

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

Выпуклый
надкай:

$$f(x_k) - f^* \leq \frac{LR^2}{2K}$$

$O\left(\frac{1}{K}\right)$ Лучше:
 $O\left(\frac{1}{K^2}\right)$

Рассматривается классическая задача выпуклой оптимизации:

$$\min_{x \in S} f(x),$$

$$S = \mathbb{R}^n$$

Подразумевается, что $f(x)$ - выпуклая функция на выпуклом множестве S . Для начала будем рассматривать задачу безусловной минимизации (БМ), $S = \mathbb{R}^n$

Вектор g называется **субградиентом** функции $f(x) : S \rightarrow \mathbb{R}$ в точке x_0 , если $\forall x \in S$:

$$f(x) \geq f(x_0) + \langle g, x - x_0 \rangle \rightarrow \langle g, x - x_0 \rangle \leq f(x) - f(x_0)$$

Градиентный спуск предполагает, что функция $f(x)$ является дифференцируемой в каждой точке задачи. Теперь же, мы будем предполагать лишь выпуклость.

Итак, мы имеем оракул первого порядка:

Вход: $x \in \mathbb{R}^n$

Выход: $\partial f(x)$ и $f(x)$

Algorithm

$$x_{k+1} = x_k - \alpha_k g_k,$$

(SD)

где g_k - произвольный субградиент функции $f(x)$ в т. x_k , $g_k \in \partial f(x_k)$

Bounds

Vanilla version

Запишем как близко мы подошли к оптимуму $x^* = \arg \min_{x \in \mathbb{R}^n} f(x) = \arg f^*$ на последней итерации:

$$g_k^2 = \|g_k\|^2$$

$$\begin{aligned} \|x_{k+1} - x^*\|^2 &= \|x_k - x^* - \alpha_k g_k\|^2 \\ &= \|x_k - x^*\|^2 + \alpha_k^2 g_k^2 - 2\alpha_k \langle g_k, x_k - x^* \rangle \end{aligned}$$

Для субградиента: $\langle g_k, x_k - x^* \rangle \leq f(x_k) - f(x^*) = f(x_k) - f^*$. Из написанного выше:

$$2\alpha_k \langle g_k, x_k - x^* \rangle = \|x_k - x^*\|^2 + \alpha_k^2 g_k^2 - \|x_{k+1} - x^*\|^2$$

Просуммируем полученное неравенство для $k = 0, \dots, T-1$

(УМ)

$$\begin{aligned}
\sum_{k=0}^{T-1} 2\alpha_k \langle g_k, x_k - x^* \rangle &= \|x_0 - x^*\|^2 - \|x_T - x^*\|^2 + \sum_{k=0}^{T-1} \alpha_k^2 g_k^2 \\
&\leq \|x_0 - x^*\|^2 + \sum_{k=0}^{T-1} \alpha_k^2 g_k^2 \\
&\leq R^2 + G^2 \sum_{k=0}^{T-1} \alpha_k^2
\end{aligned}$$

Здесь мы предположили $R^2 = \|x_0 - x^*\|^2$, $\|g_k\| \leq G$.
Предполагая $\alpha_k = \alpha$ (постоянный шаг), имеем:

$$\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle \leq \frac{R^2}{2\alpha} + \frac{\alpha}{2} G^2 T$$

Минимизация правой части по α дает $\alpha^* = \frac{R}{G} \sqrt{\frac{1}{T}}$

*К сожалению
Гензвестно
заранее
(G, R тоже)*

$$\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle \leq GR\sqrt{T}$$

Тогда (используя неравенство Йенсена и свойство субградиента $f(x^*) \geq f(x_k) + \langle g_k, x^* - x_k \rangle$)
запишем оценку на т.н. *Regret*, а именно:

$$f(\bar{x}) - f^* \leq \langle g_k, x_k - x^* \rangle$$

$$\begin{aligned}
\bar{x} &= \frac{1}{T} \sum_{k=0}^{T-1} x_k \\
f(\bar{x}) - f^* &= f\left(\frac{1}{T} \sum_{k=0}^{T-1} x_k\right) - f^* \leq \frac{1}{T} \left(\sum_{k=0}^{T-1} (f(x_k) - f^*) \right) \\
&\leq \frac{1}{T} \left(\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle \right) \\
&\leq GR \frac{1}{\sqrt{T}}
\end{aligned}$$

$$f(\bar{x}) - f^* \leq \frac{GR}{\sqrt{T}}$$

Важные моменты:

- Получение оценок не для x_T , а для среднего арифметического по итерациям \bar{x} - типичный трюк при получении оценок для методов, где есть выпуклость, но нет удобного убывания на каждой итерации. Нет гарантий успеха на каждой итерации, но есть гарантия успеха в среднем
- Для выбора оптимального шага необходимо знать (предположить) число итераций заранее.
Возможный выход: инициализировать T небольшим значением, после достижения этого количества итераций удваивать T и *рестартовать* алгоритм. Более интеллектуальный способ: адаптивный выбор длины шага.

Steepest subgradient descent

$$x_{k+1} = x_k - \alpha_k g_k$$

Попробуем выбирать на каждой итерации длину шага более оптимально. Тогда:

$$\|x_{k+1} - x^*\|^2 = \|x_k - x^*\|^2 + \alpha_k^2 g_k^2 - 2\alpha_k \langle g_k, x_k - x^* \rangle$$

Минимизируя выпуклую правую часть по α_k , получаем:

$$\text{Раньше: } \alpha_k = \arg \min_{\alpha \in \mathbb{R}} f(x_k - \alpha g_k)$$

$$\alpha_k = \frac{\langle g_k, x_k - x^* \rangle}{\|g_k\|^2}$$

Оценки изменяются следующим образом: $\|x_{k+1} - x^*\|^2 = \|x_k - x^*\|^2 + d_k^2 g_k^2 - 2d_k \langle g_k, x_k - x^* \rangle$

$$\|x_{k+1} - x^*\|^2 = \|x_k - x^*\|^2 - \frac{\langle g_k, x_k - x^* \rangle^2}{\|g_k\|^2}$$

$$\langle g_k, x_k - x^* \rangle^2 = (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2) \|g_k\|^2$$

$$\langle g_k, x_k - x^* \rangle^2 \leq (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2) G^2$$

$$\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle^2 \leq \sum_{k=0}^{T-1} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2) G^2$$

суммируем.

ЕСТЬ:

$$\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle^2 \leq (\|x_0 - x^*\|^2 - \|x_T - x^*\|^2) G^2 \rightarrow \sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle^2 \leq R^2 G^2$$

$$\left(\sum_{i=1}^n x_i \right)^2 \leq n \left(\sum x_i^2 \right)$$

$$\frac{1}{T} \left(\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle \right)^2 \leq \sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle^2 \leq R^2 G^2$$

$$(Ax)^2 \leq A^T A \cdot x^T x$$

значит,

$$(a^T b)^2 \leq a^T a \cdot b^T b$$

$$|a^T b| \leq \|a\| \cdot \|b\|$$

$$\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle \leq GR\sqrt{T}$$

НАДО:

Что приводит к абсолютно такой же оценке $\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$ на невязку по значению функции. На самом деле,

для такого класса функций нельзя получить результат лучше, чем $\frac{1}{\sqrt{T}}$ или $\frac{1}{\varepsilon^2}$ по итерациям

$$\frac{1}{\sqrt{T}}$$

$$\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

Online learning

Рассматривается следующая игра: есть игрок и природа. На каждом из $k = 0, \dots, T-1$ шагов:

- Игрок выбирает действие x_k
- Природа (возможно, враждебно) выбирает выпуклую функцию f_k , сообщает игроку значение $f(x_k), g_k \in \partial f(x_k)$
- Игрок вычисляет следующее действие, чтобы минимизировать регрет:

$$x_{k+1} = x_k - d g_k$$

$$R_{T-1} = \sum_{k=0}^{T-1} f_k(x_k) - \min_x \sum_{k=0}^{T-1} f_k(x)$$

лучший x для мин. f . (Regret).

В такой постановке цель игрока состоит в том, чтобы выбрать стратегию, которая минимизирует разницу его действия с наилучшим выбором на каждом шаге.

Несмотря на весьма сложную (на первый взгляд) постановку задачи, существует стратегия, при которой регрет растет как \sqrt{T} , что означает, что усредненный регрет $\frac{1}{T} R_{T-1}$ падает, как $\frac{1}{\sqrt{T}}$

Если мы возьмем оценку (Subgradient Bound) для субградиентного метода, полученную выше, мы имеем:

$$\sum_{k=0}^T f_k(x_k) - \sum_{k=0}^T f_k(x^*) = \sum_{k=0}^T [f_k(x_k) - f_k(x^*)]$$

$$\sum_{k=0}^{T-1} \langle g_k, x_k - x^* \rangle \leq G \|x_0 - x^*\| \sqrt{T}$$

СРФТ

Однако, в её выводе мы нигде не использовали тот факт, что $x^* = \arg \min_{x \in S} f(x)$. Более того, мы вообще не использовали никакой специфичности точки x^* . Тогда можно записать это для произвольной точки y :

$$\sum_{k=0}^{T-1} \langle g_k, x_k - y \rangle \leq G \|x_0 - y\| \sqrt{T}$$

Запишем тогда оценки для регрета, взяв $y = \arg \min_{x \in S} \sum_{k=0}^{T-1} f_k(x)$:

$$\begin{aligned} R_{T-1} &= \sum_{k=0}^{T-1} f_k(x_k) - \min_x \sum_{k=0}^{T-1} f_k(x) = \sum_{k=0}^{T-1} f_k(x_k) - \sum_{k=0}^{T-1} f_k(y) = \\ &= \sum_{k=0}^{T-1} (f_k(x_k) - f_k(y)) \leq \sum_{k=0}^{T-1} \langle g_k, x_k - y \rangle \leq \\ &\leq G \|x_0 - y\| \sqrt{T} \end{aligned}$$

$$\frac{R_T}{T} \sim \frac{1}{\sqrt{T}}$$

Итого мы имеем для нашей стратегии с постоянным шагом:

$$\boxed{\overline{R_{T-1}} = \frac{1}{T} R_{T-1} \leq G \|x_0 - x^*\| \frac{1}{\sqrt{T}}, \quad \alpha_k = \alpha = \frac{\|x_0 - x^*\|}{G} \sqrt{\frac{1}{T}}}$$

Examples

Least squares with l_1 regularization

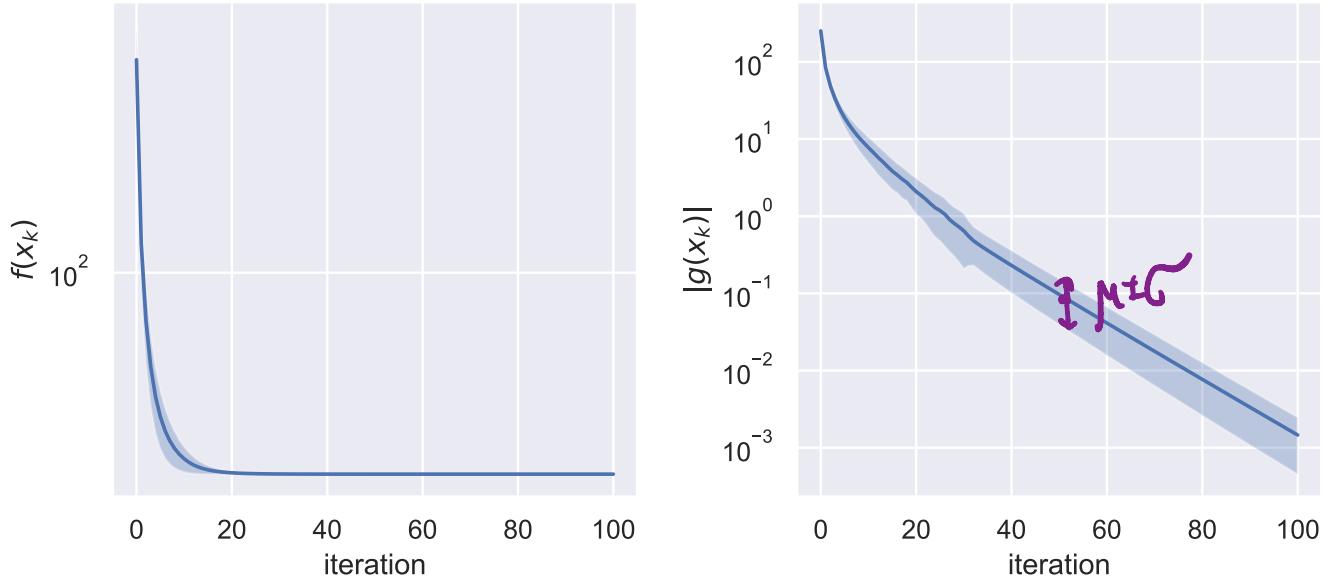
$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$

Algorithm will be written as:

$$x_{k+1} = x_k - \alpha_k (A^\top (Ax_k - b) + \lambda \text{sign}(x_k))$$

where signum function is taken element-wise.

LLS with l_1 regularization. 50 runs. $\lambda = 0.9$



Support vector machines

Let $D = \{(x_i, y_i) \mid x_i \in \mathbb{R}^n, y_i \in \{\pm 1\}\}$

We need to find $\omega \in \mathbb{R}^n$ and $b \in \mathbb{R}$ such that

$$\min_{\omega \in \mathbb{R}^n, b \in \mathbb{R}} \frac{1}{2} \|\omega\|_2^2 + C \sum_{i=1}^m \max[0, 1 - y_i(\omega^\top x_i + b)]$$

Code

- [Open in Colab](#) - Wolfe's example and why we usually have oscillations in non-smooth optimization.
- [Open in Colab](#) - Linear least squares with l_1 - regularization.

References

- [Great cheatsheet](#) by Sebastian Pokutta
- [Lecture](#) on subgradient methods @ Berkley

Немонотонный би-доп шага. Правило Армандо.

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

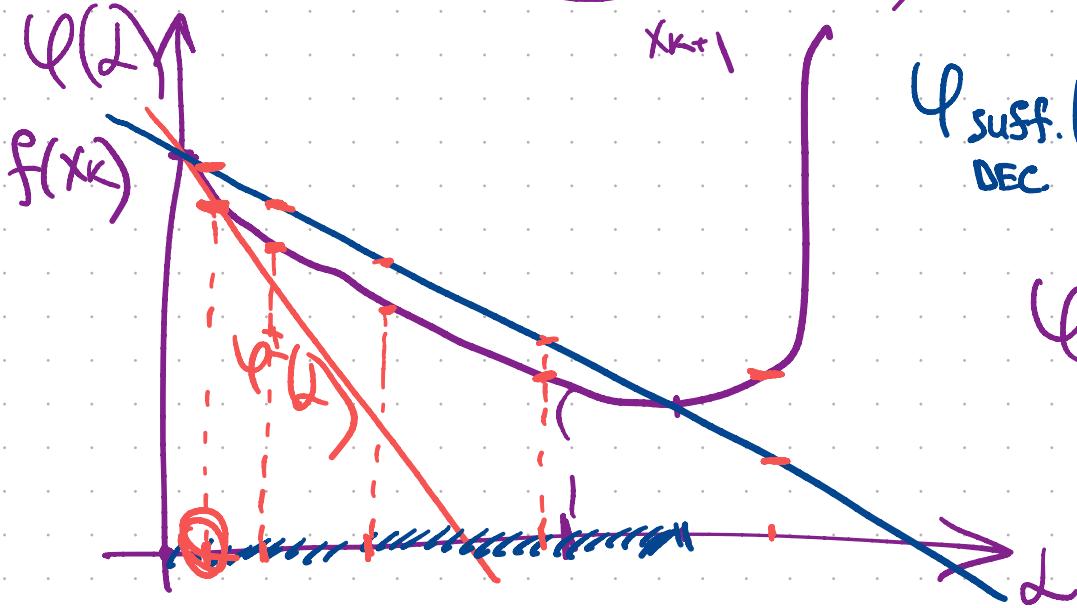
$$x_{k+1} = x_k - \lambda \nabla f(x_k)$$

Нашел - есть:

$$\lambda_k = \operatorname{arg\,min}_{\lambda \in \mathbb{R}^+} f(x_{k+1})$$

$$\varphi(\lambda) = f(x_k - \lambda \nabla f(x_k))$$

$$\varphi'(\lambda) = -\nabla f(x_{k+1})^\top \nabla f(x_k)$$



$$\varphi_{\text{sust.}}(\lambda) = f(x_k) - \frac{C}{2} \|\nabla f(x_k)\|^2 \cdot \lambda$$

гиперпл.
гетатонное
удовлетворение

$$f(x + \Delta x) \approx f(x) + \nabla f(x)^\top \Delta x$$

$$x_{k+1} \approx x_k$$

$$\varphi(\lambda + \delta\lambda) = \varphi(\lambda) + \varphi'(\lambda) \delta\lambda$$

$$\varphi(\lambda) \approx f(x_k) + (-\nabla f(x_{k+1})^\top \nabla f(x_k)) \cdot \lambda$$

$$\varphi(\lambda) \approx f(x_k) - \frac{1}{2} \|\nabla f(x_k)\|^2 \cdot \lambda$$

$$\begin{aligned} \lambda_0 &= 0 \\ \Delta\lambda &= \lambda \\ \lambda &\rightarrow 0 \end{aligned}$$

curvature condition

$$(\ell^1(\lambda)) \geq \dots$$