

ВЛИЯНИЕ МОМЕНТА В АЛГОРИТМАХ ОПТИМИЗАЦИИ

РАЗБОР РАБОТЫ GON G. WHY MOMENTUM REALLY WORKS //DISTILL. – 2017. – Т. 2. – №. 4. – С. E6.

Докладчик: Денишева Рушана

24.04.23, MIPT

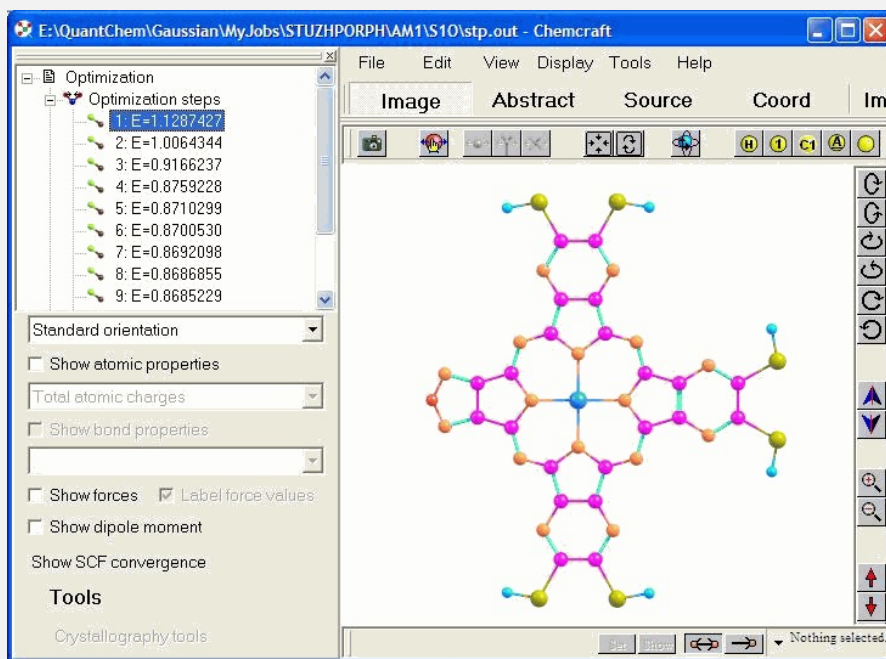
СОДЕРЖАНИЕ

- Описание задачи оптимизации
 - Описание метода градиентного спуска
- Описание решения задачи путем добавления момента
- Преимущества градиентного спуска с моментом
- Сильные и слабые стороны метода

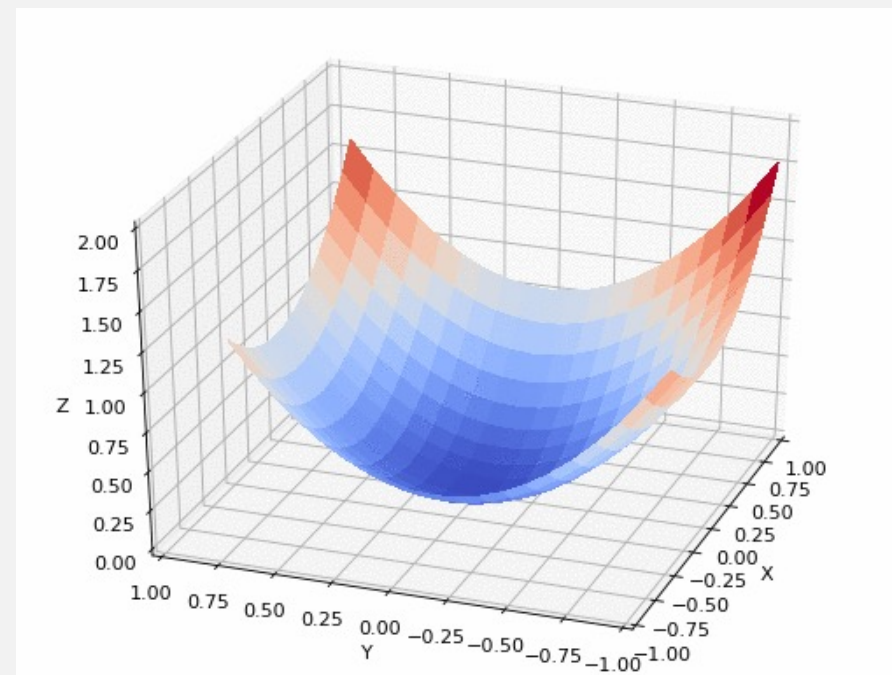
ЗАДАЧА ОПТИМИЗАЦИИ

Задача оптимизации – это задача нахождения экстремума целевой функции с учетом ограничений на управляемые переменные.

ЗАДАЧА ОПТИМИЗАЦИИ



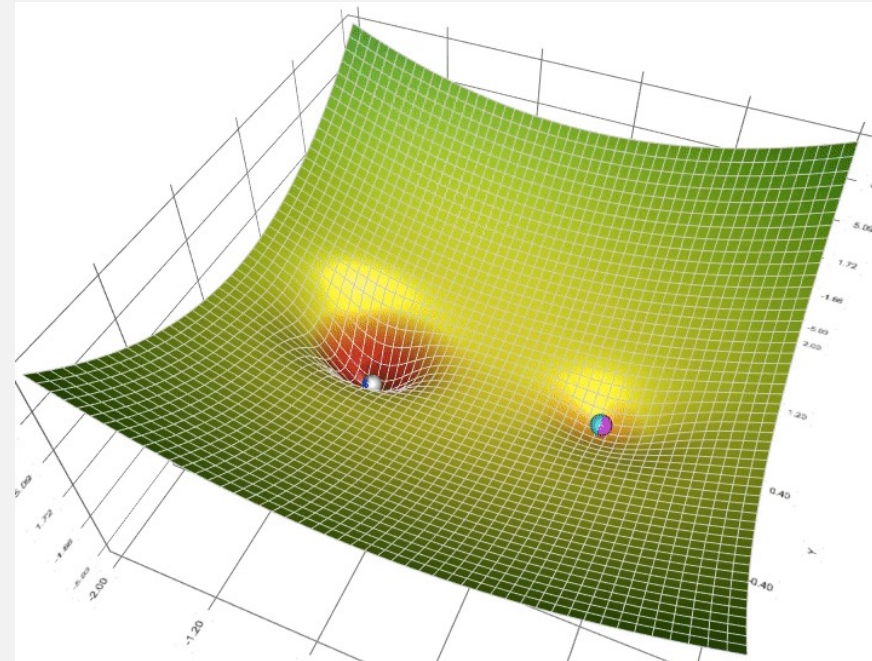
Оптимизация в медицине



Оптимизация функции потерь

ГРАДИЕНТНЫЙ СПУСК

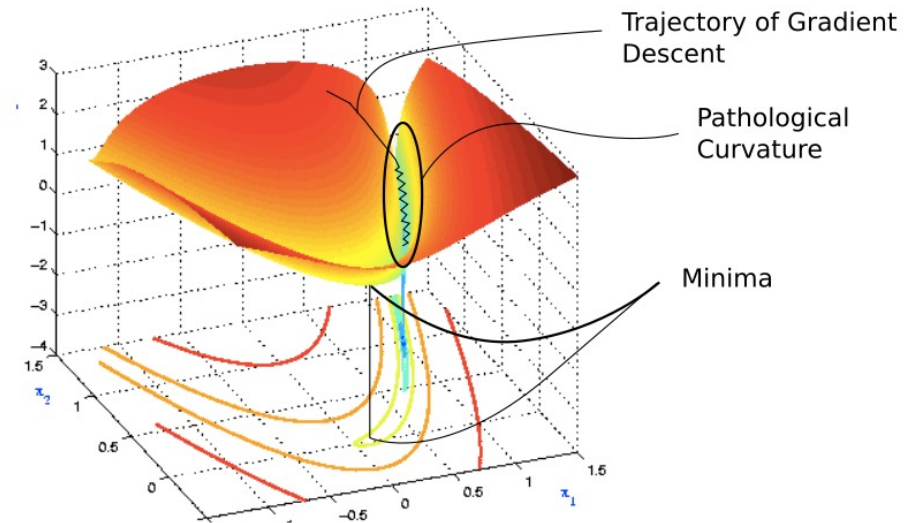
Градиентный спуск – численный метод нахождения *локального* минимума или максимума функции с помощью движения вдоль градиента, один из основных численных методов современной оптимизации.



Анимация 5 методов градиентного спуска на поверхности: градиентный спуск (голубой), импульс (пурпурный), AdaGrad (белый), RMSProp (зеленый), Adam (синий). Левая лунка - глобальный минимум; правая лунка - локальный минимум.

ПРОБЛЕМА

- “Застрять” в локальном минимуме, седле целевой функции или патологической кривизне



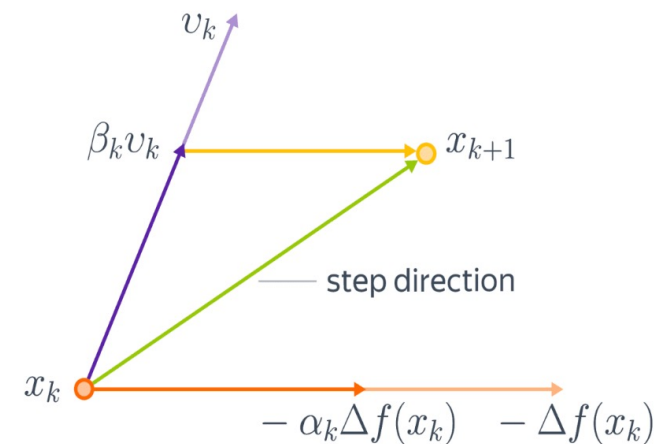
РЕШЕНИЕ – МЕТОД ИНЕРЦИИ

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

Градиентный спуск

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

Градиентный спуск с добавлением “момента”



Momentum

ПРЕИМУЩЕСТВА МОМЕНТА

- Позволяет использовать больший диапазон размеров шага и создает собственные колебания
- Дает квадратичное ускорение на многих функциях

ОЦЕНКА СКОРОСТИ СХОДИМОСТИ К МИНИМУМУ И ШАГА

ВЫПУКЛАЯ КВАДРАТИЧНАЯ ФУНКЦИЯ

Без момента

$$f(w) = \frac{1}{2}w^T A w - b^T w, \quad w \in \mathbf{R}^n.$$

$$A = Q \operatorname{diag}(\lambda_1, \dots, \lambda_n) Q^T, \quad Q = [q_1, \dots, q_n],$$

$$f(w^k) - f(w^*) = \sum (1 - \alpha \lambda_i)^{2k} \lambda_i [x_i^0]^2$$

С моментом

$$\begin{aligned} z^{k+1} &= \beta z^k + (Aw^k - b) \\ w^{k+1} &= w^k - \alpha z^{k+1}. \end{aligned}$$

$$\begin{pmatrix} y_i^k \\ x_i^k \end{pmatrix} = R^k \begin{pmatrix} y_i^0 \\ x_i^0 \end{pmatrix} \quad R = \begin{pmatrix} \beta & \lambda_i \\ -\alpha\beta & 1 - \alpha\lambda_i \end{pmatrix}.$$

где $x^k = Q(w^k - w^*)$ и $y^k = Qz^k$

ОЦЕНКА СКОРОСТИ СХОДИМОСТИ К МИНИМУМУ И ШАГА

ВЫПУКЛАЯ КВАДРАТИЧНАЯ ФУНКЦИЯ

Без момента

$$0 < \alpha \lambda_i < 2.$$

С моментом

$$0 < \alpha \lambda_i < 2 + 2\beta \quad \text{for} \quad 0 \leq \beta < 1$$



Момент позволяет увеличить размер шага в 2 раза перед расхождением!

ОЦЕНКА СКОРОСТИ СХОДИМОСТИ К МИНИМУМУ И ШАГА

ВЫПУКЛАЯ КВАДРАТИЧНАЯ ФУНКЦИЯ

Без момента

$$\begin{aligned}\text{optimal } \alpha &= \underset{\alpha}{\operatorname{argmin}} \operatorname{rate}(\alpha) = \frac{2}{\lambda_1 + \lambda_n} \\ \text{optimal rate} &= \min_{\alpha} \operatorname{rate}(\alpha) = \frac{\lambda_n/\lambda_1 - 1}{\lambda_n/\lambda_1 + 1}\end{aligned}$$

$$\text{condition number} := \kappa := \frac{\lambda_n}{\lambda_1}$$

С моментом

$$\alpha = \left(\frac{2}{\sqrt{\lambda_1} + \sqrt{\lambda_n}} \right)^2 \quad \beta = \left(\frac{\sqrt{\lambda_n} - \sqrt{\lambda_1}}{\sqrt{\lambda_n} + \sqrt{\lambda_1}} \right)^2$$

$$\text{optimal rate} = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$$

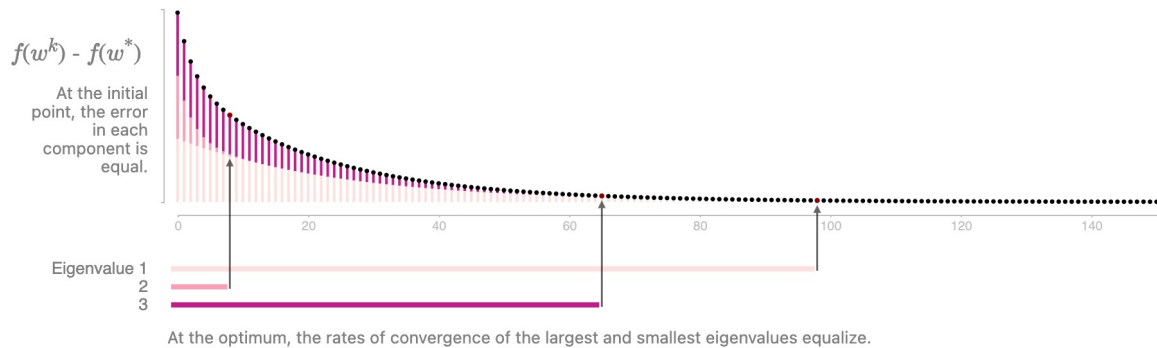


Момент дал квадратичное ускорение на нашей функции!

ОЦЕНКА СКОРОСТИ СХОДИМОСТИ К МИНИМУМУ И ШАГА

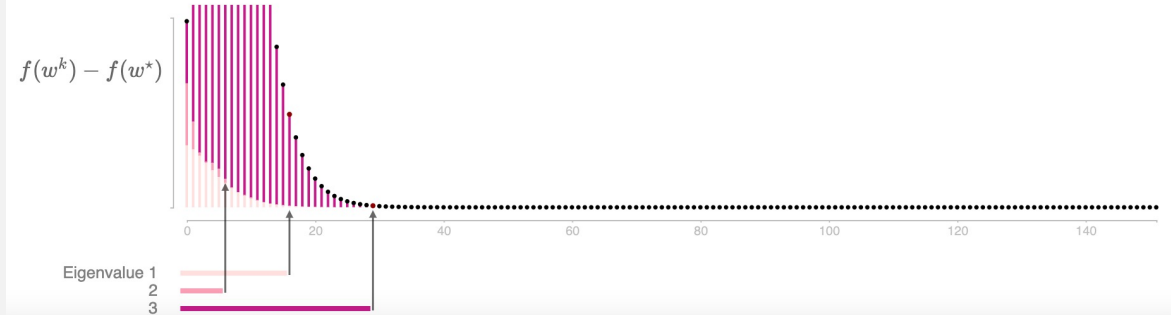
ВЫПУКЛАЯ КВАДРАТИЧНАЯ ФУНКЦИЯ

Optimization can be seen as combination of several component problems, shown here as 1 2 3 with eigenvalues $\lambda_1 = 0.01$, $\lambda_2 = 0.1$, and $\lambda_3 = 1$ respectively.



Без момента

We can do the same decomposition here with momentum, with eigenvalues $\lambda_1 = 0.01$, $\lambda_2 = 0.1$, and $\lambda_3 = 1$. Though the decrease is no longer monotonic, but significantly faster.



С моментом

ПРИМЕР: ПРОБЛЕМА КОЛОРИЗАЦИИ

G - граф с вершинами в виде пикселей,

E - множество ребер, соединяющих каждый пиксель с четырьмя соседними пикселями,

D - небольшое множество из нескольких выделенных вершин

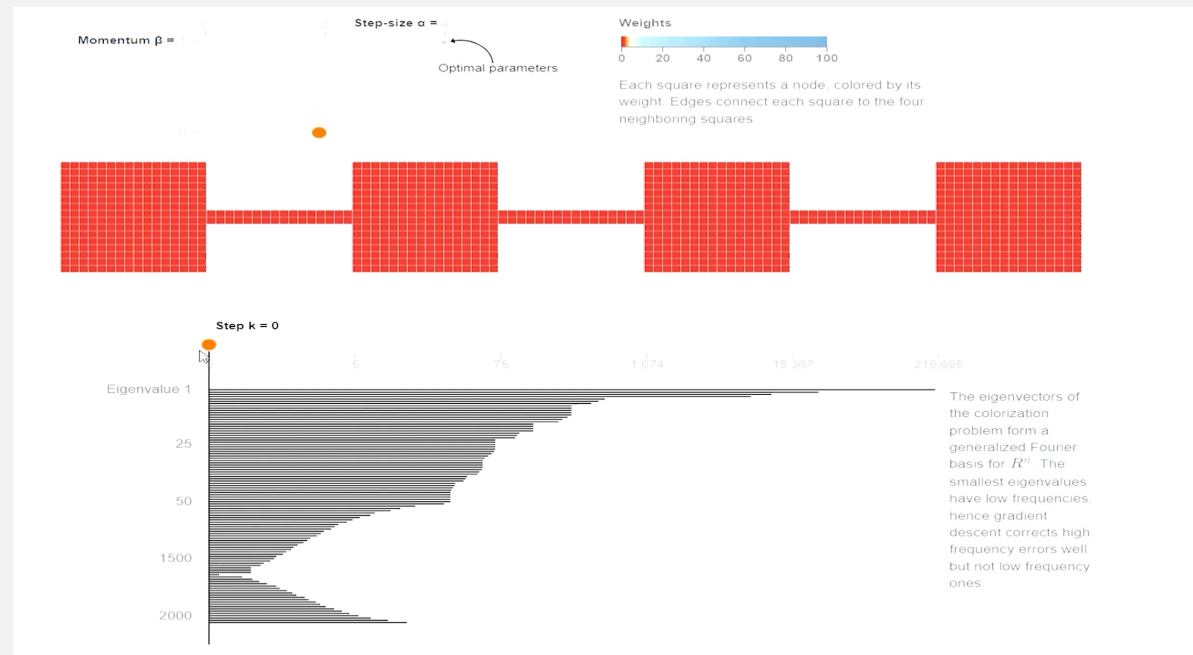
$$\text{minimize} \quad \frac{1}{2} \sum_{i \in D} (w_i - 1)^2 \quad + \quad \frac{1}{2} \sum_{i,j \in E} (w_i - w_j)^2.$$

The **colorizer** pulls
distinguished pixels
towards 1

The **smoother** spreads
out the color

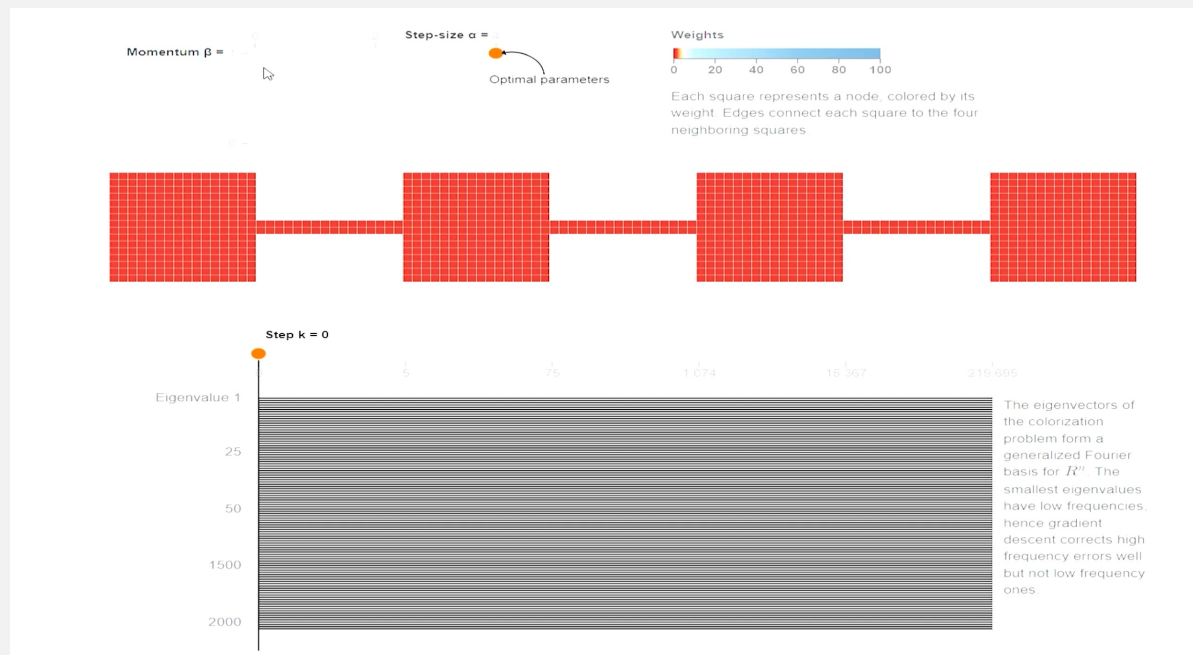
ПРИМЕР: ПРОБЛЕМА КОЛОРИЗАЦИИ

Без момента



ПРИМЕР: ПРОБЛЕМА КОЛОРИЗАЦИИ

С моментом



ГРАНИЦЫ УЛУЧШЕНИЯ СХОДИМОСТИ ГРАДИЕНТНОГО СПУСКА

$$w^{k+1} = w^0 + \sum_i^k \Gamma_i^k \nabla f(w^i) \quad \text{for some diagonal matrix } \Gamma_i^k.$$

Алгоритмическое пространство

Рассмотрим граф из одного пути

$$f^n(w) = \frac{1}{2} (w_1 - 1)^2 + \frac{1}{2} \sum_{i=1}^n (w_i - w_{i+1})^2 + \frac{2}{\kappa - 1} \|w\|^2.$$

with a colorizer of
one node

strong couplings of adjacent
nodes in the path,

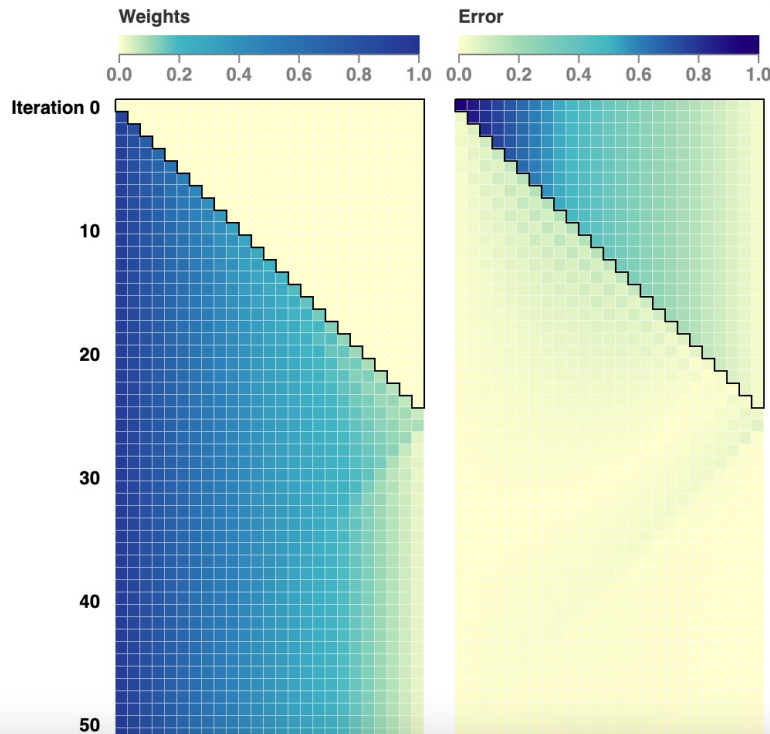
and a small
regularization term.

$$w_i^* = \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^i$$

Оптимальное решение

ГРАНИЦЫ УЛУЧШЕНИЯ СХОДИМОСТИ ГРАДИЕНТНОГО СПУСКА

Here we see the first 50 iterates of momentum on the Convex Rosenbrock for $n = 25$. The behavior here is similar to that of any Linear First Order Algorithm.



This triangle is a “dead zone” of our iterates. The iterates are always 0, no matter what the parameters.

The remaining expanding space is the “light cone” of our iterate’s influence. Momentum does very well here with the optimal parameters.



Для выпуклой квадратичной функции:

$$w^k - w^* = Qx^k = \sum_i^n x_i^0 (1 - \alpha \lambda_i)^k q_i$$

$$\begin{aligned} \|w^k - w^*\|_\infty &\geq \max_{i \geq k+1} \{|w_i^*|\} \\ &= \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{k+1} \\ &= \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|w^0 - w^*\|_\infty. \end{aligned}$$



**Максимальное улучшение сходимости
градиентного спуска на квадратичный
коэффициент!**

ЗАКЛЮЧЕНИЕ

- **Сильные стороны метода инерции:**

1. выход из локального минимума целевой функции или патологической кривизны
2. позволяет использовать больший диапазон размеров шага и создает собственные колебания
3. дает квадратичное ускорение на многих функциях

- **Слабые стороны метода инерции:**

1. фиксированный размер шага

ЗАКЛЮЧЕНИЕ

- **Практическая ценность статьи** – на методе градиентного спуска с моментом основаны более мощные алгоритмы
- **Научная ценность статьи** - получена численная оценка скорости сходимости к минимуму функции

СПАСИБО ЗА ВНИМАНИЕ!

denisheva.rr@phystech.edu