

Достоинства:

- 1 легко реализуется;
- 2 применим к любым моделям и функциям потерь;
- 3 допускает онлайнное (потокковое) обучение;
- 4 на сверхбольших выборках позволяет получать неплохие решения, даже не обработав все (x_i, y_i) ;
- 5 всё чаще применяется для Big Data.

Недостатки:

- 1 возможно застревание в локальных экстремумах;
- 2 возможна расходимость или медленная сходимость;
- 3 возможно переобучение;
- 4 подбор комплекса эвристик является искусством.

- 1 $w_j := 0$ для всех $j = 0, \dots, n$;
- 2 небольшие случайные значения:
 $w_j := \text{random} \left(-\frac{1}{2n}, \frac{1}{2n} \right)$;
- 3 $w_j := \frac{\langle y, f_j \rangle}{\langle f_j, f_j \rangle}$, $f_j = (f_j(x_i))_{i=1}^\ell$ — вектор значений признака;

эта оценка w оптимальна при квадратичной функции потерь, если признаки некоррелированы, $\langle f_j, f_k \rangle = 0$, $j \neq k$.

- 4 $w_j := \ln \frac{\sum_i [y_i=+1] f_j(x_i)}{\sum_i [y_i=-1] f_j(x_i)} \frac{\sum_i [y_i=-1]}{\sum_i [y_i=+1]}$;

эта оценка w оптимальна для задач классификации, $Y = \{-1, +1\}$, если признаки независимы.

- 5 оценки w_j по небольшой случайной подвыборке объектов;
- 6 мултистарт: многократные запуски из разных случайных начальных приближений и выбор лучшего решения.

Возможны варианты:

- 1 *перетасовка объектов (shuffling)*:
попеременно брать объекты из разных классов;
- 2 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);
- 3 вообще не брать «хорошие» объекты, у которых $M_i > \mu_+$
(при этом немного ускоряется сходимость);
- 4 вообще не брать объекты-«выбросы», у которых $M_i < \mu_-$
(при этом может улучшиться качество классификации);

Параметры μ_+ , μ_- придётся подбирать.

- ❶ 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 $h_t = 1/t$;

- ❷ *метод скорейшего градиентного спуска:*

$$\mathcal{L}_i(w - h \nabla \mathcal{L}_i(w)) \rightarrow \min_h,$$

позволяет найти *адаптивный шаг* h^* ;

при квадратичной функции потерь $h^* = \|x_i\|^{-2}$;

- ❸ периодически можно делать пробные случайные шаги для «выбивания» из локальных минимумов;
- ❹ метод Левенберга-Марквардта (второго порядка)

Метод Ньютона-Рафсона, $\mathcal{L}_i(w) \equiv \mathcal{L}(\langle w, x_i \rangle y_i)$:

$$w := w - h(\mathcal{L}_i''(w))^{-1} \nabla \mathcal{L}_i(w),$$

где $\mathcal{L}_i''(w) = \left(\frac{\partial^2 \mathcal{L}_i(w)}{\partial w_j \partial w_{j'}} \right)$ — гессиан, $n \times n$ -матрица

Эвристика: считаем, что гессиан диагонален. Тогда

$$w_j := w_j - h \left(\frac{\partial^2 \mathcal{L}_i(w)}{\partial w_j^2} + \mu \right)^{-1} \frac{\partial \mathcal{L}_i(w)}{\partial w_j},$$

h — темп обучения, можно полагать $h = 1$

μ — параметр, предотвращающий обнуление знаменателя.

Отношение h/μ есть темп обучения на ровных участках функционала $\mathcal{L}_i(w)$, где вторая производная обнуляется.