

Вычислительные аспекты оптимизации. Гладкие функционалы и пр. Метод стохастического градиента как метод оптимизации

Бакшинская Екатерина
Зенкова Наталья
Балагуров Владимир



Санкт-Петербург
2019г.

Обучающая выборка: $X^n = (x_i, y_i)_{i=1}^n$, $x_i \in \mathbb{R}^p$

В регрессии:

- Целевая переменная $y_i \in \mathbb{R}$.
- Задача:

$$Q(w) = \sum_{i=1}^n (a(x_i, w) - y_i)^2 \rightarrow \min_w$$

В классификации:

- Целевая переменная $y_i \in \{-1, +1\}$.
- Задача:

$$Q(w) = \sum_{i=1}^n [a(x_i, w)y_i < 0] \leq \sum_{i=1}^n \mathcal{L}(M_i(w)) \rightarrow \min_w,$$

где $M_i(w) = \langle x_i, w \rangle y_i$ — отступ (margin) объекта x_i

- Оптимизация часто условная. Решение проблемы — сведение условной оптимизации к безусловной.
- Оптимизируемая функция бывает негладкой (или вообще не непрерывной). Решение — аппроксимация гладкой.

Для классификации вспомним постановку задачи в SVM:

$$\sum_{i=1}^n (1 - M_i(w))_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_w,$$

здесь $M_i(w) = y_i(\langle w, x_i \rangle - w_0)$.

Эквивалентная задача безусловной минимизации:

$$C \sum_{i=1}^n (1 - M_i(w))_+ + \frac{1}{2} \|w\|^2 \rightarrow \min_w,$$

Минимизация эмпирического риска (регрессия, классификация):

$$Q(w) = \sum_{i=1}^n \mathcal{L}_i(w) \rightarrow \min_w,$$

Численная минимизация методом градиентного спуска:

Вход: $X^n = (x_i, y_i)_{i=1}^n$, h

Выход: w

❶ Инициализация $w^{(0)}$

❷ **Повторять**

❶ Вычислять: $\nabla Q(w^t)$

❷ Градиентный шаг: $w^{t+1} = w^t - h \nabla Q(w^t)$

пока $\|w^t - w^{t-1}\| > \varepsilon$

Варианты выбора градиентного шага и сходимость

- сходимость гарантируется (для выпуклых функций) при

$$h_t \rightarrow 0, \quad \sum_{t=1}^{\infty} h_t = \infty, \quad \sum_{t=1}^{\infty} h_t^2 < \infty$$

В частности можно положить $h_t = 1/t$

- метод скорейшего градиентного спуска:

$$Q(w - h \nabla Q(w)) \rightarrow \min_h$$

позволяет найти адаптивный шаг h^*

- периодически можно делать пробные случайные шаги для “выхода” из локальных минимумов

и др.

Если функционал $Q(w)$ выпуклый, гладкий и имеет минимум w^* , то имеет место следующая оценка сходимости:

$$Q(w^t) - Q(w^*) = O\left(\frac{1}{t}\right)$$

Проблема — минимизируемая функция представляет собой сумму слагаемых $\mathcal{L}_i(w)$, количество которых равно объему выборки.

Вычисление $Q(w)$ и $\nabla Q(w)$ становится трудоемким.

Метод стохастического градиента: вычисляем не точное значение градиента, а его (случайную и желательно несмещенную) оценку.

Оценить градиент суммы можно градиентом одного случайно взятого слагаемого:

$$\nabla \bar{Q}(w) \approx \nabla \mathcal{L}_i(w),$$

где i — равномерно распределены на $1, \dots, n$. Градиентный шаг переписывается в виде:

$$w^{t+1} = w^t - h \nabla \mathcal{L}_i(w)$$

После чего вычисляем оценку функционала

$$\bar{Q}(w) = \frac{1}{n} Q(w) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i, \text{ которая является несмещенной т.к. } \mathbb{E} \nabla \mathcal{L}_i = \nabla \bar{Q}(w)$$

Для выпуклого и гладкого функционала может быть получена следующая оценка:

$$\mathbb{E}[Q(w^t) - Q(w^*)] = O\left(\frac{1}{\sqrt{t}}\right)$$

Вход: $X^n = (x_i, y_i)_{i=1}^n$, h , λ

Выход: w

- ❶ Инициализация $w_j, j = 0, \dots, p$
- ❷ Инициализация $\bar{Q}(w) = \frac{1}{n}Q(w) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(w)$
- ❸ **Повторять**
 - ❶ Выбор x_i из X^n случайным образом
 - ❷ Вычисление значения функции потерь: $\varepsilon_i = \mathcal{L}_i(w)$
 - ❸ Градиентный шаг: $w^{t+1} = w^t - h \nabla \mathcal{L}_i(w)$
 - ❹ Оценка функционала: $\bar{Q}(w) = (1 - \lambda)\bar{Q}(w) + \lambda \varepsilon_i$

пока значение $\bar{Q}(w)$ и/или веса w не сойдутся

Проблема: после каждого шага w по одному объекту x_i не хотим оценивать Q по всей выборке x_1, \dots, x_n

Решение: использование рекуррентной формулы

Среднее арифметическое $\bar{Q}_m = \frac{1}{m} \sum_{i=1}^m \varepsilon_i$:

$$\bar{Q}_m = (1 - \frac{1}{m})\bar{Q}_{m-1} + \frac{1}{m}\varepsilon_m$$

Экспоненциальное скользящее среднее:

$$\bar{Q}_m = (1 - \lambda)\bar{Q}_{m-1} + \lambda\varepsilon_m$$

$$\bar{Q}_m = \lambda\varepsilon_m + \lambda(1 - \lambda)\varepsilon_{m-1} + \lambda(1 - \lambda)^2\varepsilon_{m-2} + \dots$$

Параметр $\lambda \approx \frac{1}{m}$ — *темп забывания*

Например, можно повысить точность оценки градиента, используя несколько слагаемых вместо одного (**mini-batch gradient descent**):

$$\nabla Q(w) \approx \frac{1}{m} \sum_{j=1}^m \nabla \mathcal{L}_{i_j}(w),$$

где i_j — случайно выбранные номера слагаемых из функционала, а m — параметр метода.

Stochastic average gradient:

$$\nabla Q(w) \approx \frac{1}{n} \sum_{i=1}^n z_i^{(t)},$$

$$\text{где } z_i^{(t)} = \begin{cases} \nabla \mathcal{L}_i(w^{t-1}) & \text{если } i = i_t \\ z_i^{t-1} & \text{иначе} \end{cases} \quad \text{и } z_i^{(0)} = \nabla \mathcal{L}_i(w^{(0)})$$

Сходимость метода стохастического градиента

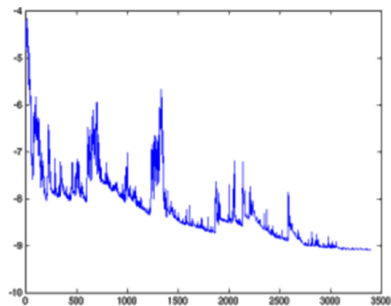
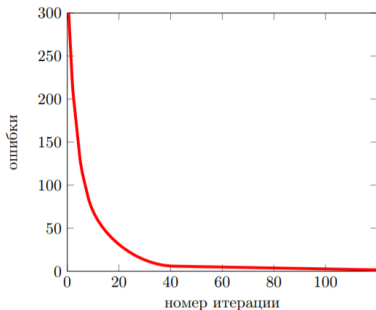


Рис. 1: Ошибка в зависимости от номера итерации

Чтобы ограничить рост абсолютных значений весов, к минимизируемому функционалу $Q(w)$ добавляется штрафное слагаемое:

$$Q_{\tau}(w) = Q(w) + \frac{\tau}{2} \|w\|^2 = Q(w) + \frac{\tau}{2} \sum_{j=1}^n w_j^2 \rightarrow \min_w$$

Аддитивная поправка в градиенте:

$$\nabla Q_{\tau}(w) = \nabla Q(w) + \tau w$$

Правило обновления весов принимает вид:

$$w^t = w^{t-1} (1 - h\tau) - h \nabla Q(w^{t-1})$$

Недостаток: параметр τ приходится подбирать в режиме скользящего контроля, что связано с большими вычислительными затратами.

- $w_j^{(0)} = 0$ для всех $j = 0, \dots, n$
- небольшие случайные значения: $w_j^{(0)} = \text{random}(-\frac{1}{2n}, \frac{1}{2n})$
- для регрессии: $w_j^{(0)} = \frac{\langle y, x_j \rangle}{\langle x_j, x_j \rangle}$
- для классификации: $w_j^{(0)} = \ln \frac{\sum_i [y_i = +1] x_{ij} \sum_i [y_i = -1]}{\sum_i [y_i = -1] x_{ij} \sum_i [y_i = +1]}$
- оценки $w_j^{(0)}$ по небольшой случайной подвыборке объектов
- мультистарт: многократные запуски из разных случайных приближений и выбор лучшего решения

Несколько вариантов выбора наблюдения (специфично для классификации), помимо выборки из выборочного распределения:

- перетасовка объектов (shuffling): попеременно брать объекты из разных классов.
- чаще брать те объекты, на которых была допущена бóльшая ошибка:
(чем меньше M_i , тем больше вероятность взять объект)
(чем меньше $|M_i|$, тем больше вероятность взять объект)
- вообще не брать объекты у которых $M_i > \mu_+$
- вообще не брать объекты у которых $M_i < \mu_-$

параметры μ_+ , μ_- придется подбирать

Преимущества и недостатки метода стохастического градиента

Преимущества:

- Легко реализуется.
- Функция потерь и семейство алгоритмов могут быть любыми.
- Легко добавить регуляризацию.
- Метод подходит для динамического обучения, когда обучающие объекты поступают потоком, и вектор весов обновляется при появлении каждого объекта.
- Подходит для задач с большими данными, иногда можно получить решение даже не обработав всю выборку.

Недостатки:

- Подбор эвристик является искусством (не забыть про переобучение, застревание, расходимость)