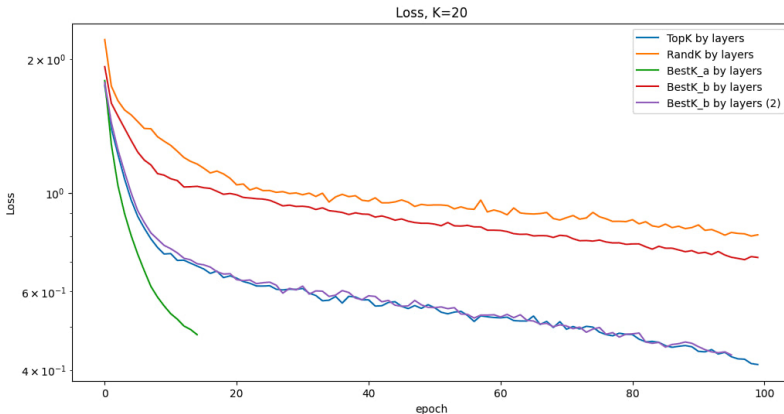


Рис.: Зависимость функции потерь от эпохи

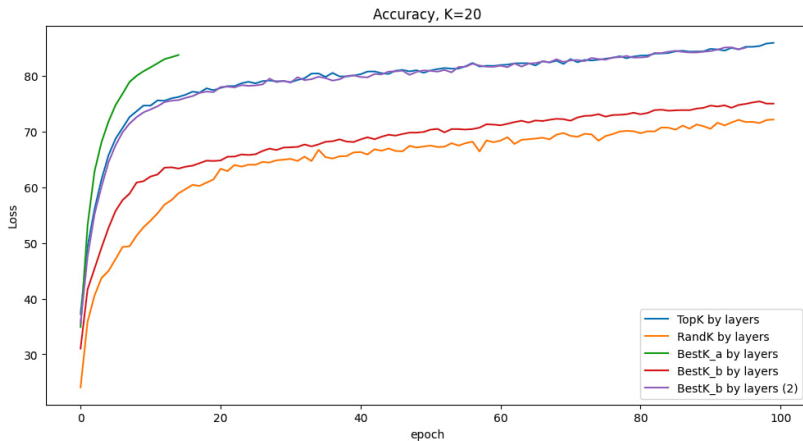


Рис.: Зависимость точности от эпохи

Хочу обратить внимание на реальное время для обучения :

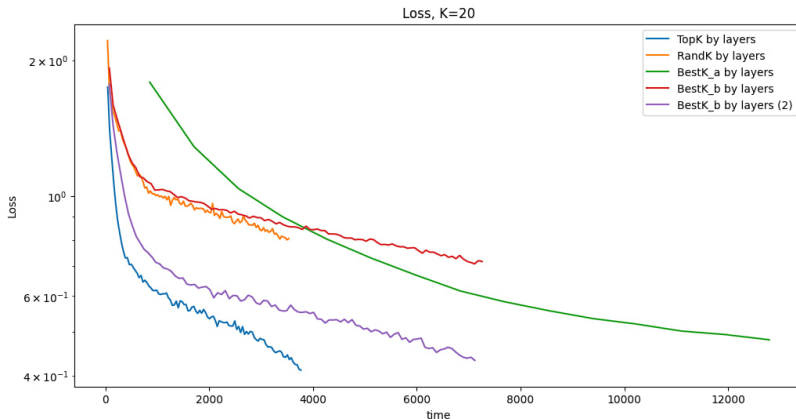


Рис.: Зависимость функции потерь от времени обучения

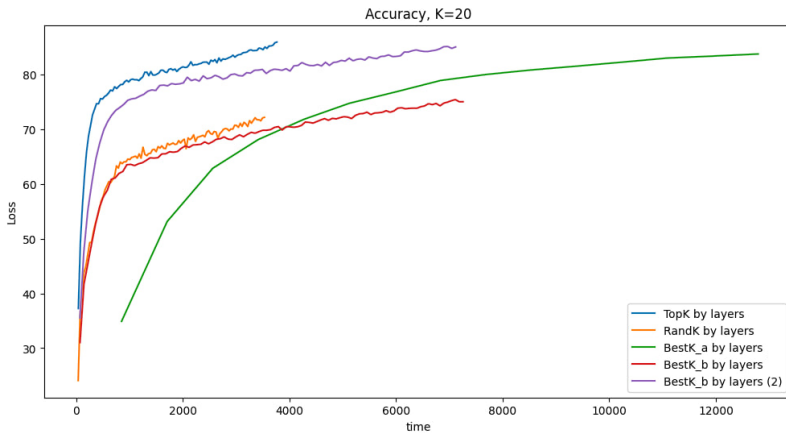


Рис.: Зависимость точности от времени

Можно сделать следующие выводы относительно результатов эксперимента

- 1° Методы действительно сходятся.
- 2° Методы показывают скорость сходимости более высокую, чем $Rand_k$, что согласуется с теорией.
- 3° Имеет место более высокая скорость сходимости методов (а), (б) и (с) в некоторых задачах. Однако они обходят вычислительно дороже, чем стандартные.

Подводя итоги работы можно заключить, что :

- 1° Предложены и опробованы 3 новых оператора сжатия
- 2° Получены гарантии сходимости этих методов
- 3° Показано превосходство новых методов над стандартными.