
Метод проекции градиента. Проксимальные методы.

Семинарист: Данилова М.

Метод проекции градиента

Проекция точки на выпуклое множество

Определение 1. Пусть $a \in \mathbb{R}^n$, $X \subseteq \mathbb{R}^n$. Тогда $\pi_X(a)$ – проекция a на X , если

$$\|\pi_X(a) - a\| \leq \|x - a\| \quad \forall x \in X$$

Замечания

- Если $a \in X$, то $\pi_X(a) = a$.
- Если X – открытое множество и $a \notin X$, то проекции точки a на X не существует.
- Если множество X – выпуклое и замкнутое, то проекция существует и единственна.
- В случае произвольного множества проекция может быть не единственна.

Теорема 1. (Критерий для нормы l_2)

$\pi_X(a)$ – проекция точки a на X тогда и только тогда, когда $\langle \pi_X(a) - a, x - \pi_X(a) \rangle \geq 0 \quad \forall x \in X$

Метод проекции градиента

Будем решать задачу условной минимизации:

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

где $Q \subseteq \mathbb{R}^d$ – выпуклое и замкнутое, а $f(x)$ – дифференцируема.

- Если x_* – точка минимума, то x_* – решение вариационного неравенства:

$$\langle \nabla f(x_*), x - x_* \rangle \geq 0 \quad \forall x \in Q.$$

- Если $f(x)$ – выпукла, то это условие является достаточным.

Метод проекции градиента

Algorithm 1 Метод проекции градиента

Require: размер шага $\gamma > 0$, стартовая точка $x^0 \in \mathbb{R}^d$, количество итераций N
for $k = 0, 1, \dots, N - 1$ **do**
 Вычислить $\nabla f(x^k)$
 $x^{k+1} = \pi_Q(x^k - \gamma \nabla f(x^k))$
end for
Ensure: x^N

Замечания

1. Пусть $f - L$ -гладкая, $x_k \in Q$

$$\begin{aligned}
 x^{k+1} &= \operatorname{argmin}_{x \in Q} \left\{ f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{L}{2} \|x - x^k\|_2^2 \right\} \\
 &= \operatorname{argmin}_{x \in Q} \left\{ \frac{1}{2L} \|\nabla f(x^k)\|_2^2 + \langle \nabla f(x^k), x - x^k \rangle + \frac{L}{2} \|x - x^k\|_2^2 \right\} \\
 &= \operatorname{argmin}_{x \in Q} \left\{ \left\| \frac{1}{\sqrt{2L}} \nabla f(x^k) + \sqrt{\frac{L}{2}} (x - x^k) \right\|_2^2 \right\} \\
 &= \operatorname{argmin}_{x \in Q} \left\{ \frac{L}{2} \left\| \frac{1}{L} \nabla f(x^k) + x - x^k \right\|_2^2 \right\} \\
 &= \operatorname{argmin}_{x \in Q} \left\{ \left\| x - \left(x^k - \frac{1}{L} \nabla f(x^k) \right) \right\|_2^2 \right\} \\
 &= \pi_Q \left(x^k - \frac{1}{L} \nabla f(x^k) \right)
 \end{aligned}$$

Таким образом, для L -гладкой функции на каждом шаге метода проекции градиента с длиной шага $\gamma = \frac{1}{L}$ мы минимизируем квадратичную аппроксимацию на множестве Q (аналогично градиентному методу, где $Q = \mathbb{R}^n$).

2. Метод проекции градиента целесообразно использовать, когда множество Q – **простое множество** (легко искать проекцию точки на данное множество).

Можно формализовать следующим образом: Q – простое множество, если решение следующей задачи

$$\min_{x \in Q} c^\top x$$

находится быстрее (чаще всего аналитически), чем решение исходной.

Пример:

$$\begin{aligned}
 &\min_{\|x-a\|_2 \leq r} c^\top x \\
 x^* &= a - \frac{rc}{\|c\|_2}.
 \end{aligned}$$

Примеры простых множеств:

- $Q = \mathbb{R}_+^n$ – неотрицательный ортант
- $Q = \{x \in \mathbb{R}^n \mid x_i \in [a_i, b_i]\}$ – параллелепипед
- $Q = B_2(0, 1) = \{x \in \mathbb{R}^n \mid \|x\|_2 \leq 1\}$ – шар в 2-норме

Сходимость метода проекции градиента

Метод проекции градиента – релаксационный метод, то есть все точки принадлежат допустимому множеству Q , и функция убывает от итерации к итерации. ($f(x)$ – L -гладкая функция, длина шага $\gamma < \frac{2}{L}$.)

Доказательство. Воспользуемся определением проекции точки $a \in \mathbb{R}^n$ на множество Q :

$$\langle \pi_Q(a) - a, x - \pi_Q(a) \rangle \geq 0 \quad \forall x \in Q$$

Пусть $\pi_Q(a) = x^{k+1}$, $a = x^k - \gamma \nabla f(x^k)$, $x = x^k$:

$$\langle x^{k+1} - (x^k - \gamma \nabla f(x^k)), x^k - x^{k+1} \rangle \geq 0.$$

Обозначим $s_k = x^{k+1} - x^k$:

$$\langle s_k + \gamma \nabla f(x^k), s_k \rangle \leq 0.$$

Разложим функцию $f(x)$ в ряд Тейлора с остаточным членом в форме Лагранжа:

$$f(x^{k+1}) = f(x^k) + \langle \nabla f(\tilde{x}^k, s_k), \tilde{x}^k \in [x^k, x^{k+1}].$$

Воспользуемся L -гладкостью функции $f(x)$ и разложением:

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L}{2} \|x^{k+1} - x^k\|^2 = \\ &= f(x^k) + \langle \nabla f(x^k), s_k \rangle + \frac{L}{2} \|s_k\|^2 = \\ &= f(x^k) + \frac{1}{\gamma} [\langle \gamma \nabla f(x^k), s_k \rangle + \|s_k\|^2] + \left(\frac{L}{2} - \frac{1}{\gamma} \right) \|s_k\|^2 \end{aligned}$$

Следовательно $f(x^{k+1}) - f(x^k) \leq 0$ тогда и только тогда, когда $\frac{L}{2} - \frac{1}{\gamma} < 0 \rightarrow \gamma < \frac{2}{L}$. \square

Теорема 2.

1. Пусть f – **выпуклая**, L -гладкая, $Q \in \mathbb{R}^n$ – выпукло и замкнуто, $\bar{x}^N = \frac{1}{N} \sum_{k=1}^N x^k$.

Тогда

$$f(\bar{x}_N) - f(x^*) \leq \frac{LR^2}{2N},$$

где $R = \|x_0 - x^*\|_2$, x^* – ближайшее к x_0 решение.

Для получения точности $\varepsilon > 0$ необходимо сделать $N = O\left(\frac{LR^2}{\varepsilon}\right)$ итераций, сублинейная скорость сходимости.

2. Пусть $f - \mu$ – **сильно выпуклая**, L -гладкая, $Q \in \mathbb{R}^n$ – выпукло и замкнуто, тогда

$$\|x^N - x^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^N \|x_0 - x^*\|_2^2.$$

Для получения точности $\varepsilon > 0$ необходимо сделать $N = O\left(\frac{L}{\mu} \ln \frac{\|x_0 - x^*\|_2^2}{\varepsilon}\right)$ итераций, линейная скорость сходимости.

Проксимальные методы

Рассмотрим общий вид задачи условной оптимизации:

$$f(x) \rightarrow \min_{x \in Q}, \quad (1)$$

где $Q \subseteq \mathbb{R}^n$ — выпуклое замкнутое множество.

$N(\varepsilon)$	выпуклость	μ -сильная выпуклость
L -гладкость	$\Omega\left(\sqrt{\frac{LR^2}{\varepsilon}}\right)$	$\Omega\left(\sqrt{\frac{L}{\mu}} \ln \frac{\mu R^2}{\varepsilon}\right)$
$\ \nabla f(x)\ _2 \leq M$	$\Omega\left(\frac{M^2 R^2}{\varepsilon^2}\right)$	$\Omega\left(\frac{M^2}{\mu \varepsilon}\right)$

Таблица 1: Нижние оценки на число итераций $N = N(\varepsilon)$, необходимых методу первого порядка для нахождения такой точки x^N , что $f(x^N) - f(x^*) \leq \varepsilon$.

Все приведённые нижние оценки точны в том смысле, что существуют методы оптимизации с верхними оценками на необходимое число итераций, соответствующими этим нижним оценкам. Казалось бы, на этом можно заканчивать изучение оптимизации: оптимальные методы мы изучили, а ничего лучше *в данной постановке задачи* получить нельзя.

Пример

Рассмотрим задачу

$$F(x) = \underbrace{\frac{1}{2}x^\top Ax}_{f(x)} + \underbrace{\lambda\|x\|_1}_{R(x)} \rightarrow \min_{x \in \mathbb{R}^n}, \quad (2)$$

где $A \in \mathbb{S}_+^n$ — симметричная положительно полуопределённая матрица, $\lambda > 0$.

- $f(x) = \frac{1}{2}x^\top Ax$ — выпуклая и L -гладкая функция с $L = \lambda_{\max}(A)$.
- $R(x) = \lambda\|x\|_1$ — выпуклая негладкая функция с ограниченными субградиентами: $\|\nabla R(x)\|_2 \leq \lambda\sqrt{n}$, где $\nabla R(x)$ — произвольный субградиент функции $R(x)$ в точке x .

Единственный класс задач из четырёх классов, рассмотренных ранее, к которому мы можем отнести задачу (2) — это класс выпуклых функций с ограниченными субградиентами. Действительно, достаточно предполагать ограниченность субградиентов только на шаре $B_{2R}(x^*) = \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 \leq 2R\}$, где $R = \|x^0 - x^*\|_2$. Поэтому можно ограничить $\nabla f(x)$ по норме некоторой константой на этом шаре. Пусть $\|\nabla F(x)\|_2 \leq M$, тогда градиентный спуск с правильно выбранным размером шага (порядка $\frac{\varepsilon}{M^2}$) будет сходиться для данной задачи со скоростью $O\left(\frac{M^2 R^2}{\varepsilon^2}\right)$.

НО полученная оценка не учитывает структуру задачи: мы полностью проигнорировали тот факт, что $f(x)$ имеет Липшицев градиент и что $R(x)$ достаточно «простая» функция. Оказывается для такого вида задач можно немного видоизменить градиентный спуск и получить

метод, который будет сходиться со скоростью $O(\frac{LR^2}{\varepsilon})$. Более того, можно получить ускоренный метод, который будет работать ещё быстрее — со скоростью $O\left(\sqrt{\frac{LR^2}{\varepsilon}}\right)$. Но для начала нам нужно формально определить, с каким новым классом задач мы имеем дело.

Задачи выпуклой композитной оптимизации

(Задачи с регуляризацией)

Рассмотрим задачу

$$F(x) = f(x) + R(x) \rightarrow \min_{x \in \mathbb{R}^n}, \quad (3)$$

где

- $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ — L -гладкая, выпуклая функция.
- $R(x) : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ — правильная выпуклая замкнутая функция. Здесь
 - правильная функция = функция, которая не всюду равна $+\infty$,
 - замкнутая функция = функция, у которой множества уровня замкнуты, т.е. для всех $\alpha \in \mathbb{R}$ множество $\{x \in \mathbb{R}^n \mid R(x) \leq \alpha\}$ замкнуто.

Пример 2.

1. $R(x) = 0 \Rightarrow \min_{x \in \mathbb{R}^n} f(x)$
2. $R(x) = \delta_Q(x) = \begin{cases} 0, & x \in Q \\ +\infty, & x \notin Q \end{cases} \Rightarrow \min_{x \in Q} f(x)$
3. $R(x) = \lambda \|x\|_1$

Проксимальный оператор и его свойства

Для функции $R(x)$, удовлетворяющей условиям (3), рассмотрим отображение из \mathbb{R}^n в \mathbb{R}^n :

$$\text{prox}_R(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{y \in \mathbb{R}^n} \left\{ R(y) + \frac{1}{2} \|y - x\|_2^2 \right\}. \quad (4)$$

Далее мы будем рассматривать только такие функции $R(x)$, относительно которых можно «быстро» вычислять проксимальный оператор.

1. $\text{prox}_R(x)$ определяется однозначно для любого $x \in \mathbb{R}^n$, т.е. $\text{prox}_R(x)$ — это отображение. Действительно, задача (4) является сильно выпуклой, а значит, имеет единственное решение.

2. $u = \text{prox}_R(x) \stackrel{\textcircled{1}}{\iff} x - u \in \partial R(u) \stackrel{\textcircled{2}}{\iff} \langle x - u, y - u \rangle \leq R(y) - R(u) \quad \forall y \in \mathbb{R}^n$. Действительно, $\textcircled{1}$ следует из необходимого и достаточного условия оптимальности первого порядка

$$u = \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ R(y) + \frac{1}{2} \|y - x\|_2^2 \right\} \iff 0 \in u - x + \partial R(u),$$

а $\textcircled{2}$ следует из определения субградиента R в точке u .

3. Для всех $x, y \in \mathbb{R}^n$ выполняются неравенства:

$$\langle x - y, \text{prox}_R(x) - \text{prox}_R(y) \rangle \geq \| \text{prox}_R(x) - \text{prox}_R(y) \|_2^2, \quad (5)$$

$$\| \text{prox}_R(x) - \text{prox}_R(y) \|_2 \leq \|x - y\|_2. \quad (6)$$

Пусть $u = \text{prox}_R(x)$, $v = \text{prox}_R(y)$, тогда по свойству 2 имеем: $\langle x - u, v - u \rangle \leq R(v) - R(u)$ и $\langle y - v, u - v \rangle \leq R(u) - R(v)$. Складывая эти неравенства, получаем: $\langle u - v, y - x + u - v \rangle \leq 0$, что эквивалентно (5). Теперь покажем (6). Если $u = v$, то неравенство (6) очевидно. Если же $u \neq v$, то из (5) и неравенства Коши-Буняковского-Шварца получаем: $\|u - v\|_2^2 \leq \langle x - y, u - v \rangle \leq \|x - y\|_2 \cdot \|u - v\|_2$, что после сокращения левой и правой частей на $\|u - v\|_2$ даёт (6).

Примеры вычисления проксимальных операторов

1. $R(x) = c$, где $c \in \mathbb{R}$. Тогда

$$\text{prox}_R(x) = \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ c + \frac{1}{2} \|y - x\|_2^2 \right\} = x.$$

2. $R(x) = \delta_Q(x) = \begin{cases} 0, & x \in Q, \\ +\infty, & x \notin Q, \end{cases}$ где $Q \subseteq \mathbb{R}^n$ — выпуклое замкнутое непустое множество.

Заметим, что минимум в определении проксимального оператора не может достигаться вне множества Q , т.к. функция под минимумом вне этого множества равняется $+\infty$. Поэтому

$$\begin{aligned} \text{prox}_R(x) &= \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ \delta_Q(y) + \frac{1}{2} \|y - x\|_2^2 \right\} = \underset{y \in Q}{\operatorname{argmin}} \left\{ \frac{1}{2} \|y - x\|_2^2 \right\} \\ &= \pi_Q(x) \quad - \text{метод проекции градиента.} \end{aligned}$$

3. $R(x) = \frac{1}{2} x^\top A x + b^\top x + c$, где $A \in \mathbb{S}_+^n$, $b \in \mathbb{R}^n$ и $c \in \mathbb{R}$. Пусть $u = \text{prox}_R(x)$. Тогда

$$u = \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ \frac{1}{2} y^\top A y + b^\top y + c + \frac{1}{2} \|y - x\|_2^2 \right\},$$

что по свойству 2 эквивалентно тому, что

$$x - u = Au + b \iff u = (A + I)^{-1}(x - b) = \text{prox}_R(x).$$

4. $R(x) = \lambda \|x\|_1$, где $\lambda > 0$. Чтобы найти прокс-оператор от данной функции, докажем вспомогательное утверждение.

Утверждение 1 (Прокс-оператор сепарабельной функции). Пусть $R(x) = R(x_1, \dots, x_r) = \sum_{i=1}^r R_i(x_i)$, где $x = (x_1^\top, \dots, x_r^\top)^\top \in \mathbb{R}^n$ и $x_i \in \mathbb{R}^{n_i}$ для $i = 1, \dots, r$. Тогда

$$\text{prox}_R(x) = \begin{pmatrix} \text{prox}_{R_1}(x_1) \\ \vdots \\ \text{prox}_{R_r}(x_r) \end{pmatrix}.$$

Доказательство. По определению мы имеем

$$\begin{aligned} \text{prox}_R(x) &= \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ \sum_{i=1}^r R_i(y_i) + \frac{1}{2} \|y - x\|_2^2 \right\} \\ &= \underset{y_i \in \mathbb{R}^{n_i}, i=1, \dots, r}{\operatorname{argmin}} \left\{ \sum_{i=1}^r \left(R_i(y_i) + \frac{1}{2} \|y_i - x_i\|_2^2 \right) \right\}. \end{aligned}$$

Отсюда следует, что задача распадается на r независимых подзадач. Используя определение $\text{prox}_{R_i}(x_i)$, получаем требуемое. \square

Возвращаясь к исходной задаче, замечаем, что $R(x) = \sum_{i=1}^n \lambda |x_i|$, то есть достаточно найти прокс-оператор функции одного аргумента $g(x) = \lambda |x|$, $x \in \mathbb{R}$:

$$u = \text{prox}_g(x) = \underset{y \in \mathbb{R}}{\operatorname{argmin}} \left\{ \lambda |y| + \frac{1}{2} (y - x)^2 \right\}.$$

- Минимум достигается при $y > 0$, тогда и только тогда, когда $\lambda + u - x = 0 \iff u = x - \lambda$. Это означает, что если $x > \lambda$, то $\text{prox}_g(x) = x - \lambda$.
- Аналогичными рассуждениями получаем, что если $x < -\lambda$, то $\text{prox}_g(x) = x + \lambda$.
- Во всех остальных случаях, т.е. при $x \in [-\lambda, \lambda]$, $\text{prox}_g(x) = 0$.

Полученный результат можно записать в следующем виде:

$$\text{prox}_g(x) = [|x| - \lambda]_+ \cdot \text{sign}(x), \text{ где } \text{sign}(x) \stackrel{\text{def}}{=} \begin{cases} 1, & x > 0, \\ 0, & x = 0, \\ -1, & x < 0 \end{cases}$$

и $[y]_+ \stackrel{\text{def}}{=} \max\{y, 0\}$. Отсюда следует, что для $R(x) = \lambda \|x\|_1$

$$\text{prox}_R(x) = [|x| - \lambda \mathbf{1}]_+ \odot \text{sign}(x),$$

где $\mathbf{1} \stackrel{\text{def}}{=} (1, \dots, 1)^\top \in \mathbb{R}^n$, модуль $|x|$, срезка $[|x| - \lambda \mathbf{1}]_+$ и сигнум (знак) $\text{sign}(x)$ применяются к векторам покомпонентно и $y \odot z \stackrel{\text{def}}{=} (y_1 z_1, \dots, y_n z_n)^\top$ обозначает произведение Адамара двух векторов (покомпонентное произведение).

Проксимальный градиентный спуск

Для задачи (3) рассмотрим следующий метод.

Algorithm 2 Проксимальный градиентный спуск

Require: размер шага $\gamma > 0$, стартовая точка $x^0 \in \mathbb{R}^d$, количество итераций N

```

1: for  $k = 0, 1, \dots, N - 1$  do
2:   Вычислить  $\nabla f(x^k)$ 
3:    $x^{k+1} = \text{prox}_{\gamma R}(x^k - \gamma \nabla f(x^k))$ 
4: end for

```

Ensure: x^N

- Внешне метод очень похож на градиентный спуск: по-прежнему требуется вычислять градиентный шаг, но теперь дополнительно от получаемой градиентным шагом точки вычисляется прокс.
- Пусть $x^* = \operatorname{argmin}_{x \in \mathbb{R}^n} F(x)$. Тогда $x^* = \operatorname{prox}_{\gamma R}(x^* - \gamma \nabla f(x^*))$.
- Пусть $f(x)$ — L -гладкая, тогда шаг проксимального метода можно выразить следующим образом:

$$\begin{aligned}
 x_{k+1} &= \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|_2^2 + R(x) \right\} \\
 &= \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ \frac{L}{2} \left\| x - x_k + \frac{1}{L} \nabla f(x_k) \right\|_2^2 + R(x) \right\} \\
 &= \operatorname{prox}_{\frac{1}{L} R(x)} \left(x_k - \frac{1}{L} \nabla f(x_k) \right)
 \end{aligned}$$

Сильно выпуклый случай

Теорема 3. Пусть $f(x)$ — μ -сильно выпуклая L -гладкая функция, $R(x)$ — правильная выпуклая замкнутая функция и $\gamma \leq \frac{1}{L}$. Тогда для любого $N > 0$ выход Алгоритма 2 удовлетворяет неравенству:

$$\|x^N - x^*\|_2^2 \leq (1 - \gamma\mu)^N \|x^0 - x^*\|_2^2. \quad (7)$$

Иными словами, для проксимального градиентного спуска с $\gamma = \frac{1}{L}$ через $N = O\left(\frac{L}{\mu} \ln \frac{R^2}{\varepsilon}\right)$ итераций, где $R = \|x^0 - x^*\|_2$, выполняется $\|x^N - x^*\|_2^2 \leq \varepsilon$.

Доказательство. Пользуясь тем, что прокс-оператор является нестягивающим (см. (6)), мы получаем

$$\begin{aligned}
 \|x^{k+1} - x^*\|_2^2 &= \|\operatorname{prox}_{\gamma R}(x^k - \gamma \nabla f(x^k)) - \operatorname{prox}_{\gamma R}(x^* - \gamma \nabla f(x^*))\|_2^2 \\
 &\stackrel{(6)}{\leq} \|x^k - x^* - \gamma (\nabla f(x^k) - \nabla f(x^*))\|_2^2 \\
 &= \|x^k - x^*\|_2^2 - 2\gamma \langle x^k - x^*, \nabla f(x^k) - \nabla f(x^*) \rangle \\
 &\quad + \gamma^2 \|\nabla f(x^k) - \nabla f(x^*)\|_2^2.
 \end{aligned}$$

Из сильной выпуклости функции f имеем

$$f(x^*) \geq f(x^k) + \langle \nabla f(x^k), x^* - x^k \rangle + \frac{\mu}{2} \|x^* - x^k\|_2^2,$$

откуда следует, что

$$-\langle x^k - x^*, \nabla f(x^k) - \nabla f(x^*) \rangle \leq -\frac{\mu}{2} \|x^k - x^*\|_2^2 - \underbrace{(f(x^k) - f(x^*) - \langle \nabla f(x^*), x^k - x^* \rangle)}_{V_f(x^k, x^*)}.$$

Из полученных неравенств имеем:

$$\|x^{k+1} - x^*\|_2^2 \leq (1 - \gamma\mu) \|x^k - x^*\|_2^2 - 2\gamma V_f(x^k, x^*) + \gamma^2 \|\nabla f(x^k) - \nabla f(x^*)\|_2^2.$$

Кроме того, из L -гладкости и выпуклости функции f следует (см. Теорему 2.1.5 из книги Ю.Е. Нестерова, 2010 года), что для любых $x, y \in \mathbb{R}^n$

$$\|\nabla f(x) - \nabla f(y)\|_2^2 \leq 2L(f(x) - f(y) - \langle \nabla f(y), x - y \rangle) = 2LV_f(x, y).$$

Используя это неравенство с $x = x^k$ и $y = x^*$, мы продолжаем наши цепочки неравенств для $\|x^{k+1} - x^*\|_2^2$:

$$\|x^{k+1} - x^*\|_2^2 \leq (1 - \gamma\mu) \|x^k - x^*\|_2^2 - 2\gamma(1 - \gamma L) V_f(x^k, x^*)$$

Из выпуклости f имеем:

$$V_f(x^k, x^*) = f(x^k) - f(x^*) - \langle \nabla f(x^*), x^k - x^* \rangle \geq 0.$$

Кроме того, т.к. $\gamma > 0$ и $\gamma \leq \frac{1}{L}$, то $2\gamma(1 - \gamma L) V_f(x^k, x^*) \geq 0$, откуда следует, что

$$\|x^{k+1} - x^*\|_2^2 \leq (1 - \gamma\mu) \|x^k - x^*\|_2^2.$$

Поскольку формула выше выполнена для всех целых $k \geq 0$, то отсюда следует, что

$$\|x^{k+1} - x^*\|_2^2 \leq (1 - \gamma\mu)^{k+1} \|x^0 - x^*\|_2^2.$$

Наконец, подставляя $\gamma = \frac{1}{L}$ и пользуясь неравенством $(1 - t)^k \leq e^{-tk}$, получаем, что для достижения $\|x^N - x^*\|_2^2 \leq \varepsilon$ достаточно $N = \frac{L}{\mu} \ln \frac{\|x^0 - x^*\|_2^2}{\varepsilon}$ итераций данного метода. \square

Выпуклый случай

Теорема 4. Пусть $f(x)$ — выпуклая L -гладкая функция, $R(x)$ — правильная выпуклая замкнутая функция и $\gamma = \frac{1}{L}$. Тогда для любого $N > 0$ выход Алгоритма 2 удовлетворяет неравенству:

$$F(x^N) - F(x^*) \leq \frac{L\|x^0 - x^*\|_2^2}{2N}. \quad (8)$$

Иными словами, для проксимального градиентного спуска с $\gamma = \frac{1}{L}$ через $N = O\left(\frac{LR^2}{\varepsilon}\right)$ итераций, где $R = \|x^0 - x^*\|_2$, выполняется $F(x^N) - F(x^*) \leq \varepsilon$.

Проксимальный ускоренный градиентный метод (FISTA)

FISTA: выпуклый случай

Пусть функция f выпукла. Тогда проксимальный градиентный метод можно ускорить следующим способом.

Algorithm 3 Проксимальный ускоренный градиентный метод (FISTA) для выпуклых функций

Require: стартовая точка $x^0 \in \mathbb{R}^d$, количество итераций N

```

1:  $y^0 = x^0, t_0 = 1$ 
2: for  $k = 0, 1, \dots, N - 1$  do
3:   Вычислить  $\nabla f(y^k)$ 
4:    $x^{k+1} = \text{prox}_{\frac{1}{L}R}(y^k - \frac{1}{L}\nabla f(y^k))$ 
5:    $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$ 
6:    $y^{k+1} = x^{k+1} + \frac{t_k - 1}{t_{k+1}}(x^{k+1} - x^k)$ 
7: end for

```

Ensure: x^N

Теорема 5. Пусть $f(x)$ — выпуклая L -гладкая функция, $R(x)$ — правильная выпуклая замкнутая функция. Тогда для любого $N > 0$ выход Алгоритма 3 удовлетворяет неравенству:

$$F(x^N) - F(x^*) \leq \frac{2L\|x^0 - x^*\|_2^2}{(N + 1)^2}. \quad (9)$$

Иными словами, для FISTA через $N = O\left(\sqrt{\frac{LR^2}{\varepsilon}}\right)$ итераций, где $R = \|x^0 - x^*\|_2$, выполняется $F(x^N) - F(x^*) \leq \varepsilon$.

FISTA: сильно выпуклый случай

Пусть теперь функция f μ -сильно выпукла.

Algorithm 4 Проксимальный ускоренный градиентный метод (FISTA) для сильно выпуклых функций

Require: стартовая точка $x^0 \in \mathbb{R}^d$, количество итераций N

```

1:  $y^0 = x^0, \varkappa = \frac{L}{\mu}$ 
2: for  $k = 0, 1, \dots, N - 1$  do
3:   Вычислить  $\nabla f(y^k)$ 
4:    $x^{k+1} = \text{prox}_{\frac{1}{L}R}(y^k - \frac{1}{L}\nabla f(y^k))$ 
5:    $y^{k+1} = x^{k+1} + \frac{\sqrt{\varkappa} - 1}{\sqrt{\varkappa} + 1}(x^{k+1} - x^k)$ 
6: end for

```

Ensure: x^N

Теорема 6. Пусть $f(x)$ — μ -сильно выпуклая L -гладкая функция, $R(x)$ — правильная выпуклая замкнутая функция. Тогда для любого $N > 0$ выход Алгоритма 4 удовлетворяет неравенству:

$$F(x^N) - F(x^*) \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right)^N \left(F(x^0) - F(x^*) + \frac{\mu}{2}\|x^0 - x^*\|_2^2\right). \quad (10)$$

Иными словами, для FISTA через $N = O\left(\sqrt{\frac{L}{\mu}} \ln \frac{F(x^0) - F(x^*) + \mu R^2}{\varepsilon}\right)$ итераций, где $R = \|x^0 - x^*\|_2$, выполняется $F(x^N) - F(x^*) \leq \varepsilon$.