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

$L$ – гладкая , $\mu$ – сильно выпукла

$$O\left(\frac{L}{\mu} \log \frac{\|x^\circ - x^*\|_2}{\varepsilon}\right) \quad \text{итераций/оракульных вызов}$$

нужно градиентного спуск
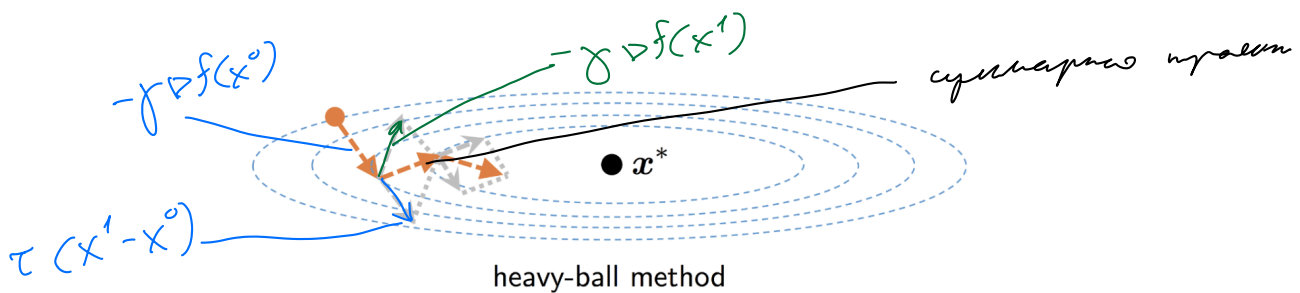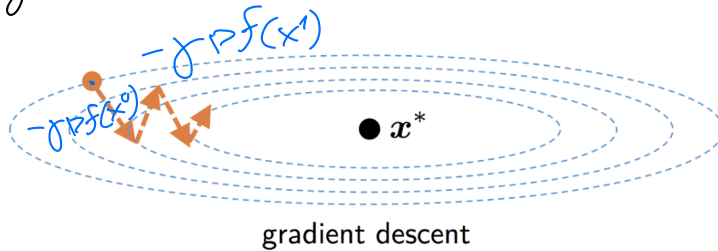
*а можно лучше?*

1964. Б.Т.Поляк

---

**Алгоритм 1** Метод тяжелого шарика

---

**Вход:** размер шагов $\{\gamma_k\}_{k=0} > 0$, моментумы $\{\tau_k\}_{k=0} \in [0;1]$, стартовая точка $x^0 = x^{-1} \in \mathbb{R}^d$, количество итераций $K$

1: **for** $k = 0, 1, \ldots, K-1$ **do**
2:     Вычислить $\nabla f(x^k)$
3:     $x^{k+1} = x^k - \gamma_k \nabla f(x^k) + \tau_k(x^k - x^{k-1})$
4: **end for**

**Выход:** $x^K$

---

Физический смысл:



gradient descent

$-\gamma \nabla f(x^0)$

$-\gamma \nabla f(x^1)$

суммарно пройдет

$\tau(x^1 - x^0)$

heavy-ball method

Pytorch (основная библиотека DL)

GD с моментум
$$\begin{cases} V^{k+1} = \beta V^k + \nabla f(x^k) \quad \beta \in [0;1] \\ x^{k+1} = x^k - \gamma V^{k+1} \end{cases}$$

$$x^{k+1} = x^k - \gamma \triangleright f(x^k) - \gamma \beta v^k$$

$$x^k = x^{k-1} - \gamma v^k \implies \gamma v^k = x^{k-1} - x^k$$

$$x^{k+1} = x^k - \gamma \triangleright f(x^k) + \beta(x^k - x^{k-1})$$

$\oplus$ :

+ интуиция и физика
+ легкость имплементации
+ дешевизна вычислений

$\ominus$ :

— два параметра для подбора $(\tau = [0,85; 0,95])$
— по памяти +1 вектор по сравнение с град. спуска
— в теории не лучше град. спуска

1983  Ю. Е. Нестеров

---

**Алгоритм 3** Ускоренный градиентный метод

**Вход:** размер шагов $\{\gamma_k\}_{k=0} > 0$, моментумы $\{\tau_k\}_{k=0} \in [0; 1]$, стартовая точка $x^0 = y^0 \in \mathbb{R}^d$, количество итераций $K$
1: **for** $k = 0, 1, \ldots, K-1$ **do**
2:   Вычислить $\nabla f(y^k)$
3:   $x^{k+1} = y^k - \gamma_k \nabla f(y^k)$
4:   $y^{k+1} = x^{k+1} + \tau_k(x^{k+1} - x^k)$
5: **end for**
**Выход:** $x^K$

---

Тяжелый шарик:
$$x^{k+1} = x^k - \gamma \triangleright f(x^k) + \tau(x^k - x^{k-1})$$

Ускоренный метод:  $y^k = x^k + \tau(x^k - x^{k-1})$

$$x^{k+1} = y^k - \gamma_k \triangleright f(y^k)$$

$$x^{k+1} = x^k + \tau(x^k - x^{k-1}) - \gamma_k \triangleright f(x^k + \tau(x^k - x^{k-1}))$$

т. ш.

Другой метод, ускоряющий градиентный спуск

---

**Алгоритм 4** Линейный каплинг: внутренний цикл

---

**Вход:** размер шагов $\{\gamma_k\}_{k=0} > 0$ и $\{\eta_k\}_{k=0} > 0$, моментумы $\{\tau_k\}_{k=0} \in [0; 1]$, стартовая точка $x^0 = y^0 = z^0 \in \mathbb{R}^d$, количество итераций $K$

1: **for** $k = 0, 1, \ldots, K-1$ **do**
2:      Вычислить $\nabla f(x^k)$
3:      $y^{k+1} = x^k - \eta_k \nabla f(x^k)$    $\leftarrow$ *стандартный шаг*
4:      $z^{k+1} = z^k - \gamma_k \nabla f(x^k)$    $\leftarrow$ *быстрый шаг*
5:      $x^{k+1} = \tau_k z^{k+1} + (1 - \tau_k) y^{k+1}$    $\leftarrow$ *вып. комб.*
6: **end for**
**Выход:** $\frac{1}{K} \sum_{k=0}^{K-1} x^k$

---

$$\frac{1 - \tau}{\tau} = \frac{\gamma}{\eta(2 - L\eta)} \quad \leftarrow \text{подбираем } \tau$$

$$\eta = \frac{1}{L} \qquad \gamma = \frac{1}{\sqrt{\mu L}}$$

$$f\left(\frac{1}{K} \sum_{k=1}^{K} x^k\right) - f^* \leq \sqrt{\frac{4L}{\mu K^2}} \left(f(x^0) - f^*\right)$$

$$\underset{\geqslant \frac{1}{2}}{}$$

$$K = 4\sqrt{\frac{L}{\mu}}$$

$$f\left(\frac{1}{K} \sum_{k=1}^{K} x^k\right) - f^* \leq \frac{1}{2}\left(f(x^0) - f^*\right)$$

Рестарт из новой стартовой точки

$$f\left(X_{final, new}\right) - f^* \leq \frac{1}{2}\left(f(X_{start, new}) - f^*\right)$$
$$\| \\ X_{final}$$

Тогда нужно $\log_2 \frac{f(x^0) - f^*}{\varepsilon}$ рестартов, чтобы гарантировать точно $\varepsilon$ по функции $f(x) - f^* \leq \varepsilon$

Результирующая сложность: $\underbrace{4\sqrt{\frac{L}{\mu}}}_{K} \log_2 \frac{f(x^0) - f^*}{\varepsilon}$

Ускор.
град.
метод
лучше же
оценки

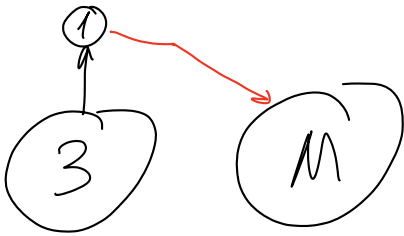лучше чем град спуск $\frac{L}{\mu} \log \ldots$

сходимость по функции:
- для сильно выпуклых задач $\approx$ эквивалент
- для ускор. метода Нестерова
  сходимость:
  $$\| x^k - x^* \|_2^2 + f(x^k) - f^*$$

Можно ли еще лучше? Нижние оценки



Класс методов:
- $x^0 = 0$ стартовая точка
  $$M_0 = \{ x^0 \} \leftarrow \text{послед. память}$$
- $\triangledown f(x^k) \qquad x^k \in M_k$
- $M_{k+1} = \text{span} \{ M_k, \triangledown f(x^k) \} \leftarrow \text{лин. оболочка}$
- $\overline{K}$ вызова оракула $\qquad x_{final} \in M_{\overline{K}}$

Плохая функция

$$f(x) = \frac{L-\mu}{8} x^\top A x + \frac{\mu}{2} \|x\|_2^2 - \frac{L-\mu}{4} e_1^\top x$$

$$A = \begin{pmatrix} 2 & -1 & & & O \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & & -1 \\ O & & & -1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{pmatrix} \quad e_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

1) $f$ — $L$-гладкая, $\mu$-сильно выпукла

$\|A\|_2 \leq 4$ $\qquad A \nleq 0$ $\qquad$ упражнение

2) $x^*$ — решение найден

$$\nabla f(x^*) = 0$$

$$\frac{L-\mu}{8} \cdot 2Ax + \frac{\mu}{2} \cdot 2x - \frac{L-\mu}{4} e_1 = 0$$

по-координатно:

$1.$ $\quad \frac{L-\mu}{4} (2x_1 - x_2) + \mu x_1 - \frac{L-\mu}{4} = 0$

другие $\quad$ $\boxed{\dfrac{L-\mu}{4} \left( -x_{k-1} + 2x_k - x_{(k+1)} \right) + \mu x_k = 0}$

последняя : $\quad \dots \dots$

$$x_{k+1} = \left( 2 + \frac{4\mu}{L-\mu} \right) x_k - x_{k-1}$$

линейная рекуррента

$$\lambda^2 = \left( 2 + \frac{4\mu}{L-\mu} \right) \lambda - 1$$

$\lambda$ решения : $\lambda_1 = \dfrac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}$ $\qquad \lambda_2 = \dots$

$$X_k = C_1 \lambda_1^k + C_2 \lambda_2^k$$

зададим нач. условия так $C_2 = 0$ $C_1 = 1$

$$X_k = \left( \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^k$$

3)
$$\nabla f(x) = \frac{L - \mu}{8} \cdot 2AX + \frac{\mu}{2} \cdot 2X - \frac{L - \mu}{4} e_1$$

1. $x^0 = 0$    $\nabla f(x^0) \in \text{span}\{e_1\}$

только 1ая коорд не нулевая

$$M_1 \subseteq \text{span}\{e_1\}$$

2. $x^1 \in M_1$    $\nabla f(x^1) \in \text{span}\{e_1, e_2\}$

только 1ая и 2ая коорд
ненулевые

$$M_2 \subseteq \text{span}\{e_1, e_2\}$$

k. $x^{k-1} \in M_{k-1}$    $\nabla f(x^{(k-1)}) \in \text{span}\{e_1 .. e_k\}$

Summary:

$$X_k^* = \left( \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^k = q^k$$

после $\bar{k}$ шагов оракула только $\bar{k}$ первых
коорд
ненулевые

$$\| X_{final}^{\bar{k}} - X^* \|^2 = \sum_{i=1}^{d} \left( X_{find, i}^{\bar{k}} - X_i^* \right)^2$$

$d = 2\bar{k}$ разбы не все коорд.

в лучшем случае между указан $\bar{k}$ коорд,

занулив 0

$$\geqslant \sum_{i=\bar{K}+1}^{2\bar{K}=d} \left(X_{final,i}^{\bar{K}} - X_i^*\right)^2$$

$$\|X_{final}^{\bar{K}} - X^*\|^2 = \sum_{i=\bar{K}+1}^{2\bar{K}} (X_i^*)^2 \overset{\|}{\underset{0}{}} = \sum_{i=\bar{K}+1}^{2\bar{K}} q^{2i}$$

$$\|X^0 - X^*\|_2^2 = \sum_{i=1}^{2\bar{K}} (X_i^*)^2 = \sum_{i=1}^{2\bar{K}} q^{2i}$$

$$\overset{\|}{\underset{0}{}}$$

$$\sum_{i=1}^{2\bar{K}} q^{2i} = \sum_{i=1}^{\bar{K}} q^{2i} + \sum_{i=\bar{K}+1}^{2\bar{K}} q^{2i} = \left(1 + q^{2\bar{K}}\right) \sum_{i=1}^{\bar{K}} q^{2i}$$

$$q^{2\bar{K}} \sum_{i=1}^{\bar{K}} q^{2i}$$

$$\sum_{i=\bar{K}+1}^{2\bar{K}} q^{2i} = q^{2\bar{K}} \sum_{i=1}^{\bar{K}} q^{2i}$$

$$\|X_{final}^{\bar{K}} - X^*\|_2^2 \geqslant \frac{q^{2\bar{K}}}{1+q^{2\bar{K}}} \|X^0 - X^*\|_2^2$$

$$\geqslant \frac{q^{2\bar{K}}}{1+1} \|X^0 - X^*\|_2^2 \sim \left(1 - \sqrt{\tfrac{\mu}{L}}\right)^K \frac{\|X^0 - X^*\|}{2}$$

$$K \quad \text{через} \quad \varepsilon \qquad K \sim \sqrt{\tfrac{L}{\mu}} \log \ldots$$

**Нижняя оценка на оракульную сложность**

Для любого метода из класса, описанного выше, существует безусловная задача оптимизации с $L$-гладкой, $\mu$-сильно выпуклой целевой функцией $f$ такая, что для решения этой задачи методу необходимо

$$\Omega\left(\sqrt{\tfrac{L}{\mu}} \log \frac{\|x^0 - x^*\|_2}{\varepsilon}\right) \text{ вызовов оракула.}$$

поэтому ⟶ оптимальность ускоренных методов