

Методы оптимизации, ФКН ВШЭ, зима 2017

Практическое задание 3: Выпуклая негладкая, условная и структурная оптимизация.

Срок сдачи: 26 марта 2017 (23:59).
Язык программирования: Python 3.

1 Алгоритмы

1.1 Субградиентный метод

Субградиентный метод — это метод решения безусловной негладкой задачи оптимизации

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

где $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ — выпуклая функция (не обязательно дифференцируемая).

Итерация субградиентного метода заключается в шаге из текущей точки x_k в направлении (любого) анти-субградиента $\phi'(x_k)$. При этом, поскольку для негладких задач норма субградиента $\|\phi'(x_k)\|_2$ не является информативной, субградиентный метод использует в качестве направления нормированный вектор $\phi'(x_k)/\|\phi'(x_k)\|_2$:

$$x_{k+1} = x_k - \alpha_k \frac{\phi'(x_k)}{\|\phi'(x_k)\|_2}.$$

Для сходимости метода необходимо, чтобы длины шагов α_k убывали к нулю, но не слишком быстро:

$$\alpha_k > 0, \quad \alpha_k \rightarrow 0, \quad \sum_{k=0}^{\infty} \alpha_k = \infty.$$

Обычно длины шагов выбирают по правилу $\alpha_k = \alpha/\sqrt{k+1}$, где $\alpha > 0$ — некоторая константа.

Нужно отметить, что последовательность $(x_k)_{k=0}^{\infty}$, построенная субградиентным методом, может не быть релаксационной последовательностью для функции ϕ , т. е. неравенство $\phi(x_{k+1}) \leq \phi(x_k)$ может не выполняться. Поэтому в субградиентном методе в качестве результата работы после N итераций метода вместо точки x_N возвращается точка $y_N := \operatorname{argmin}\{\phi(x) : x \in \{x_0, \dots, x_N\}\}$ (т. е. из всех пробных точек x_k , построенных методом, выбирается та, в которой значение функции оказалось наименьшим).¹

1.2 Проксимальный градиентный метод

Проксимальный градиентный метод используется для минимизации *композиционных функций*:

$$\phi(x) := f(x) + h(x),$$

где функция $f : \mathbb{R}^n \rightarrow \mathbb{R}$ непрерывно дифференцируемая, а функция $h : \mathbb{R}^n \rightarrow \mathbb{R}$ выпуклая (не обязательно дифференцируемая) и достаточно простая.

Под *простотой* функции h подразумевается то, что для этой функции возможно эффективно вычислить проксимальное отображение

$$\operatorname{Prox}_h^\alpha(x) := \operatorname{argmin}_{y \in \mathbb{R}^n} \left[\alpha h(y) + \frac{1}{2} \|y - x\|_2^2 \right],$$

¹Заметим, что в практической реализации метода для вычисления результата y_N сами точки x_0, \dots, x_N хранить в памяти не нужно.

где $x \in \mathbb{R}^n$, $\alpha > 0$. Например, для $h(x) = \|x\|_1$ проксимальное отображение может быть вычислено аналитически независимо для каждой из координат $1 \leq i \leq n$:

$$[\text{Prox}_{\|\cdot\|_1}^\alpha(x)]_i = \begin{cases} x_i + \alpha, & x_i < -\alpha, \\ 0, & |x_i| \leq \alpha, \\ x_i - \alpha, & x_i > \alpha. \end{cases}$$

Итерация проксимального градиентного метода имеет следующий вид:

$$x_{k+1} = \text{Prox}_h^{\alpha_k}(x_k - \alpha_k \nabla f(x_k)),$$

где $\nabla f(x_k)$ — градиент функции f в точке x_k , а $\alpha_k > 0$ — (должным образом выбранная) длина шага.

1.2.1 Схема Нестерова для подбора длины шага

Итерация проксимального градиентного метода имеет следующий геометрический смысл. Если градиент функции f удовлетворяет условию Липшица с константой $L_f > 0$, то для любых $x, y \in \mathbb{R}^n$ и $L \geq L_f$ справедлива оценка

$$\phi(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2}\|y - x\|_2^2 + h(y) =: m_L(y; x), \quad (1)$$

которая является точной в точке $y = x$. Выбрав $x = x_k$, $L = L_k := 1/\alpha_k$, получаем, что итерация проксимального градиентного метода в точности соответствует минимизации функции $m_{L_k}(\cdot; x_k)$:

$$x_{k+1} = \underset{y \in \mathbb{R}^n}{\text{argmin}} m_{L_k}(y; x_k).$$

Если $L_k \geq L_f$, то, согласно оценке (1), значение функции ϕ в новой точке x_{k+1} уменьшается относительно значения в старой точке $\phi(x_k)$ по крайней мере на величину $m_{L_k}(x_{k+1}; x_k) - \phi(x_k)$:

$$\phi(x_{k+1}) \leq m_{L_k}(x_{k+1}; x_k) \quad (\leq \phi(x_k)). \quad (2)$$

Оказывается, что именно это неравенство, а не условие $L_k \geq L_f$ отвечает за сходимость проксимального градиентного метода. Поэтому для подбора константы L_k (или, эквивалентно, длины шага α_k) на текущей итерации k можно использовать следующую простую процедуру линейного поиска: начать с некоторого значения L и увеличивать его в два раза, пока не выполнится неравенство (2). Заметим, что эта процедура определена корректно: как только значение L превысит L_f , неравенство (2) обязательно будет выполнено в силу липшицевости градиента функции f . О константе L_k удобно думать как о «локальной» константе Липшица градиента функции f .

Поскольку каждое пробное значение константы L требует полного вычисления функции ϕ , то число итераций линейного поиска необходимо, по возможности, сократить. Для этого имеет смысл инициализировать значение L_{k+1} с помощью уже найденного значения L_k . Естественным вариантом является инициализация $L_{k+1} = L_k$. Однако в этом случае константы L_k всегда будут только увеличиваться и никогда не уменьшаться; это плохо, поскольку при прочих равных лучше выбрать значение L как можно меньше, что будет соответствовать большему шагу и большему прогрессу в оптимизации (в некоторых областях пространства локальная константа Липшица может быть меньше, чем в области, содержащей начальную точку метода). Таким образом, чтобы константы Липшица на соседних итерациях были близки, но также со временем могли уменьшаться, в схеме Нестерова выполняется инициализация $L_{k+1} = L_k/2$.

Итоговый алгоритм проксимального градиентного метода с подбором длины шага по схеме Нестерова представлен в алгоритме 1. Здесь $L_0 > 0$ — параметр метода (нижняя оценка на глобальную константу Липшица L_f); его всегда можно выбрать равным единице ($L_0 = 1$) без принципиального ущерба для сходимости метода.

Можно показать, что, несмотря на то, что на отдельных итерациях проксимального градиентного метода может выполняться много шагов линейного поиска, среднее число вычислений функции ϕ за итерацию равно двум (задание 3.2).

Алгоритм 1 Проксимальный градиентный метод с подбором длины шага по схеме Нестерова

Вход: Начальная точка $x_0 \in \mathbb{R}^n$, начальная оценка константы Липшица $L_0 > 0$.

```
1:  $L \leftarrow L_0$ 
2: for  $k = 0, 1, \dots$  do
3:   while True do
4:      $y \leftarrow \operatorname{argmin}_y m_L(y; x_k)$ 
5:     if  $\phi(y) \leq m_L(y; x_k)$  then
6:       break
7:     end if
8:      $L \leftarrow 2L$ 
9:   end while
10:   $x_{k+1} \leftarrow y$ 
11:   $L \leftarrow \max\{L_0, L/2\}$ 
12: end for
```

Выход: Последняя вычисленная точка x_k

1.3 Метод барьеров

2 Модели

2.1 Модель LASSO

Модель LASSO является одной из стандартных моделей линейной регрессии. Имеется обучающая выборка $((a_i, b_i))_{i=1}^n$, где $a_i \in \mathbb{R}^n$ — вектор признаков i -го объекта, а $b_i \in \mathbb{R}$ — его регрессионное значение. Задача заключается в прогнозировании регрессионного значения b_{new} для нового объекта, представленного своим вектором признаков a_{new} .

В модели LASSO, как и в любой модели линейной регрессии, прогнозирование выполняется с помощью линейной комбинации компонент вектора a с некоторыми фиксированными коэффициентами $x \in \mathbb{R}^n$:

$$b(a) := a^T x.$$

Коэффициенты x являются параметрами модели и настраиваются с помощью решения следующей оптимизационной задачи:

$$\phi(x) := \frac{1}{2} \sum_{i=1}^n (a_i^T x - b_i)^2 + \lambda \sum_{j=1}^n |x_j| =: \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \rightarrow \min_{x \in \mathbb{R}^n}. \quad (3)$$

Здесь $\lambda > 0$ — коэффициент регуляризации (параметр модели). Особенностью LASSO является использование именно ℓ_1 -регуляризации (а не, например, ℓ_2 -регуляризации). Такая регуляризация позволяет получить разреженное решение. В разреженном решении x^* часть компонент равна нулю. (Можно показать, что при $\lambda \geq \|A^T b\|_\infty$ все компоненты будут нулевыми). Нулевые веса соответствуют исключению соответствующих признаков из модели (признание их неинформативными).

2.2 Двойственная задача и критерий останковки

Перепишем задачу (3) в следующем эквивалентном виде:

$$\min_{x \in \mathbb{R}^n, z \in \mathbb{R}^n} \left\{ \frac{1}{2} \|z\|_2^2 + \lambda \|x\|_1 : Ax - b = z \right\}. \quad (4)$$

Можно показать, что двойственной задачей к задаче (4) является

$$\max_{\mu \in \mathbb{R}^n} \left\{ -\frac{1}{2} \|\mu\|_2^2 - b^T \mu : \|A^T \mu\|_\infty \leq \lambda \right\}. \quad (5)$$

Таким образом, имея в распоряжении допустимую двойственную точку $\mu \in \mathbb{R}^n$, т. е. такую что $\|A^T \mu\|_\infty \leq \lambda$, можно вычислить следующую оценку для невязки в задаче (3):

$$\phi(x) - \phi^* \leq \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 + \frac{1}{2} \|\mu\|_2^2 + b^T \mu =: \eta(x, \mu). \quad (6)$$

Величина $\eta(x, \mu)$ называется *зазором двойственности* и обращается в ноль в оптимальных решениях x^* и μ^* задач (3) и (5). Заметим, что решения x^* и μ^* связаны между собой следующим соотношением: $Ax^* - b = \mu^*$. Поэтому для фиксированного x естественным выбором соответствующего μ будет

$$\mu(x) := \min \left\{ 1, \frac{\lambda}{\|A^T(Ax - b)\|_\infty} \right\} (Ax - b).$$

Такой выбор обеспечивает стремление зазора двойственности $\eta(x, \mu(x))$ к нулю при $x \rightarrow x^*$, что позволяет использовать условие $\eta(x, \mu(x)) < \varepsilon$ в качестве критерия останова в любом итерационном методе решения задачи (3).

В этом задании все рассматриваемые методы должны использовать критерий останова по зазору двойственности $\eta(x, \mu(x))$.

2.3 Сведение к гладкой условной задаче оптимизации

Выполнив эквивалентное преобразование задачи через надграфик, задачу (3) можно переписать в виде гладкой условной задачи:

$$\min_{x, u} \left\{ \frac{1}{2} \|Ax - b\|_2^2 + 1_n^T u \right\} \quad \text{s. t.} \quad x \preceq u, \quad x \succeq -u. \quad (7)$$

где $x, u \in \mathbb{R}^n$ и $1_n := (1, \dots, 1) \in \mathbb{R}^n$. Эта задача является задачей QP, и ее можно решать, например, с помощью метода барьеров.

3 Формулировка задания

- 1 Скачайте коды, прилагаемые к заданию:

<https://github.com/mihaha/hse-optimization-course/tree/master/task3>

Эти файлы содержат прототипы функций, которые Вам нужно будет реализовать. Некоторые процедуры уже частично или полностью реализованы.

- 2 Реализуйте субградиентный метод (функция `subgradient_method` в модуле `optimization`).
- 3 Реализуйте негладкий оракул для функции (3) (классы `LeastSquaresOracle` и `LassoNonsmoothOracle` в модуле `oracles`).
- 4 Реализуйте процедуру вычисления зазора двойственности (6) (функция `lasso_duality_gap` в модуле `oracles`).
- 5 Реализуйте проксимальный градиентный метод (функция `proximal_gradient_descent` в модуле `optimization`).
- 6 Реализуйте композитный оракул для функции (3) (классы `L1RegOracle` и `LassoProxOracle` в модуле `oracles`).
- 7 Выпишите для задачи (7) вспомогательную функцию $\phi_t(x, u)$, минимизируемую в методе барьеров. Выпишите систему линейных уравнений, задающую ньютоновское направление $d_k = (d_k^x, d_k^u)$. Предложите свой способ решения этой системы и прокомментируйте его достоинства и недостатки. Для текущей точки (x, u) и направления d_k определите максимальную допустимую длину шага α . Какую начальную точку (x_0, u_0) можно выбрать для метода барьеров?

- 8 Реализуйте метод барьеров для задачи (7) (функция `barrier_method_lasso` в модуле `optimization`). Для подбора длины шага на внутренних итерациях используйте бэктрекинг. (Не забудьте выбрать начальное значение длины шага равным единице, если оно допустимо.)
- 9 Проведите эксперименты, описанные ниже. Напишите отчет.

3.1 Эксперимент: Выбор длины шага в субградиентном методе

Исследуйте работу субградиентного метода в зависимости от выбора константы α_0 в формуле для длины шага. Обязательно попробуйте различные начальные точки.

3.2 Эксперимент: Среднее число итераций линейного поиска в схеме Нестерова

Для проксимального метода построить график суммарного (кумулятивного) числа итераций линейного поиска в зависимости от номера итерации. Действительно ли среднее число итераций линейного поиска не превышает (примерно) двух?

3.3 Эксперимент: Сравнение методов

Провести экспериментальное сравнение трех реализованных методов для различных значений числа переменных n ; рассмотреть, как минимум, случаи $n \leq 100$ и $n \geq 1000$. Какие методы быстрее достигают низкой точности ($\varepsilon = 10^{-2}$), а какие высокой ($\varepsilon = 10^{-10}$)?

Значение коэффициента регуляризации λ выбрать стандартным образом: $\lambda = 1$.

Для сравнения методов рисовать графики 1) числа итераций против гарантируемой точности по зазору двойственности (6) и 2) реального времени работы против зазора двойственности. Для зазора двойственности использовать логарифмическую шкалу.

4 Оформление задания

Результатом выполнения задания являются

1. Файлы `optimization.py` и `oracles.py` с реализованными методами и оракулами.
2. Полные исходные коды для проведения экспериментов и рисования всех графиков. Все результаты должны быть воспроизводимыми. Если вы используете случайность — зафиксируйте `seed`.
3. Отчет в формате PDF о проведенных исследованиях.

Выполненное задание следует отправить письмом на почту своей группы² с заголовком

Практическое задание 3, Фамилия Имя.

Задание следует присылать только один раз с окончательным вариантом.

Каждый проведенный эксперимент следует оформить как отдельный раздел в PDF документе (название раздела — название соответствующего эксперимента). Для каждого эксперимента необходимо сначала написать его описание: какие функции оптимизируются, каким образом генерируются данные, какие методы и с какими параметрами используются. Далее должны быть представлены результаты соответствующего эксперимента — графики, таблицы и т.д. Наконец, после результатов эксперимента должны быть написаны Ваши выводы — какая зависимость наблюдается и почему.

²Для 141 группы: `opt.homework+141@gmail.com`. Для 142 группы: `opt.homework+142@gmail.com`. Для 145 группы: `opt.homework+145@gmail.com`.

Важно: Отчет не должен содержать никакого кода. Каждый график должен быть прокомментирован — что на нем изображено, какие выводы можно сделать из этого эксперимента. Обязательно должны быть подписаны оси. Если на графике нарисовано несколько кривых, то должна быть легенда. Сами линии следует рисовать достаточно толстыми, чтобы они были хорошо видимыми.

5 Проверка задания

Перед отправкой задания обязательно убедитесь, что Ваша реализация проходит автоматические *предварительные* тесты `presubmit_tests.py`, выданные вместе с заданием. Для этого запустите команду

```
python presubmit_tests.py
```

Важно: Файл с тестами может измениться. Перед отправкой обязательно убедитесь, что ваша версия `presubmit_tests.py` — последняя.

Важно: Решения, которые не будут проходить тесты `presubmit_tests.py`, будут автоматически оценены в 0 баллов. Проверяющий не будет разбираться, почему Ваш код не работает и читать Ваш отчет.

Оценка за задание будет складываться из двух частей:

1. Правильность и эффективность реализованного кода.
2. Качество отчета.

Правильность и эффективность реализованного кода будет оцениваться автоматически с помощью независимых тестов (отличных от предварительных тестов). Качество отчета будет оцениваться проверяющим. При этом оценка может быть субъективной и апелляции не подлежит.

За реализацию модификаций алгоритмов и хорошие дополнительные эксперименты могут быть начислены дополнительные баллы. Начисление этих баллов является субъективным и безапелляционным.

Важно: Практическое задание выполняется самостоятельно. Если вы получили ценные советы (по реализации или проведению экспериментов) от другого студента, то об этом должно быть явно написано в отчёте. В противном случае «похожие» решения считаются плагиатом и все задействованные студенты (в том числе те, у кого списали) будут сурово наказаны.