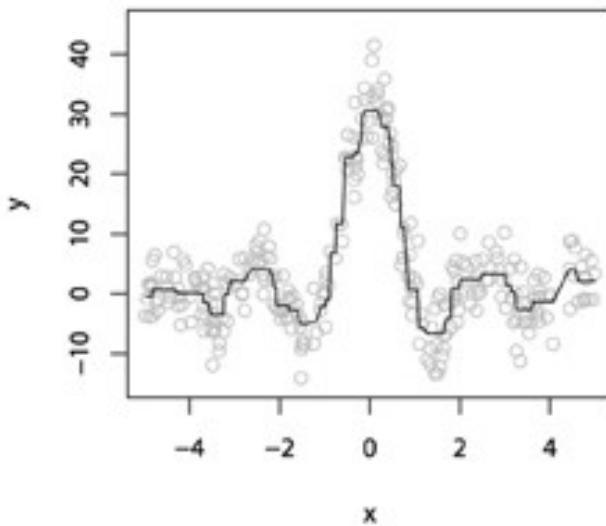
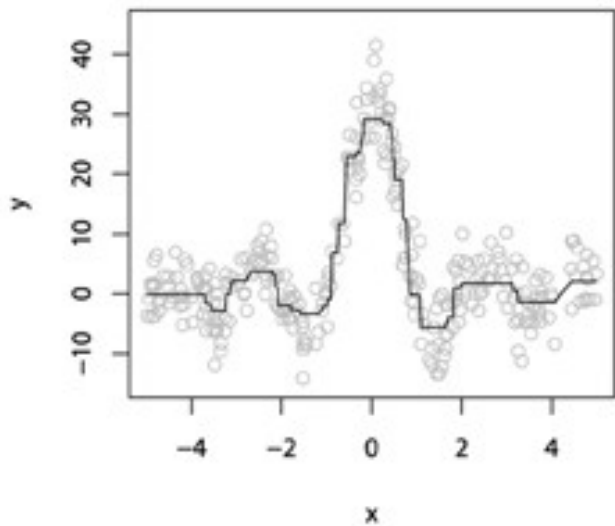
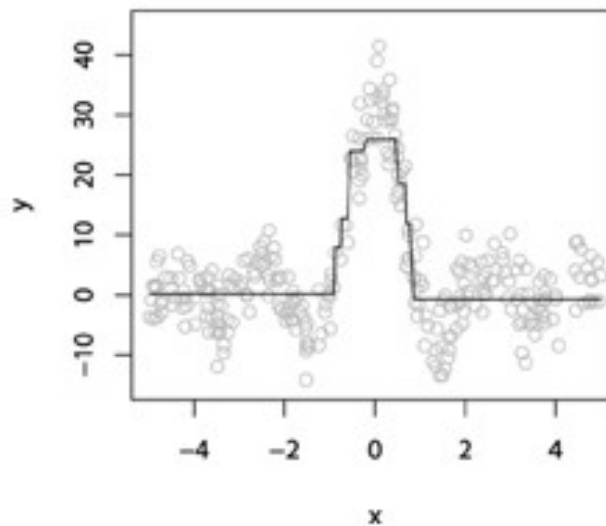
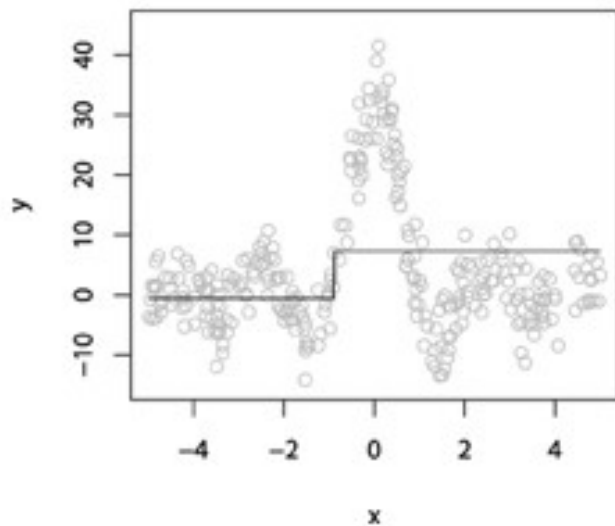


Машинное обучение

Композиции алгоритмов



Содержание лекции

- Определение
- AdaBoost
- AnyBoost
- Градиентный бустинг
- Бэггинг и метод случайных подпространств

Определение

- Алгоритмы классификации часто представимы в виде: $a(x) = C(b(x))$, где функция $b : X \rightarrow R$ называется **алгоритмическим оператором**, $C : R \rightarrow Y$ — **решающим правилом**.
- **Композицией** T алгоритмов $a_t(x) = C(b_t(x))$, $t = 1, \dots, T$ называется суперпозиция алгоритмических операторов $b_t : X \rightarrow R$, корректирующей операции $F : R^T \rightarrow R$ и решающего правила $C : R \rightarrow Y$:

$$a(x) = C(F(b_1(x), \dots, b_T(x)))$$

Примеры

- Простое голосование

$$F(b_1(x), \dots, b_T(x)) = \frac{1}{T} \sum_{t=1}^T b_t(x).$$

- Взвешенное голосование

$$F(b_1(x), \dots, b_T(x)) = \sum_{t=1}^T \alpha_t b_t(x)$$

Общий алгоритм построения композиции

$$a(x) = C(F(b_1(x), \dots, b_T(x))) = \text{sign} \left(\sum_{t=1}^T \alpha_t b_t(x) \right)$$

$$Q_T = \sum_{i=1}^{\ell} \left[y_i \sum_{t=1}^T \alpha_t b_t(x_i) < 0 \right] \rightarrow \mathbf{min}$$

$$Y = \{\pm 1\}, \quad b_t: X \rightarrow \{-1, 0, +1\}, \quad C(b) = \text{sign}(b)$$

$b_t(x) = 0$ — отказ (лучше промолчать, чем соврать)

- При добавлении в композицию слагаемого $\alpha_t b_t(x)$ оптимизируется только базовый алгоритм b_t и коэффициент при нём α_t , а все предыдущие слагаемые $\alpha_1 b_1(x), \dots, \alpha_{t-1} b_{t-1}(x)$ полагаются фиксированными
- Пороговая функция потерь в функционале Q_T аппроксимируется (заменяется) непрерывно дифференцируемой оценкой сверху.

Общий алгоритм построения композиции

- Итерационный процесс:

$$b_1 = \arg \min_b Q(b, X^\ell);$$

$$b_2 = \arg \min_{b, F} Q(F(b_1, b), X^\ell);$$

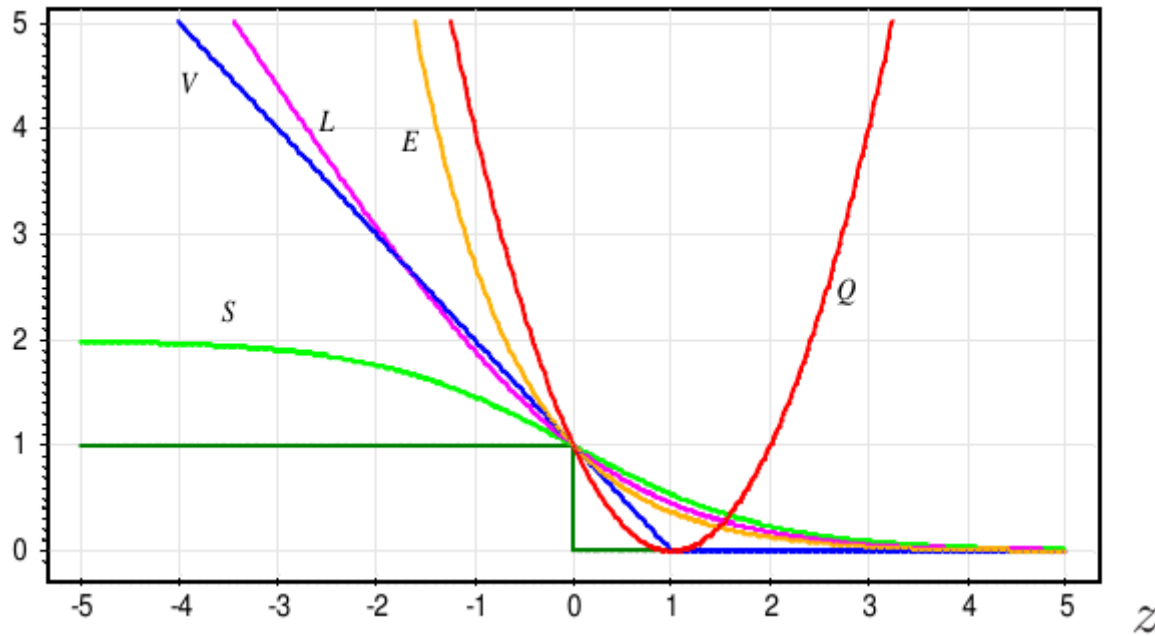
...

$$b_t = \arg \min_{b, F} Q(F(b_1, \dots, b_{t-1}, b), X^\ell)$$

- Алгоритм сильно упрощается, если задача для b_t сводится к исходной с весами объектов и с, возможно, другой функцией потерь:

$$b_t = \arg \min_b \sum_{i=1}^{\ell} w_i \tilde{\mathcal{L}}(b(x_i), y_i)$$

Аппроксимации функции потерь



$S(z) = 2(1 + e^z)^{-1}$ — сигмоидная;

$L(z) = \log_2(1 + e^{-z})$ — логарифмическая;

$V(z) = (1 - z)_+$ — кусочно-линейная;

$E(z) = e^{-z}$ — экспоненциальная;

$Q(z) = (1 - z)^2$ — квадратичная.

AdaBoost

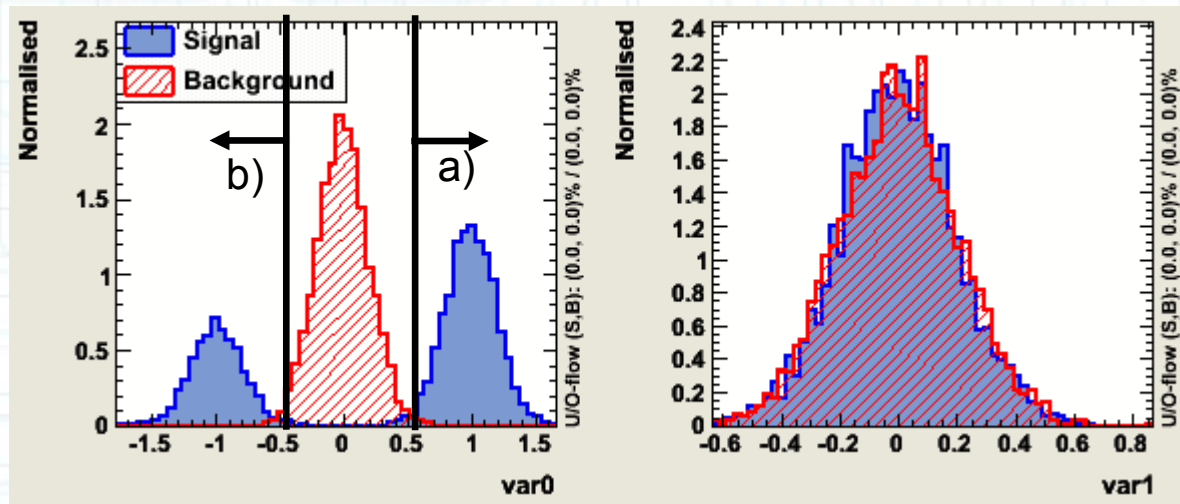
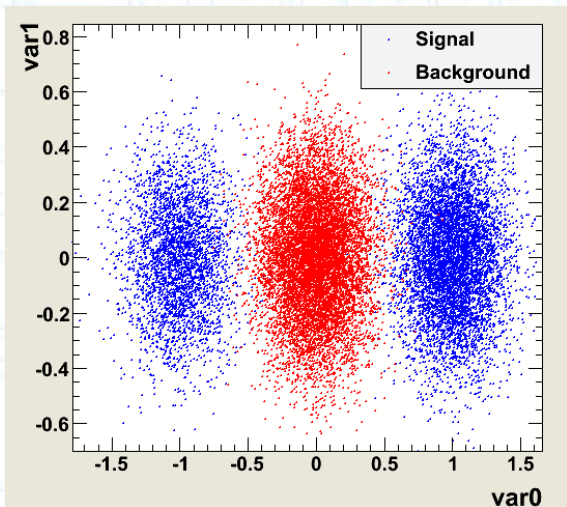
- Исторически первый
- Использует экспоненциальную аппроксимацию $[y_i b(x_i) < 0] \leq e^{-y_i b(x_i)}$
- Оценим функционал качества сверху

$$\begin{aligned} Q_T &\leq \tilde{Q}_T = \sum_{i=1}^{\ell} \exp\left(-y_i \sum_{t=1}^T \alpha_t b_t(x_i)\right) = \\ &= \sum_{i=1}^{\ell} \underbrace{\exp\left(-y_i \sum_{t=1}^{T-1} \alpha_t b_t(x_i)\right)}_{w_i} e^{-y_i \alpha_T b_T(x_i)} \rightarrow \mathbf{min} \end{aligned}$$

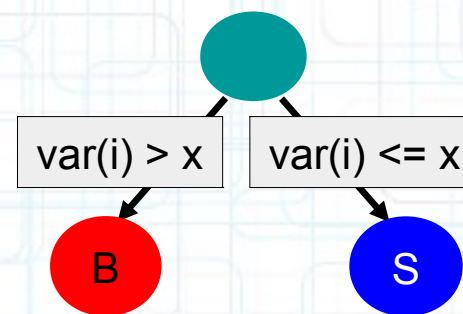
- Таким образом, задача сведена к исходной, но с весами и другой функцией потерь

AdaBoost: демонстрация

Обучающая выборка:



Базовый классификатор:



Всего два нормальных пороговых предиката:

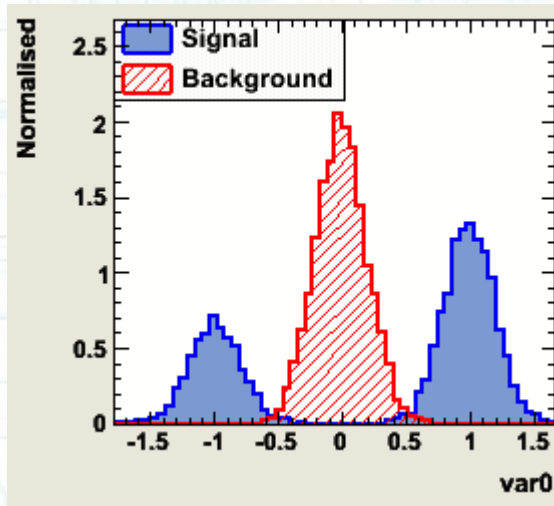
a) $\text{Var0} > 0.5 \rightarrow \epsilon_{\text{signal}} = 66\% \quad \epsilon_{\text{bkg}} \approx 0\%$ - общая ошибка 16.5%

b) $\text{Var0} < -0.5 \rightarrow \epsilon_{\text{signal}} = 33\% \quad \epsilon_{\text{bkg}} \approx 0\%$ - общая ошибка 33%

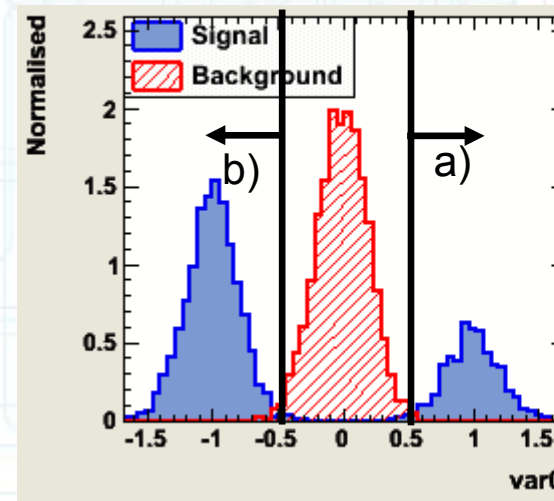
На первой итерации бустинга будет выбрано дерево с разделителем a)

AdaBoost: демонстрация

- Первый классификатор, выбрав разрез a) ошибается в 16.5 % случаев
- Перед построением следующего дерева присвоим объектам веса $\exp(-M)$



Новый
вес



Теперь разрез b) приводит к меньшей ошибке, чем a).
Второе дерево разделит выборку по предикату $\text{Var0} < -0.5$

AdaBoost

- Нормируем веса, чтобы их сумма = 1
- Определим две метрики качества:

$$N(b; U^\ell) = \sum_{i=1}^{\ell} u_i [b(x_i) = -y_i]; \quad P(b; U^\ell) = \sum_{i=1}^{\ell} u_i [b(x_i) = y_i].$$

- Теорема (Freund, Schapire, 1996)

Пусть для любого нормированного вектора весов U^ℓ существует алгоритм $b \in B$, классифицирующий выборку хотя бы немного лучше, чем наугад: $P(b; U^\ell) > N(b; U^\ell)$.

Тогда минимум функционала \tilde{Q}_T достигается при

$$b_T = \arg \max_{b \in B} \sqrt{P(b; \tilde{W}^\ell)} - \sqrt{N(b; \tilde{W}^\ell)}.$$

$$\alpha_T = \frac{1}{2} \ln \frac{P(b_T; \tilde{W}^\ell)}{N(b_T; \tilde{W}^\ell)}.$$

Доказательство

Воспользуемся тождеством $\forall \alpha \in \mathbb{R}, \forall b \in \{-1, 0, +1\}$:

$$e^{-\alpha b} = e^{-\alpha} [b=1] + e^{\alpha} [b=-1] + [b=0].$$

Положим для краткости $\alpha = \alpha_T$ и $b_i = b_T(x_i)$. Тогда

$$\begin{aligned} \tilde{Q}_T &= \left(\underbrace{e^{-\alpha} \sum_{i=1}^{\ell} \tilde{w}_i [b_i = y_i]}_P + \underbrace{e^{\alpha} \sum_{i=1}^{\ell} \tilde{w}_i [b_i = -y_i]}_N + \underbrace{\sum_{i=1}^{\ell} \tilde{w}_i [b_i = 0]}_{1-P-N} \right) \underbrace{\sum_{i=1}^{\ell} w_i}_{\tilde{Q}_{T-1}} \\ &= (e^{-\alpha} P + e^{\alpha} N + (1 - P - N)) \tilde{Q}_{T-1} \rightarrow \min_{\alpha, b}. \end{aligned}$$

$$\frac{\partial}{\partial \alpha} \tilde{Q}_T = (-e^{-\alpha} P + e^{\alpha} N) \tilde{Q}_{T-1} = 0 \Rightarrow e^{-\alpha} P = e^{\alpha} N \Rightarrow e^{2\alpha} = \frac{P}{N}.$$

Получили требуемое: $\alpha_T = \frac{1}{2} \ln \frac{P}{N}$.

Доказательство

Подставим оптимальное значение $\alpha = \frac{1}{2} \ln \frac{P}{N}$ обратно в \tilde{Q}_T :

$$\begin{aligned}\tilde{Q}_T &= (e^{-\alpha} P + e^{\alpha} N + (1 - P - N)) \tilde{Q}_{T-1} = \\ &= (1 + \sqrt{\frac{N}{P}} P + \sqrt{\frac{P}{N}} N - P - N) \tilde{Q}_{T-1} = \\ &= \left(1 - (\sqrt{P} - \sqrt{N})^2\right) \tilde{Q}_{T-1} \rightarrow \min_b.\end{aligned}$$

Поскольку \tilde{Q}_{T-1} не зависит от α_T и b_T , минимизация \tilde{Q}_T эквивалентна либо максимизации $\sqrt{P} - \sqrt{N}$ при $P > N$, либо максимизации $\sqrt{N} - \sqrt{P}$ при $P < N$, однако второй случай исключён условием теоремы.

Получили $b_T = \arg \max_b \sqrt{P} - \sqrt{N}$. Теорема доказана.

Классический AdaBoost

- Пусть отказов нет, $b_t : X \rightarrow \{\pm 1\}$.
Тогда $P = 1 - N$
- Тогда минимум функционала Q_T достигается при

$$b_T = \arg \min_{b \in B} N(b; \widetilde{W}^\ell).$$

$$\alpha_T = \frac{1}{2} \ln \frac{1 - N(b_T; \widetilde{W}^\ell)}{N(b_T; \widetilde{W}^\ell)}.$$

Алгоритм AdaBoost

Вход: обучающая выборка X^ℓ ; параметр T ;

Выход: базовые алгоритмы и их веса $\alpha_t b_t$, $t = 1, \dots, T$;

1: инициализировать веса объектов:

$$w_i := 1/\ell, \quad i = 1, \dots, \ell;$$

2: **для всех** $t = 1, \dots, T$

3: обучить базовый алгоритм:

$$b_t := \arg \min_b N(b; W^\ell);$$

$$4: \quad \alpha_t := \frac{1}{2} \ln \frac{1 - N(b_t; W^\ell)}{N(b_t; W^\ell)};$$

5: обновить веса объектов:

$$w_i := w_i \exp(-\alpha_t y_i b_t(x_i)), \quad i = 1, \dots, \ell;$$

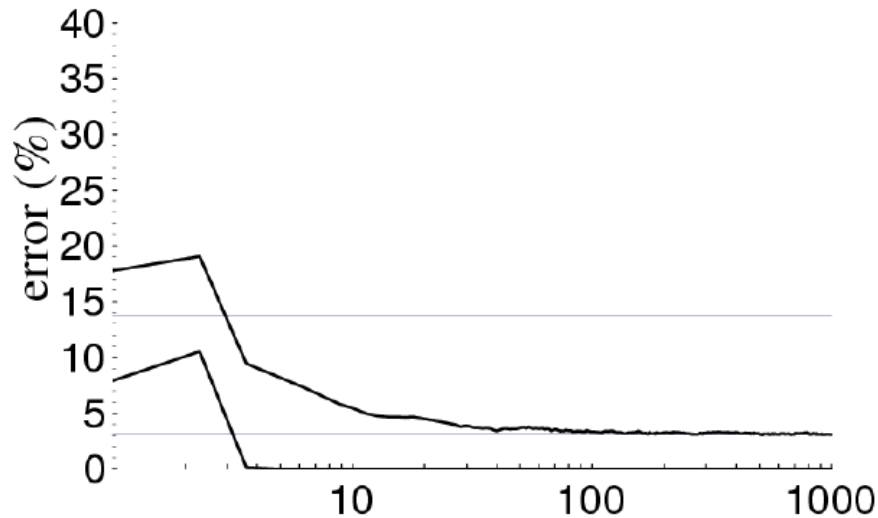
6: нормировать веса объектов:

$$w_0 := \sum_{j=1}^{\ell} w_j;$$

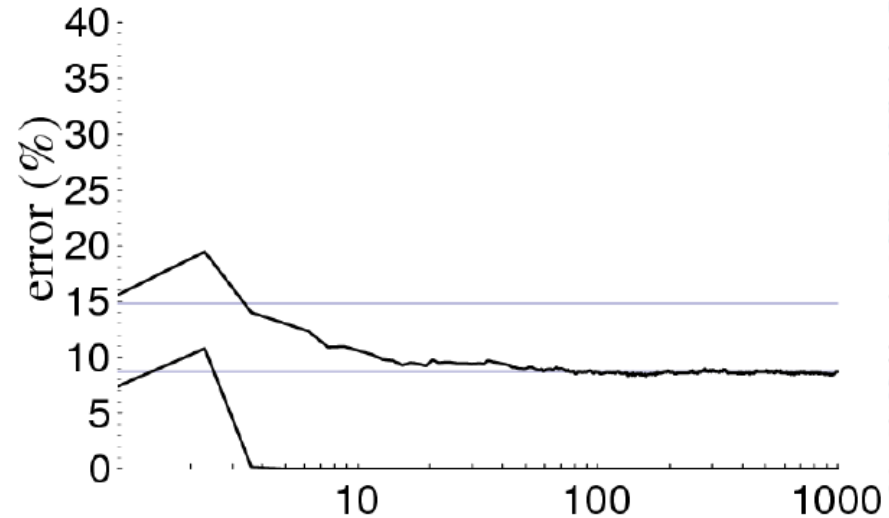
$$w_i := w_i / w_0, \quad i = 1, \dots, \ell;$$

Обобщающая способность не ухудшается с ростом сложности T

Задача UCI:letter



Задача UCI:satimage



Schapire, Freund, Lee, Bartlett (1998) Boosting the margin: a new explanation for the effectiveness of voting methods // Annals of Statistics Vol.26, No.5, Pp. 1651–1686.

Замечания

- Базовые классификаторы (weak classifiers):
 - решающие деревья — используются чаще всего
 - пороговые правила (data stumps)
$$B = \left\{ b(x) = [f_j(x) \leq \theta] \mid j = 1, \dots, n, \theta \in \mathbb{R} \right\};$$
- Базовые классификаторы должны быть слабыми, из сильных хорошую композицию не построить
- Недостатки AdaBoost:
 - чрезмерная чувствительность к выбросам из-за e^{-M}

Обобщение: $\mathcal{L} = \mathcal{L}(M)$.

AnyBoost

Возьмём $Y = \{\pm 1\}$, $b_t: X \rightarrow \mathbb{R}$, $C(b) = \text{sign}(b)$;

$\mathcal{L}(M)$ — функция потерь, гладкая функция отступа M ;

$M_T(x_i) = y_i \sum_{t=1}^T \alpha_t b_t(x_i)$ — отступ композиции на объекте x_i ;

Оценка сверху для числа ошибок композиции:

$$Q_T \leq \tilde{Q}_T = \sum_{i=1}^{\ell} \mathcal{L}(M_{T-1}(x_i) + \alpha y_i b(x_i)) \rightarrow \min_{\alpha, b}.$$

Линеаризация функции потерь по α в окрестности $\alpha = 0$:

$$\tilde{Q}_T \approx \sum_{i=1}^{\ell} \mathcal{L}(M_{T-1}(x_i)) - \alpha \sum_{i=1}^{\ell} \underbrace{-\mathcal{L}'(M_{T-1}(x_i))}_{w_i} y_i b(x_i) \rightarrow \min_b,$$

где w_i — веса объектов.

Обобщение: $\mathcal{L} = \mathcal{L}(M)$.

AnyBoost

Минимизация линейризованного \tilde{Q}_T при фиксированном α

$$\tilde{Q}_T \approx \sum_{i=1}^{\ell} \mathcal{L}(M_{T-1}(x_i)) - \alpha \sum_{i=1}^{\ell} w_i y_i b(x_i) \rightarrow \min_b.$$

приводит к принципу *явной максимизации отступов* (direct optimization of margin, DOOM):

$$\sum_{i=1}^{\ell} w_i y_i b(x_i) \rightarrow \max_b.$$

Затем α определяется путём одномерной минимизации \tilde{Q}_T .

Итерации этих двух шагов приводят к алгоритму AnyBoost.

Замечание. AnyBoost переходит в AdaBoost в частном случае, при $b_t: X \rightarrow \{-1, 0, +1\}$ и $\mathcal{L}(M) = e^{-M}$.

Алгоритм AnyBoost

Вход: обучающая выборка X^ℓ ; **параметр** T ;

Выход: базовые алгоритмы и их веса $\alpha_t b_t$, $t = 1, \dots, T$;

1: инициализировать отступы: $M_i := 0$, $i = 1, \dots, \ell$;

2: **для всех** $t = 1, \dots, T$

3: вычислить веса объектов:

$$w_i = -\mathcal{L}'(M_i), \quad i = 1, \dots, \ell;$$

4: обучить базовый алгоритм согласно принципу DOOM:

$$b_t := \arg \max_b \sum_{i=1}^{\ell} w_i y_i b(x_i);$$

5: решить задачу одномерной минимизации:

$$\alpha_t := \arg \min_{\alpha > 0} \sum_{i=1}^{\ell} \mathcal{L}(M_i + \alpha b_t(x_i) y_i);$$

6: обновить значения отступов:

$$M_i := M_i + \alpha_t b_t(x_i) y_i; \quad i = 1, \dots, \ell;$$

Обобщение: $\mathcal{L} - \forall$.

Градиентный бустинг

Линейная (выпуклая) комбинация базовых алгоритмов:

$$a(x) = \sum_{t=1}^T \alpha_t b_t(x), \quad x \in X, \quad \alpha_t \in \mathbb{R}_+.$$

Функционал качества с произвольной функцией потерь $\mathcal{L}(a, y)$:

$$Q(\alpha, b; X^\ell) = \sum_{i=1}^{\ell} \mathcal{L} \left(\underbrace{\sum_{t=1}^{T-1} \alpha_t b_t(x_i)}_{f_{T-1,i}} + \alpha b(x_i), y_i \right) \rightarrow \min_{\alpha, b}.$$

$\underbrace{\hspace{10em}}_{f_{T,i}}$

$f_{T-1} = (f_{T-1,i})_{i=1}^{\ell}$ — текущее приближение

$f_T = (f_{T,i})_{i=1}^{\ell}$ — искомый вектор, решение задачи $Q(f) \rightarrow \min$

Friedman G. Greedy Function Approximation: A Gradient Boosting Machine. 1999.

Обобщение: $\mathcal{L} - \forall$.

Градиентный бустинг

Градиентный метод минимизации $Q(f) \rightarrow \min, f \in \mathbb{R}^\ell$:

f_0 := начальное приближение;

$$f_{T,i} := f_{T-1,i} - \alpha g_i, \quad i = 1, \dots, \ell;$$

$g_i = \mathcal{L}'(f_{T-1,i}, y_i)$ — компоненты вектора градиента,
 α — градиентный шаг.

Наблюдение: это очень похоже на одну итерацию бустинга!

$$f_{T,i} := f_{T-1,i} + \alpha b(x_i), \quad i = 1, \dots, \ell$$

Идея: будем искать такой базовый алгоритм b_T , чтобы вектор $(b_T(x_i))_{i=1}^\ell$ приближал вектор градиента $(-g_i)_{i=1}^\ell$:

$$b_T := \arg \min_b \sum_{i=1}^{\ell} (b(x_i) + g_i)^2$$

Алгоритм градиентного бустинга

Вход: обучающая выборка X^ℓ ; **параметр** T ;

Выход: базовые алгоритмы и их веса $\alpha_t b_t$, $t = 1, \dots, T$;

1: инициализация: $f_i := 0$, $i = 1, \dots, \ell$;

2: **для всех** $t = 1, \dots, T$

3: найти базовый алгоритм, приближающий градиент:

$$b_t := \arg \min_b \sum_{i=1}^{\ell} (b(x_i) + \mathcal{L}'(f_i, y_i))^2;$$

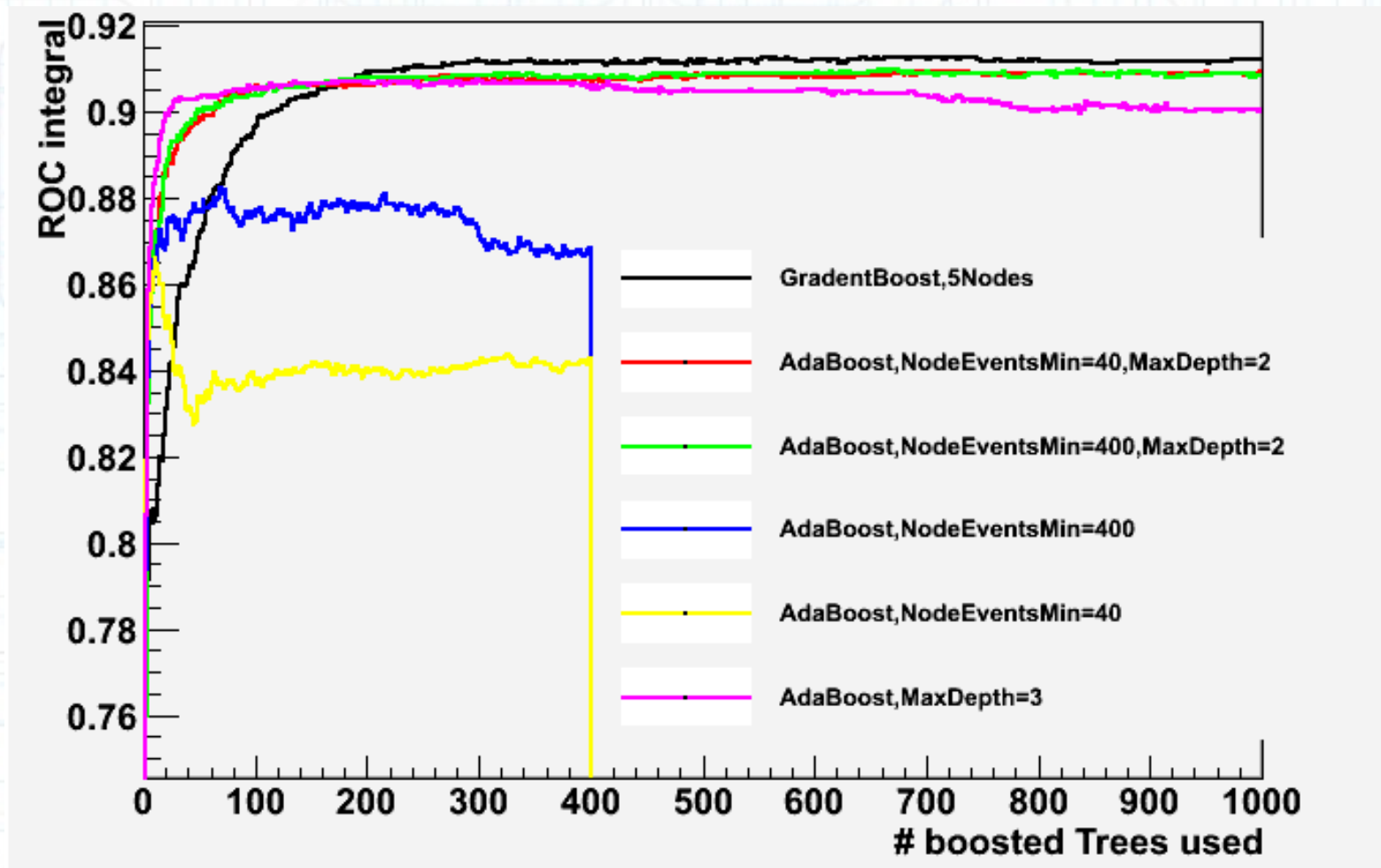
4: решить задачу одномерной минимизации:

$$\alpha_t := \arg \min_{\alpha > 0} \sum_{i=1}^{\ell} \mathcal{L}(f_i + \alpha b_t(x_i), y_i);$$

5: обновить значения композиции на объектах выборки:

$$f_i := f_i + \alpha_t b_t(x_i); \quad i = 1, \dots, \ell;$$

Эксперименты



Бустинг работает лучше всего для “слабых” базовых классификаторов
Необходимо подбирать параметры деревьев

Выводы

- Градиентный бустинг — наиболее общий из всех бустингов:
 - произвольная функция потерь
 - произвольное пространство оценок R
 - подходит для регрессии, классификации, ранжирования
- Интересная интерпретация бустинга: добавление базового алгоритма — это одна итерация градиентного спуска
- Обычно GB применяется к решающим деревьям
- Градиентный бустинг над ODT = Yandex.MatrixNet

Стохастические методы построения композиций

- Bagging – обучает базовые алгоритмы по случайным подвыборкам (подмножество строк матрицы О-П)
- Метод случайных подпространств (RSM) – обучает базовые алгоритмы по случайным подмножествам признаков (подмножество строк матрицы О-П)

Алгоритм

Вход: обучающая выборка X^ℓ ; **параметры:** T
 ℓ' — длина обучающих подвыборок;
 n' — длина признакового подописания;
 ε_1 — порог качества базовых алгоритмов на обучении;
 ε_2 — порог качества базовых алгоритмов на контроле;

Выход: базовые алгоритмы b_t , $t = 1, \dots, T$;

- 1: **для всех** $t = 1, \dots, T$
- 2: $U :=$ случайное подмножество X^ℓ длины ℓ' ;
- 3: $\mathcal{G} :=$ случайное подмножество \mathcal{F} длины n' ;
- 4: $b_t := \mu(\mathcal{G}, U)$;
- 5: **если** $Q(b_t, U) > \varepsilon_1$ или $Q(b_t, X^\ell \setminus U) > \varepsilon_2$ **то**
- 6: не включать b_t в композицию;

Композиция — простое голосование: $a(x) = C \left(\sum_{t=1}^T b_t(x) \right)$.

Сравнение

- Бустинг лучше для больших обучающих выборок и для классов с границами сложной формы.
- Бэггинг и RSM лучше для коротких обучающих выборок.
- RSM лучше в тех случаях, когда признаков больше, чем объектов, или когда много неинформативных признаков.
- Бэггинг и RSM эффективно распараллеливаются, бустинг выполняется строго последовательно.

Случайный лес (Random forest)

- каждое дерево $b_t(x)$ обучается по случайной подвыборке с повторениями (Bagging)
- в каждой вершине каждого дерева рассматривается случайное подмножество из \sqrt{n} признаков (напоминает RSM)
- для ветвления выбирается признак с наилучшим пороговым предикатом по энтропийному критерию
- усечений (pruning) деревьев нет