

Вводная лекция

Методы оптимизации

Александр Безносиков

Московский физико-технический институт

5 сентября 2024



Команда курса: лектор

- Безносиков Александр Николаевич
- почта: beznosikov.an@phystech.edu, anbeznosikov@gmail.com
- tg: @abeznosikov

Команда курса: семинаристы

- Андреев Артем Викторович
tg: @artyomandreyev
- Богданов Александр Иванович
tg: @d0dos
- Былинкин Дмитрий Андреевич
tg: @lxstsvund
- Кормаков Георгий Владимирович
tg: @gkormakov
- Корнилов Никита Максимович
tg: @Tugnir

Команда курса: лектор и семинаристы

- Моложавенко Александр Александрович
tg: @MetaMelon
- Ткаченко Светлана
tg: @Aikiseito
- Чежегов Савелий Андреевич
tg: @Savochak
- Юдин Никита Евгеньевич
tg: @nikitayudin

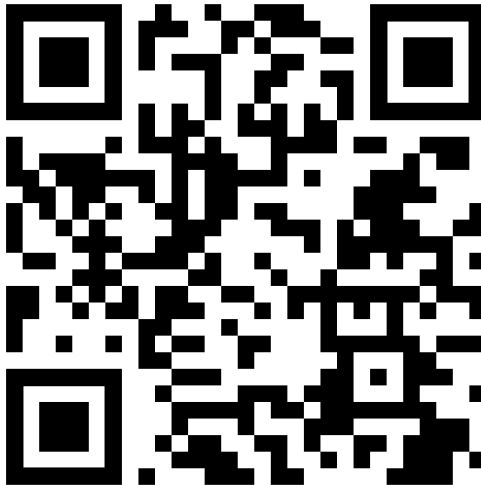
Команда курса: ассистенты

- Кузнецова Алина
- Лисов Петр
- Пичугин Александр
- Ребриков Алексей
- Торопин Иван
- Хафизов Фанис
- Янко Иван

Команда курса: помощь с материалами

- Баширов Наиль
- Подлипнова Ирина
- Прохоров Борис
- Шестаков Александр

Telegram беседа



Диск с материалами



Форма обратной связи



Таблица с оценками



Правила игры: комментарии

- Коллоквиум проходит в конце семестра во время последнего семинара и на зачетной неделе (на выбор). Программа коллоквиума соответствует всей программе курса, изученной в рамках лекций и семинаров. Принимают коллоквиум семинарист и несколько приглашенных преподавателей. Процедура коллоквиума соответствует процедуре проведения обычного устного экзамена на Физтехе с билетами, дополнительными вопросами/задачами и беседой в рамках курса.
- Разбор статьи предполагает полный разбор популярной статьи в области оптимизации. Необходимо разобраться в идее, в доказательствах, а также воспроизвести численные эксперименты, добавив к ним свои. Отчетность: 15 минутное выступление с презентацией на семинарском занятии. Заявки на разбор статьи принимаются до 31 октября.

Задача оптимизации

- **Вопрос:** какие приложения задач оптимизации знаете/уже встречали?

Машинное обучение

Вопрос: Как формулируется задача оптимизации в машинном обучении?

- На практике обычно имеется только конечная выборка/конечный набор данных $(a_1, b_1), \dots, (a_n, b_n)$, на котором формулируется задача минимизации эмпирической функции потерь:

$$\min_{x \in \mathbb{R}^d} \left\{ \hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \ell(g(x, a_i), b_i) \right\}. \quad (1)$$

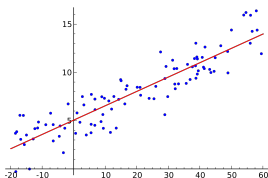
- Однако предполагается, что выборка взята из некоторого распределения \mathcal{D} , и хочется получить хорошее приближенное решение задачи минимизации ожидаемой функции потерь:

$$\min_{x \in \mathbb{R}^d} \{f(x) = \mathbb{E}_{(a,b) \sim \mathcal{D}}[\ell(g(x, a), b)]\}. \quad (2)$$

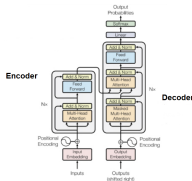
Машинное обучение

Вопрос: Какие примеры ℓ и g уже встречались?

- Линейная регрессия: $g(x, a) = \langle x, a \rangle$, $\ell(b_1, b_2) = (b_1 - b_2)^2$.



- NLP: g – NLP transformer (LLM), $l(b_1, b_2) = -\int b_1(x) \log b_2(x) dx$ (cross-entropy на предсказание следующего слова)



Машинное обучение

- **Вопрос:** Возникает естественный вопрос - как связаны эти задачи? (1) является Монте-Карло аппроксимацией интеграла (2).

Теорема

Если $\ell(g(x, a_i), b_i)$ выпукла по x для любых a_i, b_i , M -Липшицева по x , и мы ограничим область поиска минимума множеством Q с диаметром D , тогда для $\hat{x}^* = \operatorname{argmin}_{x \in Q} \hat{f}(x)$ с вероятностью хотя бы $1 - \delta$

$$\left| f(\hat{x}^*) - \min_{x \in Q} f(x) \right| = O \left(\sqrt{\frac{M^2 D^2 d \ln(n) \ln(d/\delta)}{n}} \right).$$

Машинное обучение

- Популярность оптимизации в машинном обучении

[Adam: A method for stochastic optimization](#)

DP Kingma, J Ba
arXiv preprint arXiv:1412.6980

193709

2014

- Одна из ключевых секций на ведущих ML конференциях.



Статистика

Предположим, что некоторая переменная b зависит от переменных $a_1, a_2, a_3, \dots, a_{d-1}$ линейным образом:

$$b(a_1, \dots, a_{d-1}) = x_0 + x_1 a_1 + x_2 a_2 + \dots + x_{d-1} a_{d-1},$$

где коэффициенты x_0, \dots, x_{d-1} нам неизвестны. Предположим, что мы хотим найти эти коэффициенты, измеряя переменную b при различных значениях a_1, \dots, a_{d-1} . Казалось бы, в этом нет ничего сложного, ведь для решения системы достаточно провести d измерений. **Вопрос:** какая проблема? В действительности измерения производятся с некоторой погрешностью: для заданного набора $a_1^i, a_2^i, \dots, a_{d-1}^i$ мы измеряем

$$b_i = x_0 + x_1 a_1 + \dots + x_{d-1} a_{d-1} + \xi_i,$$

где $\xi_i \sim \mathcal{N}(0, \sigma^2)$.

Статистика

- Другими словами, мы предполагаем, что $b_i \sim \mathcal{N}(x_0 + x_1 a_1 + \dots + x_{d-1} a_{d-1}, \sigma^2)$, где параметры $x = (x_0, \dots, x_{d-1})^\top$ должны быть найдены по выборке $\{b_i\}_{i=1}^n$ (мы будем считать, что b_1, \dots, b_n – независимые случайные величины). **Вопрос:** Из каких соображений выбрать параметры x_0, \dots, x_{d-1} ?
- Можно, например, рассмотреть оценку максимального правдоподобия (для краткости введем вектор $x^i = (1, x_1^i, \dots, x_{d-1}^i)^\top$):

$$\begin{aligned} x^* &= \arg \max_{x \in \mathbb{R}^d} \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{1}{2\sigma^2} (b_i - \langle a^i, x \rangle)^2 \right) \\ &= \arg \max_{x \in \mathbb{R}^d} \ln \left(\prod_{i=1}^d \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{1}{2\sigma^2} (b_i - \langle a^i, x \rangle)^2 \right) \right). \end{aligned}$$

Статистика

- Поскольку логарифм произведения равен сумме логарифмов, а аддитивные и мультипликативные константы не меняют точку оптимума, получаем:

$$\begin{aligned}x^* &= \arg \max_{x \in \mathbb{R}^d} \left\{ \text{Const} + \sum_{i=1}^n -\frac{1}{2\sigma^2} (b_i - \langle a^i, x \rangle)^2 \right\} \\&= \arg \min_{x \in \mathbb{R}^d} \sum_{i=1}^n (b_i - \langle a^i, x \rangle)^2 \\&= \arg \min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2,\end{aligned}$$

где A составлена из строк $(a^i)^\top$.

Рюкзак

- На производстве могут выпускать товары d типов, нужно оптимально организовать производство (разработать план $\{x_i\}_{i=1}^d$, сколько товаров каждого типа производить). Чтобы
 - 1 выполнялись ограничения на рабочую силу (для производства 1 у.е. товара i го типа в месяц нужно a_i работников, всего работников A):
$$\sum_{i=1}^d a_i x_i \leq A,$$
 - 2 была возможность закупить сырье (для производства 1 у.е. товара i го типа в месяц нужно b_j сырья j при этом цена формируется в зависимости от размера заказа по функции $c_j(\cdot)$, всего средств на закупку C):
$$\sum_{i=1}^d \sum_{j=1}^n c_j(b_j x_i) \leq C,$$
 - 3 а главное максимизировалась прибыль от продажи (цена за 1 у.е. товара i го типа d_i):
$$\max_{\{x_i\} \geq 0} d_i x_i.$$

Вопрос: по-хорошему x_i дискретны, а мы решаем непрерывную задачу это сильно страшно? Зависит от масштабов производства: округление может дать хороший результат, но существуют и более хитрые трюки. как решать дискретные задачи через непрерывные.

Задача коммивояжера

- На плоскости заданы N точек с координатами $\{x_i, y_i\}_{i=1}^N$, нужно построить путь кратчайшей длины, проходящий через все точки. Формально:

$$\min_{\text{path}} \left[\sum_{i=1}^{N-1} (x_{\text{path}(i+1)} - x_{\text{path}(i)})^2 + (y_{\text{path}(i+1)} - y_{\text{path}(i)})^2 \right].$$

Существует масса способов ее решения, один из самых популярных – линейное программирование.



Методы оптимизации

- Нет смысла искать лучший метод для решения конкретной задачи. Например, лучший метод для решений задачи $\min_{x \in \mathbb{R}^d} \|x\|^2$ сходится за 1 итерацию: этот метод просто всегда выдаёт ответ $x^* = 0$. Очевидно, что для других задач такой метод не пригоден.
- Эффективность метода определяется для класса задач, т.к. обычно численные методы разрабатываются для *приближённого* решения множества однотипных задач.
- Метод разрабатывается для класса задач \implies метод не может иметь с самого начала полной информации о задаче. Вместо этого метод использует модель задачи, например, формулировку задачи, описание функциональных компонент, множества, на котором происходит оптимизация и т.д.

- Предполагается, что численный метод может накапливать специфическую информацию о задаче при помощи некоторого *оракула*. Под оракулом можно понимать некоторое устройство (программу, процедуру), которое отвечает на последовательные вопросы численного метода.

Вопрос: Какого рода вопросы хочется задавать оракулу?

Примеры оракулов

- **Оракул нулевого порядка** в запрашиваемой точке x возвращает значение целевой функции $f(x)$.
- **Оракул первого порядка** в запрашиваемой точке возвращает значение функции $f(x)$ и её градиент в данной точке
$$\nabla f(x) = \left(\frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_n} \right).$$
- **Оракул второго порядка** в запрашиваемой точке возвращает значение и градиент функции $f(x)$, $\nabla f(x)$, а также её гессиан в данной точке $(\nabla^2 f(x))_{ij} = \frac{\partial^2 f(x)}{\partial x_j \partial x_i}$.

Примеры итерационных методов. Градиентный спуск

Рассмотрим задачу оптимизации

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

где функция $f(x)$ дифференцируема. Предположим, что в любой точке мы можем посчитать её градиент.

Алгоритм 1 Градиентный спуск с постоянным размером шага

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

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

Выход: x^K

Вопрос: в чем Алгоритм 1 отличается от определения общей итеративной схемы? В итеративной схеме использовался \mathcal{T}_ϵ .

Критерии останова

- По аргументу: $\|x^k - x^*\| \leq \varepsilon$.

Вопрос: какие проблемы тут видим?

- x^* – неизвестно, но можно так:

$$\|x^{k+1} - x^k\| \leq \|x^{k+1} - x^*\| + \|x^k - x^*\|.$$

Тогда если $\|x^{k+1} - x^*\| \leq \|x^k - x^*\| \leq \varepsilon$, следует $\|x^{k+1} - x^k\| \leq 2\varepsilon$ (в обратную сторону, очевидно, неверно). $\|x^{k+1} - x^k\| \leq \varepsilon$ – это скорее практический вариант критерия, который работает, если есть понимание (интуиция), что $\|x^k - x^*\| \rightarrow 0$.

- x^* – не уникально. Тогда можно поменять следующий критерий
- По функции: $f(x^k) - f^* \leq \varepsilon$.

Часто f^* известно, например, для $f(x) = \|Ax - b\|^2$. На практике можно использовать $|f(x^k) - f(x^{k+1})|$.

- По норме градиента: $\|\nabla f(x^k)\| \leq \varepsilon$.

Вопрос: когда такой критерий можно использовать? В безусловной оптимизации

Сложность методов оптимизации

- **Аналитическая/Оракульная сложность** — число обращений к оракулу, необходимое для решения задачи с точностью ε .
- **Арифметическая/Временная сложность** — общее число вычислений (включая работу оракула), необходимых для решения задачи с точностью ε .