

Градиентные методы в задачах стохастической оптимизации

Юрий Владимирович Иванский

спбгу

Семинар по оптимизации, машинному обучению
и искусственному интеллекту

25 мая 2023

- 1 Достижение консенсуса
- 2 Алгоритм стохастической аппроксимации с пробным одновременным возмущением (SPSA)
- 3 Ускоренные градиентные методы (метод Поляка, метод Нестерова)

Определение

n агентов в сетевой системе достигают *консенсуса* в момент времени t , если $x_t^i = x_t^j \quad \forall i, j \in N, i \neq j$.

Определение

n агентов в сетевой системе достигают *асимптотического среднеквадратического ε -консенсуса* для $\varepsilon > 0$, если

$$\overline{\lim}_{t \rightarrow \infty} E \|x_t^i - x_t^j\|^2 \leq \varepsilon.$$

Динамика состояния агента (узла):

$$x_{t+1}^i = x_t^i + u_t^i$$

Протокол локального голосования:

$$u_t^i = \gamma \sum_{j \in N^i} b^{i,j} (x_t^j - x_t^i)$$

N^i — множество соседей, $b^{i,j}$ — вес ребра между узлами i и j

Консенсус, протокол локального голосования

Динамика состояния агента (узла):

$$x_{t+1}^i = x_t^i + u_t^i + z_t^i$$

Протокол локального голосования:

$$u_t^i = \gamma \sum_{j \in N_t^i} b_t^{i,j} (x_{t-d_t^{i,j}}^j + w_t^{i,j} - x_t^i - w_t^{i,i})$$

N_t^i — множество соседей узла i в момент времени t ,

z_t^i — внешнее неконтролируемое воздействие в момент времени t ,

$w_t^{i,j}$ — шум, возникающий при передаче данных по каналу связи от узла j на узел i в момент времени t ,

$b_t^{i,j}$ — весовой коэффициент матрицы смежности графа связей сетевой системы

$d_t^{i,j}$ — целочисленная задержка, возникающая при передаче данных по каналу связи от узла j на узел i в момент времени t

Динамика состояния агента (узла):

$$x_{t+1}^i = x_t^i + \gamma \sum_{j \in N^i} b^{ij} (x_t^j - x_t^i)$$

Динамика состояния сетевой системы:

$$X_{t+1} = X_t + \gamma B X_t - \gamma D(B) X_t$$

$$X_{t+1} = (I - \gamma L) X_t$$

Лапласиан:

$$L = D - B, \quad D = \text{diag}\{B * \mathbf{1}_n\}$$

Обобщенная постановка задачи

Цель распределенной оптимизации (как правило) — вычислить минимум некоторой целевой функции за счет взаимодействия между агентами

$$\bar{F}(x) = \sum_{i=1}^n F^i(x)$$

где $x \in \mathbb{R}^d$, а $F^i(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ — целевая функция агента i , известная только ему

Пусть $x_1, x_2, \dots \in \mathbb{R}^d$ — точки измерений

$$y_t = f(x_t, w_t) + v_t,$$

где $f : \mathbb{R}^d \times \mathbb{R}^p \rightarrow \mathbb{R}$,

w_t — случайный вектор, v_t — произвольная внешняя помеха.

Задача стохастической оптимизации:

$$F(x) = \int f(x, w)P(dw) \rightarrow \min_x.$$

Процедура Кифера-Вольфовица

Число наблюдений на итерации $M = 2d$

$$\hat{\theta}_0 \in \mathbb{R}^d$$

$$\hat{\theta}_n = \hat{\theta}_{n-1} - \frac{\alpha_n}{2\beta_n} (Y_n^+ - Y_n^-),$$

$$x_n^{(i,\pm)} = \hat{\theta}_{n-1} \pm \beta_n e_i$$

$$Y_n^\pm = \begin{pmatrix} f(x_n^{(1,\pm)}, w_n^{(1,\pm)}) + v_n^{(1,\pm)} \\ f(x_n^{(2,\pm)}, w_n^{(2,\pm)}) + v_n^{(2,\pm)} \\ \vdots \\ f(x_n^{(d,\pm)}, w_n^{(d,\pm)}) + v_n^{(d,\pm)} \end{pmatrix}$$

Рандомизированные алгоритмы стохастической аппроксимации (Simultaneous Perturbation Stochastic Approximation)

Сокращение числа наблюдений до 1 или 2 вместо $2d$ (!)

- Алгоритм с одним измерением

$$x_n = \hat{\theta}_{n-1} + \beta_n \Delta_n, \quad \Delta_n = \begin{pmatrix} \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{pmatrix}$$

$$y_n = f(x_n, w_n) + v_n$$

$$\hat{\theta}_n = \hat{\theta}_{n-1} - \frac{\alpha_n}{\beta_n} \mathcal{H}_n(\Delta_n) y_n$$

- Алгоритм с двумя измерениями

$$x_n^\pm = \hat{\theta}_{n-1} \pm \beta_n^\pm \Delta_n, \quad y_n^\pm = f(x_n^\pm, w_n^\pm) + v_n^\pm$$

$$\hat{\theta}_n = \hat{\theta}_{n-1} - \frac{\alpha_n}{\beta_n^+ + \beta_n^-} \mathcal{H}_n(\Delta_n) (y_n^+ - y_n^-)$$

Обоснование «псевдоградиентности» при произвольных внешних помехах

$$\begin{aligned} & E\left\{\hat{\theta}_{n-1} - \frac{\alpha_n}{\beta_n} \Delta_n y_n \mid \mathcal{F}_{n-1}\right\} = \\ &= \hat{\theta}_{n-1} - \frac{\alpha_n}{\beta_n} (E\{\Delta_n f(x_n) \mid \mathcal{F}_{n-1}\} + E\{\Delta_n\} E\{v_n \mid \mathcal{F}_{n-1}\}) = \\ &= \hat{\theta}_{n-1} - \frac{\alpha_n}{\beta_n} (E\{\Delta_n f(\hat{\theta}_{n-1} + \beta_n \Delta_n) \mid \mathcal{F}_{n-1}\}) \approx \\ &\approx \hat{\theta}_{n-1} - \frac{\alpha_n}{\beta_n} (E\{\Delta_n f(\hat{\theta}_{n-1}) + \frac{\beta_n \Delta_n \Delta_n^T \nabla f(\hat{\theta}_{n-1})}{2} \mid \mathcal{F}_{n-1}\}) = \\ &= \hat{\theta}_{n-1} - \frac{\alpha_n}{2} \nabla f(\hat{\theta}_{n-1}) \end{aligned}$$

- Состоятельность при почти произвольных внешних помехах [Граничин, 1989]
- Асимптотически-оптимальная скорость сходимости [Поляк, Цыбаков, 1990]
- Минимальное число наблюдений на итерации [Spall, 1992, 1997]
- Применимость в задаче об отслеживании изменений параметров [Granichin et al., 2009]

Отслеживание изменений параметров (трекинг)

Пусть Ξ — множество, $\{f_{\xi}(x, w)\}_{\xi \in \Xi}$ — семейство дифференцируемых функций

$$y_t = f_{\xi_t}(x_t, w_t) + v_t$$

$$F_t(x) = \int f_{\xi_t}(x, w_t) P(dw_t) P(d\xi_t) \rightarrow \min_x$$

Комбинированный алгоритм: SPSA + консенсус

[Granichin O. et al. 2021 Trans. Autom. Cont.]

Введем Δ_k^i , $k = 1, 2, \dots$, $i \in N$ — последовательность независимых случайных векторов из \mathbb{R}^d , называемых *одновременным пробным возмущением*, с симметричными функциями распределения $P_k^i(\cdot)$; и множество вектор-функций (ядер) $\mathcal{K}_k^i(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $k = 1, 2, \dots$

Выберем начальное приближение $\hat{\theta}_0^i \in \mathbb{R}^d$, шаг $\alpha > 0$, коэффициент усиления γ , и последовательности неотрицательных чисел $\{\beta_k^+\}$ и $\{\beta_k^-\}$ такие, что $\beta_k = \beta_k^+ + \beta_k^- > 0$. Рассмотрим алгоритм с двумя измерениями распределенных функций $f_{\xi_t}^i(\theta)$ для каждого агента $i \in N$, строящий последовательность точек измерения $\{x_t^i\}$ и оценки $\{\hat{\theta}_t^i\}$:

$$\begin{cases} x_{2k}^i = \hat{\theta}_{2k-2}^i + \beta_k^+ \Delta_k^i, x_{2k-1}^i = \hat{\theta}_{2k-2}^i - \beta_k^- \Delta_k^i, \\ \hat{\theta}_{2k-1}^i = \hat{\theta}_{2k-2}^i, \\ \hat{\theta}_{2k}^i = \hat{\theta}_{2k-1}^i - \alpha \left(\frac{y_{2k}^i - y_{2k-1}^i}{\beta_k} \mathcal{K}_k^i(\Delta_k^i) + \right. \\ \left. \gamma \sum_{j \in N_{2k-1}^i} b_{2k-1}^{i,j} (\tilde{\theta}_{2k-1}^{i,j} - \hat{\theta}_{2k-1}^i) \right) \end{cases}$$

Основной результат

Предположения:

- Градиент $\nabla f_{\xi}^i(x)$ удовлетворяет условию Липшица $\forall x_1, x_2 \in \mathbb{R}^d$

$$\|\nabla f_{\xi}^i(x_1) - \nabla f_{\xi}^i(x_2)\| \leq M\|x_1 - x_2\|$$

- Последовательные разности $\tilde{v}_k^i = v_{2k}^i - v_{2k-1}^i$ помех наблюдения ограничены $|\tilde{v}_k^i| \leq c_v < \infty$, or $\mathbb{E}(\tilde{v}_k^i)^2 \leq c_v^2$ если посл-ть $\{\tilde{v}_k^i\}$ случайная
- Помехи в канале связи $w_t^{i,j}$ случайные, независимые, одинаково распределенные $\mathbb{E}w_t^{i,j} = 0$, $\mathbb{E}\|w_t^{i,j}\|^2 \leq \sigma_w^2$. Все случайные векторы и переменные $w_t^{i,j}$, $b_t^{i,j}$, ξ_t , and ξ_{t+1} взаимно независимы (если они случайные)

Теорема

Последовательность оценок $\{\bar{\theta}_{2k}^i\}$ имеет асимптотически эффективную верхнюю границу $\bar{L} > 0$ невязки оценок: $\forall \varepsilon > 0 \exists \bar{k}$:

$$\forall k > \bar{k}: \sqrt{E\|\bar{\theta}_{2k}^i - 1_n \otimes \theta_{2k}\|^2} \leq \bar{L} + \varepsilon.$$

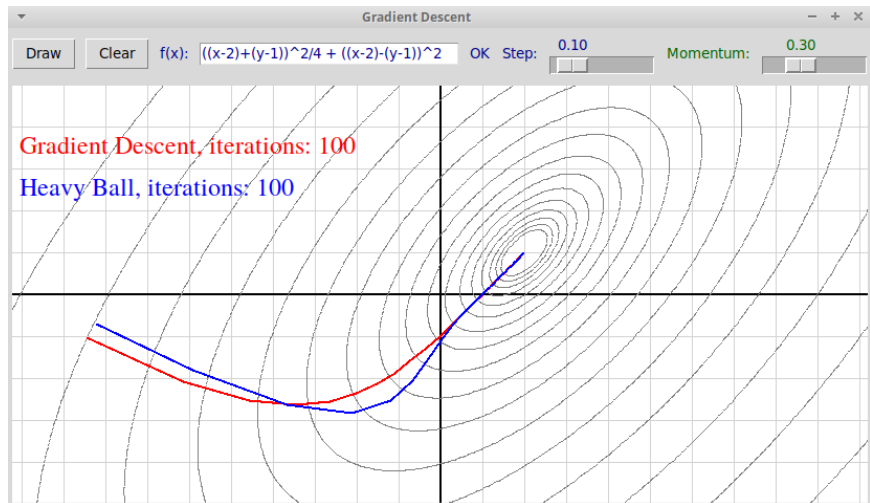
[Поляк, 1964, Ж. вычисл. матем. и матем. физ.]

Метод тяжелого шарика, метод инерции (momentum method, heavy ball method, momentum-gd)

Пусть $f(x)$ — непрерывно дифференцируемая функция, удовлетворяет условию Липшица для градиента с константой L :

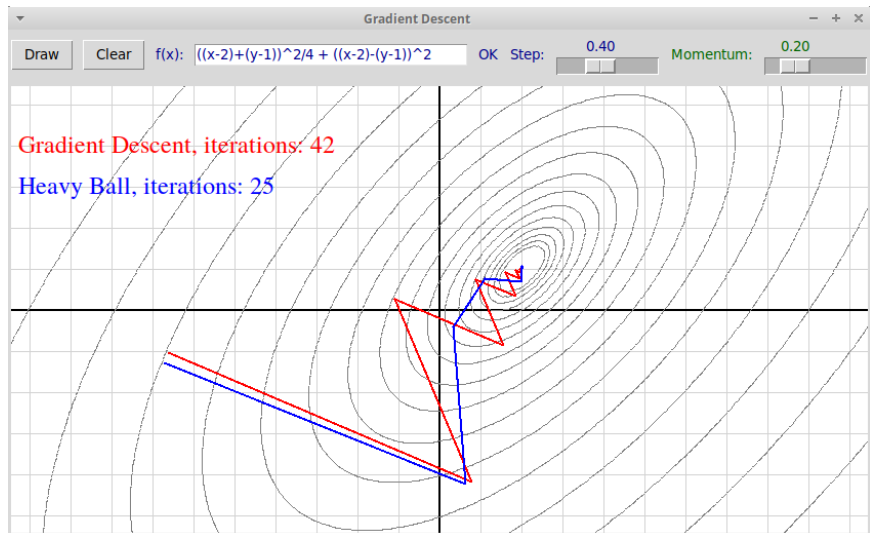
$$\|\nabla f(x_1) - \nabla f(x_2)\| \leq L\|x_1 - x_2\|$$

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}), \quad 0 \leq \beta < 1, \quad 0 \leq \alpha < 2(1 - \beta)/L$$



<http://mech.math.msu.su/~vvb/MasterAI/GradientDescent.html>

Метод Поляка



<http://mech.math.msu.su/~vvb/MasterAI/GradientDescent.html>

[Нестеров, 1983, Докл. АН СССР]

Пусть $f(x)$ — непрерывно дифференцируемая функция, удовлетворяет условию Липшица для градиента с константой L :

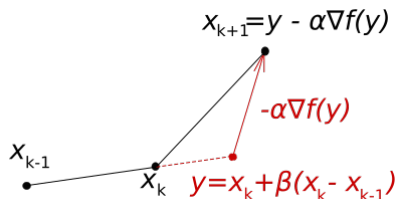
$$\|\nabla f(x_1) - \nabla f(x_2)\| \leq L\|x_1 - x_2\|$$

$$y_0 = x_0$$

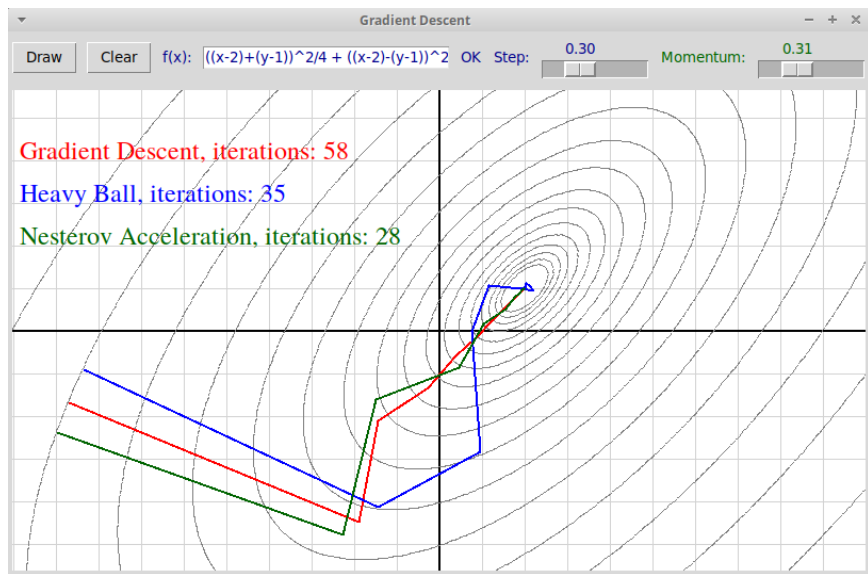
$$x_{k+1} = y_k - \alpha \nabla f(y_k)$$

$$y_{k+1} = x_{k+1} + \beta(x_{k+1} - x_k)$$

$$0 \leq \beta < 1, 0 \leq \alpha < 2(1 - \beta)/L$$



Метод Нестерова



<http://mech.math.msu.su/~vvb/MasterAI/GradientDescent.html>

[Erofeeva et al. 2022. American Control Conference]

$$\left\{ \begin{array}{l} \tilde{x}_{2k-2}^i = \frac{1}{\gamma_{k-1}^i + \alpha_k^i (\mu - \eta)} \left(\alpha_k^i \gamma_{k-1}^i z_{2k-2}^i + \gamma_k^i \hat{\theta}_{2k-2}^i \right), \\ x_{2k}^i = \tilde{x}_{2k-2}^i + \beta \Delta_k^i, \quad x_{2k-1}^i = \tilde{x}_{2k-2}^i - \beta \Delta_k^i, \\ \tilde{x}_{2k-1}^i = \tilde{x}_{2k-2}^i, \quad \hat{\theta}_{2k-1}^i = \hat{\theta}_{2k-2}^i, \\ g_{2k}^i = \Delta_k^i \frac{y_{2k}^i - y_{2k-1}^i}{2\beta} + \omega \sum_{j \in N^i} b^{i,j} (\tilde{x}_{2k-1}^i - \tilde{x}_{2k-1}^{i,j}), \\ \hat{\theta}_{2k}^i = \tilde{x}_{2k-1}^i - h g_{2k}^i, \\ z_{2k}^i = \frac{1}{\gamma_k^i} \left[(1 - \alpha_k^i) \gamma_{k-1}^i z_{2k-2}^i + \right. \\ \left. \alpha_k^i (\mu - \eta) \tilde{x}_{2k-1}^i - \alpha_k^i g_{2k}^i \right], \end{array} \right.$$

где $\alpha_k^i, \beta, \omega, \gamma_k^i, h, \mu, \eta$ – параметры и константы, x_t^i – точки наблюдения, Δ_k^i – пробное возмущение.

- Выбор параметров итеративной процедуры
- Оценка скорости сходимости
- Оптимизация размера шага, распространение оценки в децентрализованной сети

- Б.Т. Поляк, (1964) О некоторых способах ускорения сходимости итерационных методов // Журнал вычислительной математики и математической физики, т. 4, № 5, с. 791-803.
- Ю.Е. Нестеров, (1983) Метод минимизации выпуклых функций со скоростью сходимости $O(1/k^2)$ // Доклады Академии наук СССР, т. 269, № 3, с. 543-547.
- О.Н. Граничин, (1989) Об одной стохастической рекуррентной процедуре при зависимых помехах в наблюдении, использующей на входе пробные возмущения // Вестник Ленинградского университета. Серия 1: Математика, механика, астрономия, №4, с. 19-21.
- Б.Т. Поляк, А.Б. Цыбаков (1990) Оптимальные порядки точности поисковых алгоритмов стохастической оптимизации // Проблемы передачи информации, 26:2, с. 45–53.
- Spall, J.C. (1992) Multivariate Stochastic Approximation Using a Simultaneous Perturbation Gradient Approximation // IEEE Transactions on Automatic Control, vol. 37, pp. 332-341.

- Spall, J.C. (1997) A One-Measurement Form of Simultaneous Perturbation Stochastic Approximation // Automatica, vol. 33, pp. 109-112.
- Granichin, O., Gurevich, L., and Vakhitov, A. (2009). Discrete-time minimum tracking based on stochastic approximation algorithm with randomized differences // In Proceedings of the 48h IEEE Conference on Decision and Control (CDC), pp. 5763-5767. IEEE.
- Granichin, O., Erofeeva, V., Ivanskiy, Y., and Jiang, Y. (2020). Simultaneous perturbation stochastic approximation-based consensus for tracking under unknown-but-bounded disturbances. // IEEE Transactions on Automatic Control, 66(8), pp. 3710-3717.
- Erofeeva, V., Granichin, O., Tursunova, M., Sergeenko, A., and Jiang, Y. (2022). Accelerated simultaneous perturbation stochastic approximation for tracking under unknown-but-bounded disturbances. // In Proceedings of 2022 American Control Conference (ACC), pp. 1582-1587. IEEE.
- <https://distill.pub/2017/momentum/>