

① Коэффициенты в системе глубины

$L(y, z)$ - эмпирич. ф. потерь

$$a_N(x) = \sum_{j=1}^N \gamma_j b_j(x)$$

можно задать
когда система
как деревьями
выбор γ_j зависит
на значениях листьев

Обучение N-ой модели:

$$S_i = - \frac{\partial L(y, z)}{\partial z} \Big|_{z = a_{N-1}(x)} \text{ "псевдоостатки"}$$

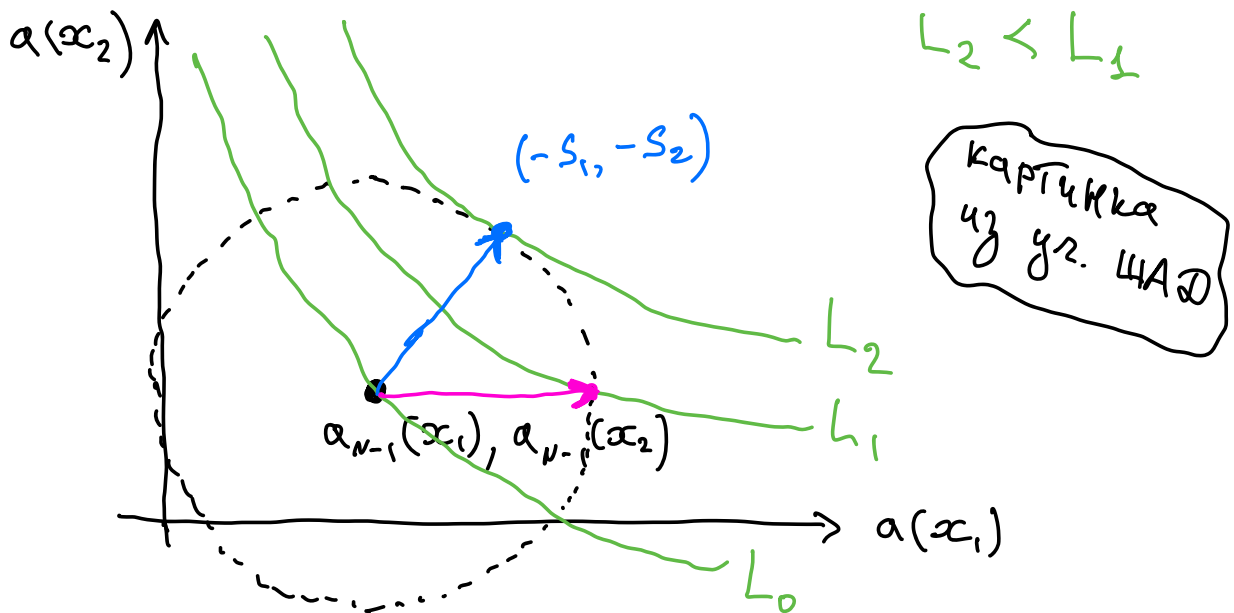
$$① \frac{1}{n} \sum_{i=1}^n (b_N(x) - S_i)^2 \rightarrow \min_{b_N}$$

$$② \frac{1}{n} \sum_{i=1}^n L(y_i; a_{N-1}(x) + \gamma \cdot b_N(x)) \rightarrow \min_{\gamma}$$

① Почему $b_N(x)$ обучается на S_i , а не $y_i - a_{N-1}(x_i)$?

→ Высокий риск переобучиться на $y_i - a_{N-1}(x_i)$

→ Итерирование исходной ф. потерь



$$(b(x_i) - \varepsilon_i)^2$$

$$(b(x_i) - s_i)^2$$

$$\varepsilon_i = a_N(x_i) - y_i$$

② Почему базовая модель груба на MSE?

а) Разложение в разности ф. в ряд

Тейлора по 2^{ому} слагаемому

б) Сдвиги s_i уже включают инфу про L

Разумно аппроксимировать s_i с одинаковой
взвешенностью

$$L(y, z) = \frac{1}{2} \cdot (10[z \geq y] + 1 - [z < y]) \cdot (y - z)^2$$

$$y_1 = 0 \quad a_{n-1}(x_1) = 5 \quad L = 125 \quad S_1 = -50$$

$$y_2 = 0 \quad a_{n-1}(x_2) = -5 \quad L = 12.5 \quad S_2 = 5$$

Если мы будем считать b_n как L , то тогда
при $z \geq y$ будет правильное числ. выражение

$$(y - z)^2 \quad |y - z|$$

$$\sum_{i=1}^n (b_n(x_i) - s_i)^2 = \sum_{i=1}^n (b_n^2(x_i) - 2 \cdot s_i \cdot b_n(x_i) + s_i^2)$$

$$\sum_{i=1}^n b_n^2(x_i) - 2 \cdot \|s\| \cdot \|b_n(x)\| \cdot \cos(s, b_n(x))$$

$$\langle x, y \rangle = \sum x_i y_i$$

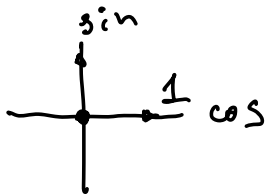
$$\rightarrow \min_{b_n(x)}$$

$$\cos(x, y) = \frac{\langle x, y \rangle}{\|s\| \cdot \|b_n(x)\|}$$

$$\|x\| = \sqrt{\sum x_i^2}$$

Т.е. мы хотим с некоторыми оговорками,

$$\text{чтобы } \cos(s, b_n) \uparrow \Rightarrow \angle(s, b_n) = 0$$



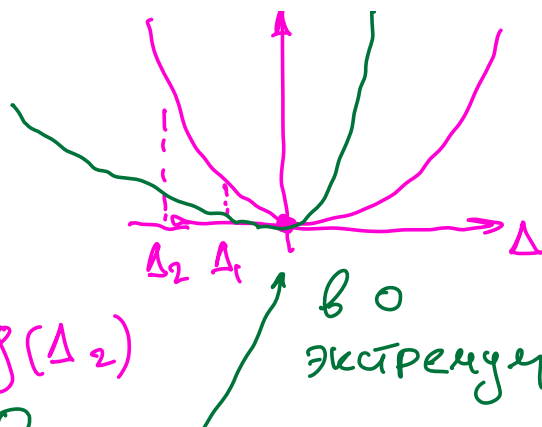
.L.

Типо а)

$$\bullet L(y, z) = g(y - z)$$

$$\bullet g(0) = 0$$

$$\bullet |\Delta_1| \leq |\Delta_2| \Rightarrow g(\Delta_1) \leq g(\Delta_2)$$



$$\begin{aligned} L(y, z) = g(y - z) &\approx \overbrace{g(0)}^0 + \underbrace{g'(0)}_0 \cdot (y - z) + \\ &+ \frac{g''(0)}{2!} (y - z)^2 = \\ &= C \cdot (y - z)^2 \end{aligned}$$

③ Почему $v_n(x)$ обычно используют неглубокие деревья?

→ глубокие деревья больше учить

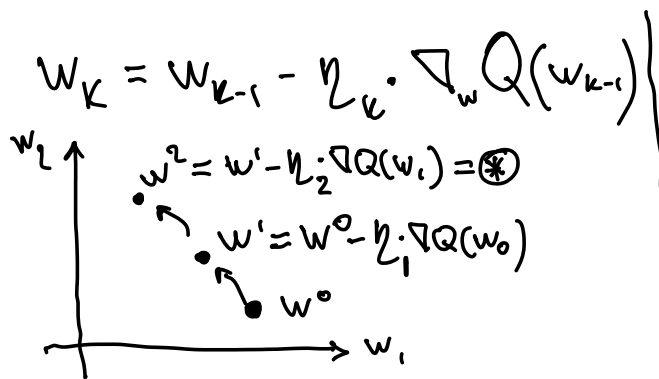
→ если базовая модель сложная, её можно переобучить.

Мы учим "качественный" ансамбль \Rightarrow переобучение уменьшается

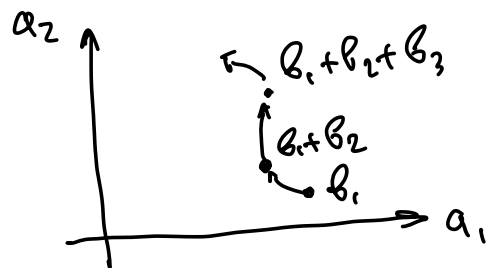
④ Почему говорят, что градиентный бустинг — это градиентный спуск в пространстве прогнозов композиции

на обучающей выборке?

прототип S_i



$$a(x) = \sum_{j=1}^N \overbrace{b_j(x)}^{\text{прототип } S_i} \cdot \gamma_j$$



$$* = w^0 - \eta_1 \nabla Q(w^0) - \eta_2 \nabla Q(w^1)$$

$$\hat{w} = \sum_{k=1}^N (-\eta_k \nabla Q(w_{k-1}))$$

А можно ли интегрировать модификации
град. спуска для системы?

Да, можно! [статья от 2020 года]

Мomentum:

$$h_0 = 0$$

$$h_k = \alpha \cdot h_{k-1} + \eta_k \cdot \nabla_w Q(w_{k-1})$$

$$w_k = w_{k-1} - h_k$$

$a_N(x)$ — композиция алгоритмов

$b_N(x) \simeq \nabla Q$ базовые модели

$h_N(x)$ — «икерция»

$$a_0(x) \quad b_0(x) \quad h_0(x) - \text{init}$$

$$b_n(x) = \arg \min_{b_n} \sum_{i=1}^n \left[b_n(x_i) - \frac{\partial L(z_i, z)}{\partial z} \right]_{z=a_{n-1}(x)}^2$$

$$a_n(x) = a_{n-1}(x) + h_n(x)$$

$$h_n(x) = \lambda \cdot h_{n-1}(x) + \eta \cdot b_n(x)$$

Быстрее сойдётся
хранить больше моделей

⊕ В статье есть
пара хитростей

⑤ Зачем вообще грузить модели, если
бустинг такой классный?

→ Линейные модели интерпретируются

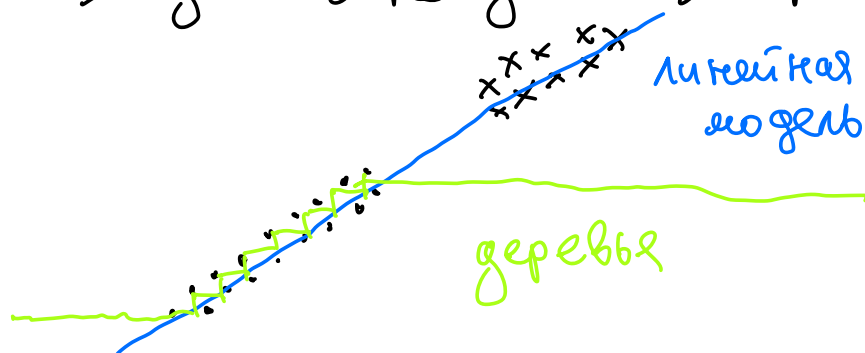
→ Если много признаков и выборка
маленькая бустинг может работать плохо

→ Бустинг Vs Нейросети

ТАБЛИЦЫ

картинки/тексты/видео

→ Бустинг не умеет экстраполировать



② Функция для классификации

$$\begin{array}{l} y_i \in \{-1, 1\} \\ y_i - a_{n-1}(x_i) \in \{-2, 0, 2\} \end{array} \quad \left| \begin{array}{l} ? \\ \cdot \end{array} \right.$$

Давайте договоримся:

$a_n(x_i)$ — уверенность в 1^{ом} классе (логит)

$$a_n(x_i) = \log \frac{P(y=1|x)}{1 - P(y=1|x)}$$

$$P(y \neq 1|x) = \frac{1}{1 + e^{-a_n(x_i)}}$$

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

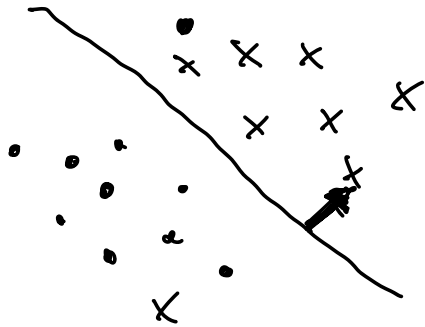
$$\sigma(a_n(x_i)) = P(y=1|x)$$

Давайте играть на logloss

$$L(y, z) = \log(1 + \exp(-yz)) \quad y \in \{-1, 1\}$$

$$y \ln \frac{1}{1 + \exp(-z)} + (1-y) \cdot (1 \dots) \quad y \in \{0, 1\}$$

Играем на $y_i - a_{n-1}(x_i)$ — площадь узора



$$M_i = y_i \cdot a(x_i) \\ \langle w, x \rangle$$

$$y_i = a(x_i) \quad a(x_i) \rightarrow y_i$$

$$y_i a(x_i) \rightarrow 0$$

$$S_i = - \frac{\partial L(y, z)}{\partial z} \bigg|_{z=a_{N-1}(x_i)} = \frac{y}{1 + \exp(-yz)} \cdot \exp(-yz) =$$

$$= \frac{y}{\exp(yz) + 1} \bigg|_{z=a_{N-1}(x_i)} \quad \exp(yz)$$

$$\frac{1}{n} \sum_{i=1}^n \left[b_N(x_i) - \frac{y_i}{1 + \exp(y_i \cdot a_{N-1}(x_i))} \right] \xrightarrow{z} \min_{b_N}$$

$$\frac{y_i \cdot a_{N-1}(x_i)}{0} \simeq +\infty$$

$$S_i \simeq 0$$

ке каго ↑ отступ

$$|S_i| = 1/2$$

сомнения \Rightarrow двичен

$$|S_i| = 1$$

выброс

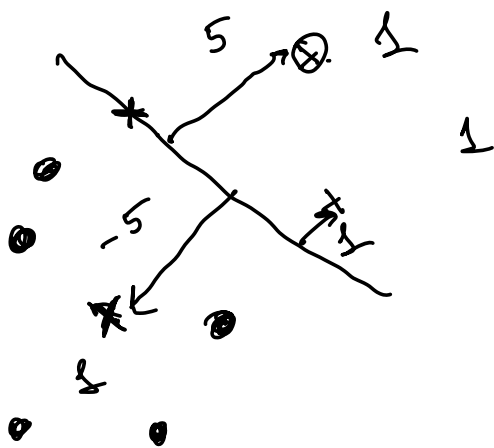
уверен

в неверном
ответе

мы не делаем
на выбросах
сильный акцент

$\exp(-y_i \cdot a_{N-1}(x_i))$ — мера важности

объекта x_i на N -ой итерации бустинга



$$\exp(-5)$$

$$\exp(-1)$$

$$\exp(5)$$

