

- Введение в машинное обучение •

## Композиционные методы машинного обучения

Воронцов Константин Вячеславович

`k.v.vorontsov@phystech.edu`

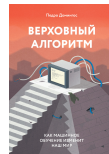
`http://www.MachineLearning.ru/wiki?title=User:Vokov`

Этот курс доступен на странице вики-ресурса

`http://www.MachineLearning.ru/wiki`

«Введение в машинное обучение (курс лекций, К.В.Воронцов)»

- 1 *символизм* – поиск логических закономерностей
    - Decision Tree, Rule Induction
  - 2 *коннекционизм* – обучаемые нейронные сети
    - BackPropagation, Deep Belief Nets, Deep Learning CNN, ResNet, LSTM, GRU, Attention, Transformer
  - 3 *эволюционизм* – саморазвитие сложных моделей
    - Genetic Algorithms, Genetic Programming, Symbolic Regression
  - 4 *байесионизм* и вероятностно-статистические методы
    - MLE, EM, GLM, LR, OBC, Naive Bayes, QD, LDF Bayesian Networks, Bayesian Learning, Graphical Models
  - 5 *аналогизм* – «близким объектам близкие ответы»
    - kNN, RBF, SVM, KDE, Kernel Smoothing
- ⊕ *композиционизм* – кооперация моделей
- Weighted Voting, Boosting, Bagging, Stacking, Random Forest, Яндекс.CatBoost



- 1 Простое голосование**
  - Исторический экскурс
  - Простое голосование. Алгоритм бэггинга
  - Случайные леса
- 2 Взвешенное голосование**
  - Градиентный бустинг
  - Обобщающая способность бустинга
  - Анализ смещения и разброса
- 3 Развитие идеи ансамблирования**
  - Комитетный бустинг
  - Смеси экспертов
  - Философия ансамблирования

## Задачи на малых данных

Особенности геологических данных в задачах поиска месторождений редкого типа (золото, уран, алмазы и т.д.)

- объектов мало (7 + 11), признаков много (более сотни)
- надёжной геофизической модели не существует
- в данных бывают «пропуски» — неизмеренные значения

Группа признаков		Пространственно-временные						Вещественные						
		2	3	4	6	..	..	69	21	22	23	24	..	..
Месторождения М <sub>i</sub>	Витватерсранд (M <sup>1</sup> )	1	0	0	1	0	1	1	1	0	1	0	0	1
	Блайнд-Ривер (M <sup>2</sup> )	1	0	0	1	0	0	0	1	1	1	0	0	1
	Жакобина (M <sup>3</sup> )	0	0	1	0	1	0	1	1	1	1	0	1	1
	Муанга, Габон (M <sup>4</sup> )	0	1	0	0	1	0	—	1	1	0	0	1	0
	Тарква, Гапа (M <sup>5</sup> )	1	0	0	0	0	0	—	1	1	1	0	0	1
	Австралия (M <sup>6</sup> )	0	0	1	—	0	1	—	1	—	—	—	1	1
	Эпо-Кюлия, Финляндия (M <sup>7</sup> )	1	0	0	1	1	1	1	—	—	—	1	—	1

Кренделев Ф. П., Дмитриев А. Н., Журавлев Ю. И. Сравнение геологического строения зарубежных месторождений докембрийских конгломератов с помощью дискретной математики. Доклады АН СССР. 1967

Дмитриев А. Н., Журавлев Ю. И., Кренделев Ф. П. Об одном принципе классификации и прогноза геологических объектов и явлений. 1968.

## Алгоритмы вычисления оценок, АВО (Ю. И. Журавлёв)

Объединение основных на тот момент нестрогих (эвристических) принципов:

- символизм (вывод правил из данных)
- эволюционизм (отбор наилучших правил)
- аналогизм (оценки сходства объектов)
- байесионизм (классы — смеси плотностей)
- коннекционизм (взвешенное голосование)

$$a(x) = \arg \max_{y \in Y} \sum_{i: y_i=y} \sum_{\omega \in \Omega} w_{\omega i} B_{\omega i}(x, x_i)$$

где  $B_{\omega i}$  — бинарные функции сходства по наборам признаков  $\omega$ :

$$B_{\omega i}(x, x_i) = \bigwedge_{j \in \omega} [ |f_j(x) - f_j(x_i)| < \varepsilon ]$$



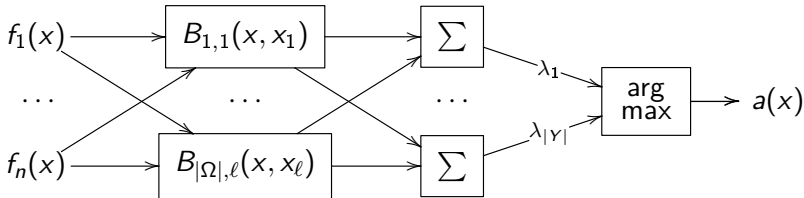
Юрий  
Иванович  
Журавлёв  
(1935–2022)

Журавлёв Ю. И., Никифоров В. В. Алгоритмы распознавания, основанные на вычислении оценок, 1971.

## АВО объединяет многие эвристические принципы

≈ трёхслойная нейросеть RBF (Radial Basis Function):

$$a(x) = \arg \max_{y \in Y} \lambda_y \sum_{i, \omega} [y_i = y] w_{\omega i} B_{\omega i}(x, x_i)$$



≈ метод потенциальных функций  $K(x, x_i) = B_{\omega i}(x, x_i)$

≈ линейный классификатор SVM с радиальным ядром

≈ байесовский классификатор с плотностями-смесями  $p(x|y)$

≈ отбор эталонов:  $w_{\omega i} = 0$  для не-эталонов  $x_i$

≈ отбор признаков в бинарных функциях сходства

## Принципы информативности, непротиворечивости, тупиковости

- *информативность* предиката  $R(x)$  класса  $y \in Y$ :  
$$\begin{cases} p_y(R) = \#\{x_i: R(x_i)=1 \text{ и } y_i=y\} \rightarrow \max \\ n_y(R) = \#\{x_i: R(x_i)=1 \text{ и } y_i \neq y\} \rightarrow \min \end{cases}$$
- *информативность* функции сходства  $B(x, x')$ :  
$$\begin{cases} p(B) = \#\{(x_i, x_j): B(x_i, x_j)=1 \text{ и } y_i=y_j\} \rightarrow \max \\ n(B) = \#\{(x_i, x_j): B(x_i, x_j)=1 \text{ и } y_i \neq y_j\} \rightarrow \min \end{cases}$$
- *непротиворечивость*:  $n_y(R) = 0, n(B) = 0$ 
  - *тест*  $\omega$ :  $B_\omega(x_i, x_j) = 0, \forall i, j: y_i \neq y_j$
  - *представительный набор*  $(\omega, i)$ :  $B_\omega(x_i, x_j) = 0, \forall j: y_i \neq y_j$
- *тупиковость*: никакое подмножество признаков  $\omega' \subset \omega$  не является тестом (или представительным набором)

---

Дмитриев А. Н., Журавлев Ю. И., Кренделев Ф. П. Об одном принципе классификации и прогноза геологических объектов и явлений. 1968.

Журавлёв Ю. И., Никифоров В. В. Алгоритмы распознавания, основанные на вычислении оценок, 1971.

## Задача обучения ансамбля (композиции) моделей

**Дано:**  $X^\ell = (x_i, y_i)_{i=1}^\ell \subset X \times Y$  — обучающая выборка  
 $a_t(x, w) = C(b_t(x, w))$  — «слабые» обучаемые базовые модели  
 $b_t: X \rightarrow R$  — алгоритмические операторы с параметрами  $w$   
 $C: R \rightarrow Y$  — решающее правило простого вида (без параметров)  
 $R$  — удобное пространство оценок

**Найти:** ансамбль  $a(x) = C(F(b_1(x, w_1), \dots, b_T(x, w_T), x, \alpha))$   
 $F: R^T \times X \rightarrow R$  — корректирующая функция с параметрами  $\alpha$

**Критерий** обучения «сильного» алгоритма как ансамбля из  $T$  по-отдельности «слабых» базовых алгоритмов:

$$Q(a, X^\ell) = \sum_{i=1}^{\ell} \mathcal{L}(a(x_i), y_i) \rightarrow \min_{w_1, \dots, w_T, \alpha}$$

Ю.И.Журавлёв. Об алгебраическом подходе к решению задач распознавания или классификации. Проблемы кибернетики, 1978.

M.Kearns, L.G.Valiant. Cryptographic limitations on learning Boolean formulae and finite automata. 1989.

## Пространства оценок $R$ и решающие правила $C$

- **Пример 1:** классификация,  $Y$  — конечное множество,  $R = Y$ ,  $C(b) \equiv b$  — решающее правило не используется.

- **Пример 2:** бинарная классификация,  $Y = \{-1, +1\}$ ,

$$a(x, w) = \text{sign}(b(x, w)),$$

$$R = \mathbb{R}, \quad b: X \rightarrow \mathbb{R} \text{ — real-valued classifier, } C(b) \equiv \text{sign}(b).$$

- **Пример 3:** классификация на  $M$  классов  $Y = \{1, \dots, M\}$ ,

$$a(x, w) = \arg \max_{y \in Y} b_y(x, w_y),$$

$$R = \mathbb{R}^M, \quad b: X \rightarrow \mathbb{R}^M, \quad C(b_1, \dots, b_M) \equiv \arg \max_{y \in Y} b_y.$$

- **Пример 4:** регрессия,  $Y = R = \mathbb{R}$ ,  
 $C(b) \equiv b$  — решающее правило не используется.

## Корректирующие (агрегирующие) функции

Общие требования к агрегирующей функции  $F(b_1, \dots, b_T, x, \alpha)$ :

- $F \in [\min_t b_t, \max_t b_t]$  — среднее по Коши  $\forall x$
- $F$  монотонно не убывает по всем  $b_t$

Примеры агрегирующих функций:

- простое голосование (simple voting):

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

- взвешенное голосование (weighted voting):

$$F(b_1, \dots, b_T, \alpha) = \sum_{t=1}^T \alpha_t b_t, \quad \sum_{t=1}^T \alpha_t = 1, \quad \alpha_t \geq 0$$

- смесь моделей-экспертов (mixture of experts)  
с функциями компетентности (gating function)  $g_t: X \rightarrow \mathbb{R}$

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

## Проблема разнообразия (diversity) базовых алгоритмов

Измерение с.в.  $\xi$  по независимым наблюдениям  $\{\xi_t\}$ :

- $E\frac{1}{T}(\xi_1 + \dots + \xi_T) = E\xi$  — матожидание среднего
- $D\frac{1}{T}(\xi_1 + \dots + \xi_T) = \frac{1}{T}D\xi$  — дисперсия  $\rightarrow 0$  при  $T \rightarrow \infty$

Но базовые алгоритмы не являются независимыми с.в.:

- решают одну и ту же задачу
- настраиваются на один целевой вектор ( $y_i$ )
- обычно выбираются из одной и той же модели

Способы повышения *разнообразия* базовых алгоритмов:

- обучение по различным (случайным) подвыборкам
- обучение по различным (случайным) наборам признаков
- обучение из разных параметрических моделей
- обучение с использованием рандомизации
- (иногда даже) обучение по зашумлённым данным

## Методы стохастического ансамблирования

Способы повышения разнообразия с помощью рандомизации:

- bagging (bootstrap aggregating) — подвыборки обучающей выборки «с возвращением», в каждую выборку попадает  $1 - (1 - \frac{1}{\ell})^\ell \rightarrow 1 - \frac{1}{e} \approx 63.2\%$  объектов, при  $\ell \rightarrow \infty$
- pasting — случайные обучающие подвыборки
- random subspaces — случайные подмножества признаков
- random patches — случ. подмн-ва и объектов, и признаков
- cross-validated committee — выборка разбивается на  $k$  блоков ( $k$ -fold) и делается  $k$  обучений без одного блока

Пусть  $\mu: (G, U) \mapsto b$  — метод обучения по подвыборке  $U \subseteq X^\ell$ , использующий только признаки из  $G \subseteq F^n = \{f_1, \dots, f_n\}$

$$Q(b, U) = \sum_{x_i \in U} \mathcal{L}(b(x_i, w), y_i) \rightarrow \min_w \text{ — миним. эмпирич. риска}$$

---

*Tin Kam Ho.* The random subspace method for constructing decision forests. 1998.  
*Leo Breiman.* Bagging predictors // Machine Learning. 1996.

## Методы стохастического ансамблирования в одном псевдо-коде

**Вход:** обучающая выборка  $X^\ell$ ; параметры:  $T$ ,  
 $\ell'$  — объём обучающих подвыборок,  
 $n'$  — размерность признаков подпространств,  
 $\varepsilon_1$  — порог качества базовых алгоритмов на обучении,  
 $\varepsilon_2$  — порог качества базовых алгоритмов на контроле;

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

для всех  $t = 1, \dots, T$ :

$U_t :=$  случайная подвыборка объёма  $\ell'$  из  $X^\ell$ ;

$G_t :=$  случайное подмножество мощности  $n'$  из  $F^n$ ;

$b_t := \mu(G_t, U_t)$ ;

если  $Q(b_t, U_t) > \varepsilon_1$  то не включать  $b_t$  в ансамбль;

если  $Q(b_t, X^\ell \setminus U_t) > \varepsilon_2$  то не включать  $b_t$  в ансамбль;

**Ансамбль** — простое голосование:  $b(x) = \frac{1}{T} \sum_{t=1}^T b_t(x)$

## Несмещённая оценка ошибок

*Out-of-bag* — несмещённая оценка ансамбля на объекте:

$$\text{OOB}(x_i) = \frac{1}{|T_i|} \sum_{t \in T_i} b_t(x_i), \quad T_i = \{t: x_i \notin U_t\}$$

Несмещённая оценка ошибки ансамбля на обучающей выборке:

$$\text{OOB}(X^\ell) = \sum_{i=1}^{\ell} \mathcal{L}(\text{OOB}(x_i), y_i),$$

где  $\mathcal{L}(b(x_i), y_i)$  — значение функции потерь на объекте  $x_i$ .

Оценивание важности признаков  $f_j$ ,  $j = 1, \dots, n$ :

$$\text{importance}_j = \frac{\text{OOB}^j(X^\ell) - \text{OOB}(X^\ell)}{\text{OOB}(X^\ell)} \cdot 100\%,$$

где при вычислении  $b_t(x_i)$  для  $\text{OOB}^j$  значения признака  $f_j$  случайным образом перемешиваются на всех объектах  $x_i \notin U_t$ .

## Случайный лес (random forest)

### Грубое обучение деревьев для случайного леса:

- бэггинг над решающими деревьями, без pruning
- признак в каждой вершине дерева выбирается из случайного подмножества  $k$  из  $n$  признаков. По умолчанию  $k = \lfloor n/3 \rfloor$  для регрессии,  $k = \lfloor \sqrt{n} \rfloor$  для классификации

### Параметры, которые можно настраивать (в частности, по ООВ):

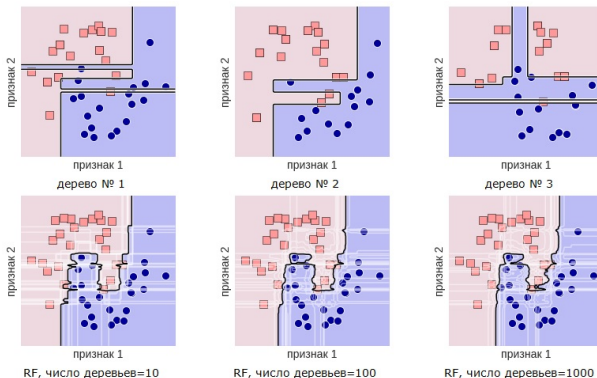
- число  $T$  деревьев
- число  $k$  случайно выбираемых признаков
- максимальная глубина деревьев
- минимальное число объектов в расщепляемой подвыборке
- минимальное число объектов в листьях
- критерий расщепления внутренних вершин дерева

---

Breiman L. Random Forests. Machine Learning, 2001.

## Постепенное сглаживание разделяющей поверхности

Пример разделения выборки с помощью отдельных деревьев (показаны соответствующие бутстреп-подвыборки) и случайного леса с числом деревьев 10, 100, 1000:



<https://dyakonov.org/2019/04/19/ансамбли-в-машинном-обучении>

## Преимущества и ограничения стохастического ансамблирования

### Преимущества:

- метод-обёртка (envelop) над базовым методом обучения
- подходит для классификации, регрессии и других задач
- простая реализация и простое распараллеливание
- возможность получения несмещённых оценок ООВ
- возможность оценивания важности признаков
- RF — один из лучших универсальных методов в ML

### Ограничения:

- требуется ооооочень много базовых алгоритмов
- трудно агрегировать устойчивые базовые методы обучения

---

<https://dyakonov.org/2016/11/14/случайный-лес-random-forest>

## Градиентный бустинг с произвольной функцией потерь

**Дано:**  $X^\ell = (x_i, y_i)_{i=1}^\ell \subset X \times Y$ ,  $b_t: X \rightarrow \mathbb{R}$

**Найти** линейный ансамбль базовых алгоритмов  $b_t$ :

$$a_T(x) = \sum_{t=1}^T \alpha_t b_t(x), \quad x \in X, \quad \alpha_t \geq 0$$

**Эвристика:** обучаем  $\alpha_T, b_T$  при фиксированных предыдущих  
**Критерий** с заданной гладкой функцией потерь  $\mathcal{L}(b, y)$ :

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

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

$(a_{T,i})_{i=1}^\ell$  — вектор текущего приближения

## Параметрическая аппроксимация градиентного шага

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

$a_{0,i} :=$  начальное приближение;

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

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

Это очень похоже на добавление одного базового алгоритма:

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

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

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

## Алгоритм градиентного бустинга (Gradient Boosting)

**Вход:** обучающая выборка  $X^\ell$ ; **параметр**  $T$ ;

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

инициализация:  $a_{0,i} := 0$ ,  $i = 1, \dots, \ell$ ;

**для всех**  $t = 1, \dots, T$

базовый алгоритм, приближающий антиградиент:

$$b_t := \arg \min_{b \in \mathcal{B}} \sum_{i=1}^{\ell} (b(x_i) + \mathcal{L}'(a_{t-1,i}, y_i))^2;$$

задача одномерной минимизации:

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

обновление вектора значений на объектах выборки:

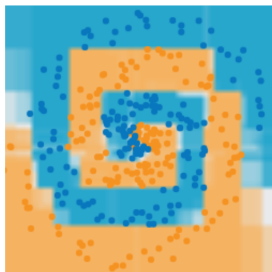
$$a_{t,i} := a_{t-1,i} + \alpha_t b_t(x_i); \quad i = 1, \dots, \ell;$$

Каждый следующий базовый алгоритм обучается так, чтобы по возможности исправить ошибки предыдущих алгоритмов.

## Пример. Классификация синтетической выборки

100 деревьев глубины 5

Prediction:

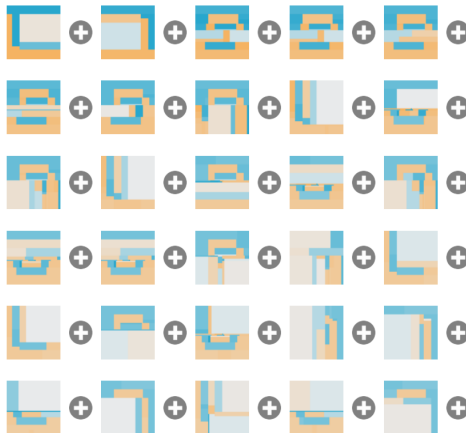


↑  
predictions of GB (all 100 trees)

train loss: 0.022    test loss: 0.218



Decision functions of first 30 trees

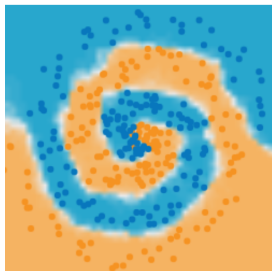


[http://arogozhnikov.github.io/2016/07/05/gradient\\_boosting\\_playground.html](http://arogozhnikov.github.io/2016/07/05/gradient_boosting_playground.html)

## Пример. Классификация синтетической выборки

100 деревьев глубины 5, с подбором вращения каждого дерева

Prediction:

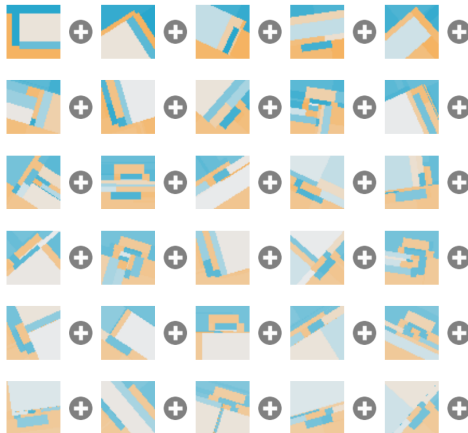


↑  
predictions of GB (all 100 trees)

train loss: 0.013    test loss: 0.092



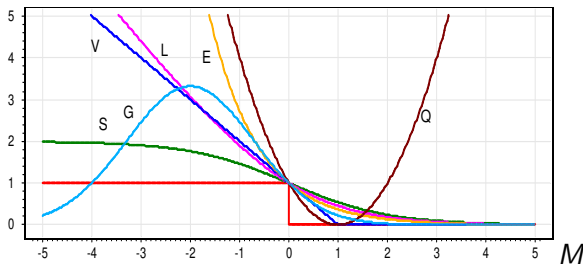
Decision functions of first 30 trees



[http://arogozhnikov.github.io/2016/07/05/gradient\\_boosting\\_playground.html](http://arogozhnikov.github.io/2016/07/05/gradient_boosting_playground.html)

## Варианты бустинга для двухклассовой классификации

Гладкие аппроксимации пороговой функции потерь [ $M < 0$ ]:



$E(M) = e^{-M}$  — экспоненциальная (AdaBoost);

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

$Q(M) = (1 - M)^2$  — квадратичная (GentleBoost);

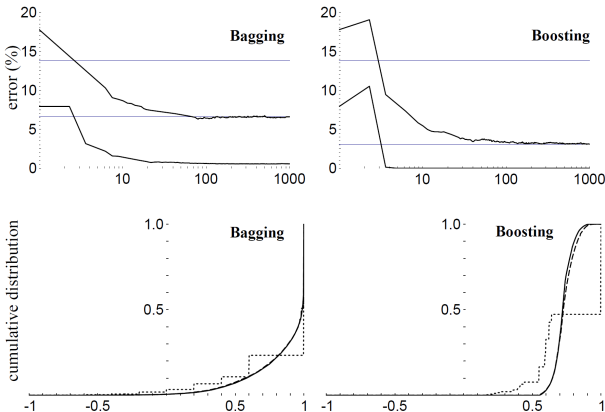
$G(M) = \exp(-cM(M + s))$  — гауссовская (BrownBoost);

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

$V(M) = (1 - M)_+$  — кусочно-линейная (из SVM);

## Удивительно: линейные ансамбли почти не переобучаются!

Ошибки на обучении и тесте. Снизу распределение отступов.



*R.E.Schapire, Y.Freund, Wee Sun Lee, P.Bartlett. Boosting the margin: a new explanation for the effectiveness of voting methods. Annals of Statistics, 1998.*

## Обоснование бустинга (случай бинарной классификации)

Усиленная частота ошибок классификатора  $\text{sign } b(x)$ ,  $b \in \mathcal{B}$ :

$$\nu_\theta(b, X^\ell) = \frac{1}{\ell} \sum_{i=1}^{\ell} [b(x_i)y_i \leq \theta], \quad \theta > 0.$$

Обычная частота ошибок  $\nu_0(b, X^\ell) \leq \nu_\theta(b, X^\ell)$  при  $\theta > 0$ .

### Теорема (Freund, Schapire, Lee, Bartlett, 1998)

Если  $|\mathcal{B}| < \infty$ , то  $\forall \theta > 0$ ,  $\forall \eta \in (0, 1)$  с вероятностью  $1 - \eta$

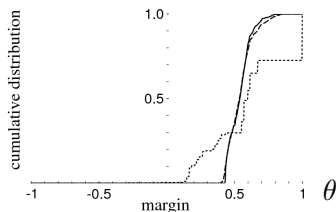
$$P[y a(x) < 0] \leq \nu_\theta(a, X^\ell) + C \sqrt{\frac{\ln |\mathcal{B}| \ln \ell}{\ell \theta^2} + \frac{1}{\ell} \ln \frac{1}{\eta}}$$

**Основной вывод:** оценка зависит от  $|\mathcal{B}|$ , но не от  $T$ .  
Голосование не увеличивает сложность базовой модели,  
а лишь усредняет ответы базовых алгоритмов.

**Пример.** Для семейства решающих пней  $|\mathcal{B}| \leq \ell n$

## Обоснование бустинга: что же всё-таки происходит?

**Распределение отступов:**  
доля объектов, имеющих отступ меньше заданного  $\theta$  после 5, 100, 1000 итераций (задача UCI:vehicle)



- С ростом  $T$  распределение отступов сдвигается вправо, то есть бустинг «раздвигает» классы в пространстве векторов растущей размерности  $(b_1(x), \dots, b_T(x))$
- Значит, в оценке можно уменьшать второй член, увеличивая  $\theta$  при неизменной  $\nu_\theta(a, X^\ell) = \nu_0(a, X^\ell)$
- Можно уменьшить второй член, если уменьшить  $|\mathcal{B}|$ , то есть взять простую модель базовых алгоритмов

*R.E.Schapire, Y.Freund, Wee Sun Lee, P.Bartlett.* Boosting the margin: a new explanation for the effectiveness of voting methods. *Annals of Statistics*, 1998.

## Анализ смещения–разброса (bias–variance)

Задача регрессии:  $Y = \mathbb{R}$

Квадратичная функция потерь:  $\mathcal{L}(a, y) = (a(x) - y)^2$

Вероятностная постановка:  $X^\ell = (x_i, y_i)_{i=1}^\ell \sim p(x, y)$

Метод обучения:  $\mu: 2^X \rightarrow A$ , т.е. выборка  $\mapsto$  алгоритм

Задача минимизации среднеквадратичного риска:

$$R(a) = E_{x,y}(a(x) - y)^2 = \int_X \int_Y (a(x) - y)^2 p(x, y) dx dy \rightarrow \min_a$$

Идеальный минимизатор среднеквадратичного риска:

$$a^*(x) = E(y|x) = \int_Y y p(y|x) dy$$

Основная мера качества метода обучения  $\mu$ :

$$Q(\mu) = E_{X^\ell} R(\mu(X^\ell)) = E_{X^\ell} E_{x,y} (\mu(X^\ell)(x) - y)^2$$

## Разложение ошибки на шум, смещение и разброс

$a^*(x) = E(y|x)$  — неизвестная идеальная зависимость  $y$  от  $x$

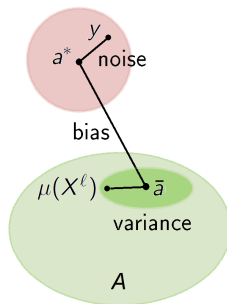
$y(x) \sim p(y|x)$  — наблюдаемый ответ на объекте  $x$

$a = \mu(X^\ell)$  — аппроксимация, выбранная по  $X^\ell$  из семейства  $A$

$\bar{a}(x) = E_{X^\ell}(a(x))$  — средний ответ обученного алгоритма

**Теорема.** При квадратичной функции потерь для любого  $\mu$

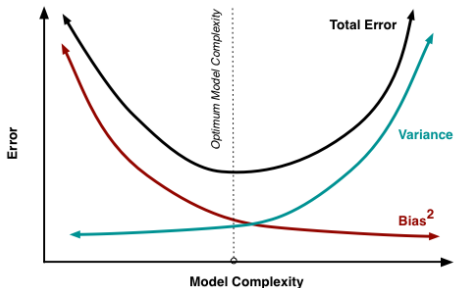
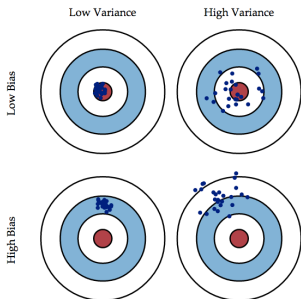
$$\begin{aligned}
 Q(\mu) = & \underbrace{E_{x,y} (a^*(x) - y)^2}_{\text{шум (noise)}} + \\
 & + \underbrace{E_{x,y} (\bar{a}(x) - a^*(x))^2}_{\text{смещение (bias)}} + \\
 & + \underbrace{E_{x,y} E_{X^\ell} (\mu(X^\ell)(x) - \bar{a}(x))^2}_{\text{разброс (variance)}}
 \end{aligned}$$



## Разложение ошибки на шум, смещение и разброс

Качественное понимание: по мере роста сложности модели

- смещение (bias) уменьшается
- разброс (variance) увеличивается



## Анализ смещения–разброса для простого голосования

Обучение базовых алгоритмов по случайным подвыборкам:

$$b_t = \mu(X_t^k), \quad X_t^k \sim X^\ell, \quad t = 1, \dots, T$$

Ансамбль — простое голосование:  $a_T(x) = \frac{1}{T} \sum_{t=1}^T b_t(x)$

**Смещение** ансамбля совпадает со смещением отдельного базового алгоритма:

$$\text{bias} = E_{x,y} (a^*(x) - E_{X^\ell} b_t(x))^2$$

**Разброс** состоит из дисперсии и различности (ковариации):

$$\begin{aligned} \text{variance} = & \frac{1}{T} E_{x,y} E_{X^\ell} (b_t(x) - E_{X^\ell} b_t(x))^2 + \\ & + \frac{T-1}{T} E_{x,y} E_{X^\ell} (b_t(x) - E_{X^\ell} b_t(x)) (b_s(x) - E_{X^\ell} b_s(x)) \end{aligned}$$

## Почему сложные ансамбли не переобучаются?

### С позиций анализа отступов:

- ансамблирование не увеличивает сложность модели
- но с каждой итерацией увеличивает зазор между классами

### С позиций анализа смещения–разброса:

- разнообразие базовых алгоритмов уменьшает разброс
- бэггинг уменьшает только разброс
- бустинг уменьшает и смещение, и разброс

### Практическое сравнение: boosting / bagging / RSM

- бустинг лучше для классов с границами сложной формы
- бэггинг и RSM лучше для коротких обучающих выборок
- RSM лучше, когда много неинформативных признаков
- бэггинг параллельно обучает базовые алгоритмы  $b_t$
- бустинг обучает каждый  $b_t$  параллельно по частям выборки

## Недостатки бэггинга и бустинга

- задача минимизировать число  $T$  вообще не ставится
- композиция из сотен алгоритмов не интерпретируема
- не удаётся строить короткие композиции из «сильных» алгоритмов типа SVM (только длинные из «слабых»)

### Несколько эмпирических наблюдений:

- веса алгоритмов не важны для оптимизации отступов
- веса объектов не важны для обеспечения различности

### Предлагается:

- отказаться от аппроксимации пороговой функции потерь,
- оптимизировать распределение отступов композиции,
- обучать базовые алгоритмы последовательно (как бустинг),
- обучать их на подвыборках (как бэггинг), но не случайных,
- использовать простое голосование (комитет большинства)

## Оптимизация распределения отступов на каждом шаге

**Идея:** явно управлять распределением отступов, максимизируя различность базовых алгоритмов и минимизируя их число.

**Дано:** задача бинарной классификации,  $X^\ell$ ,  $Y = \{\pm 1\}$

**Найти:** ансамбль  $b(x) = \frac{1}{T} \sum_{t=1}^T b_t(x)$ ,  $C(b) = \text{sign}(b)$ .

**Критерий** — минимум числа ошибок ансамбля на обучении:

$$Q(a, X^\ell) = \sum_{i=1}^{\ell} [y_i a(x_i) < 0] = \sum_{i=1}^{\ell} \underbrace{[y_i b_1(x_i) + \dots + y_i b_T(x_i)]}_{M_{iT}} < 0],$$

$M_{it} = y_i b_1(x_i) + \dots + y_i b_t(x_i)$  — отступ (margin) объекта  $x_i$ .

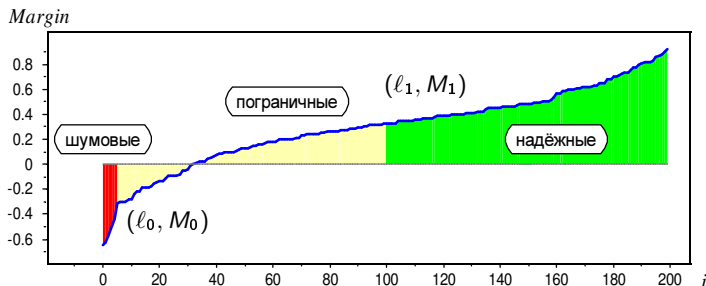
**Эвристика:**  $b_t$  компенсирует ошибки ансамбля,

$$Q(b_t, U_t) = \sum_{x_i \in U_t} [y_i b_t(x_i) < 0] \rightarrow \min_{b_t}$$

$U_t = \{x_i : M_0 < M_{i,t-1} \leq M_1\}$ ,  $M_0, M_1$  — параметры метода

## Формирование выборки для обучения базового алгоритма

Упорядочим объекты по возрастанию отступов  $M_{i,t-1}$ :



**Принцип выравнивания распределения отступов**

два случая, когда  $b_t$  на объекте  $x_i$  обучать не надо:

$M_{i,t-1} < M_0$ ,  $i < l_0$  — объект  $x_i$  шумовой

$M_{i,t-1} > M_1$ ,  $i > l_1$  — объект  $x_i$  надёжный

## Алгоритм ComBoost (Committee Boosting)

**Вход:** выборки  $X^\ell, X^k$ ; **параметры**  $T, \ell_0, \ell_1, \ell_2, \Delta\ell$ ;

**Выход:**  $b_1, \dots, b_T$ ;

$b_1 := \arg \min_b Q(b, X^\ell)$ ; отступы  $M_i = y_i b_1(x_i)$ ,  $i = 1, \dots, \ell$ ;

**для всех**  $t = 2, \dots, T$ :

упорядочить выборку  $X^\ell$  по возрастанию отступов  $M_i$ ;

**для всех**  $\ell' = \ell_1, \dots, \ell_2$  с шагом  $\Delta\ell$ :

$U_t = \{x_i \in X^\ell : \ell_0 \leq i \leq \ell'\}$ ;

$b_{t\ell'} := \arg \min_b Q(b, U_t)$  — инкрементное обучение;

выбрать наилучший  $b_t \in \{b_{t\ell'}\}$  по критерию  $Q(a, X^k)$ ;

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

**пока**  $Q$  существенно улучшается.

Маценов А. А. Комитетный бустинг: минимизация числа базовых алгоритмов при простом голосовании. ММРО-13, 2007.

## Результаты эксперимента на 4 задачах из репозитория UCI

По 50 случайным разбиениям «обучение : контроль» = 4 : 1

	ionosphere	pima	bupa	votes
SVM	12,9	24,2	42,0	4,6
ComBoost <sub>0</sub> [SVM] (T)	12,6 (4)	23,1 (2)	34,2 (5)	4,0 (2)
ComBoost [SVM] (T)	<b>12,3 (5)</b>	<b>22,5 (2)</b>	30,9 (5)	<b>3,8 (3)</b>
AdaBoost [SVM] (T)	15,0 (65)	22,7 (18)	<b>30,6 (15)</b>	4,0 (8)
Parzen	6,3	25,1	41,6	6,9
ComBoost <sub>0</sub> [Parzen]	6,1	25,0	38,1	6,8
ComBoost [Parzen]	<b>5,8</b>	<b>24,7</b>	30,6	<b>6,2</b>
AdaBoost [Parzen]	6,0	24,8	<b>30,5</b>	6,5

ComBoost<sub>0</sub> — без подбора длины подвыборки  $U_t$  в цикле  $l' = l_1, \dots, l_2$

Parzen — метод окна Парзена с подбором ширины окна по leave-one-out

**Результат:** ComBoost способен строить короткие ансамбли из сильных и устойчивых базовых алгоритмов

Маценов А. А. Комитетный бустинг: минимизация числа базовых алгоритмов при простом голосовании. ММРО-13, 2007.

## Обобщение для задач с произвольным числом классов

$Y = \{1, \dots, M\}$ , ансамбль — простое голосование, причём каждый базовый алгоритм  $b_{yt}$  голосует только за свой класс  $y$ :

$$a(x) = \arg \max_{y \in Y} \Gamma_y(x); \quad \Gamma_y(x) = \frac{1}{|T_y|} \sum_{t \in T_y} b_{yt}(x).$$

В алгоритме ComBoost три небольших изменения:

- обобщённое определение отступа  $M_i$ :

$$M_i = \Gamma_{y_i}(x_i) - \max_{y \in Y \setminus \{y_i\}} \Gamma_y(x_i).$$

- придётся решать, для какого класса строить очередной  $b_{yt}$  (например, для того  $y$ , на котором доля ошибок больше)
- изменится пересчёт отступов в конце итерации

---

Allwein E. L., Schapire R. E., Singer Y. Reducing multiclass to binary: A unifying approach for margin classifiers. 2000

## Смесь алгоритмов (Mixture of Experts)

$b_t: X \rightarrow \mathbb{R}$  — базовые алгоритмы,  $t = 1, \dots, T$

$g_t: X \rightarrow \mathbb{R}$  — функция компетентности (шлюз, gate) для  $b_t(x)$

$$b(x) = \sum_{t=1}^T g_t(x)b_t(x)$$

Чем больше  $g_t(x)$ , тем выше доверие к ответу  $b_t(x)$ .

Условие нормировки:  $\sum_{t=1}^T g_t(x) = 1$  для любого  $x \in X$ .

Нормировка «мягкого максимума» SoftMax:  $\mathbb{R}^T \rightarrow \mathbb{R}^T$ :

$$\tilde{g}_t(x) = \text{SoftMax}_t(g_1(x), \dots, g_T(x); \gamma) = \frac{e^{\gamma g_t(x)}}{e^{\gamma g_1(x)} + \dots + e^{\gamma g_T(x)}}$$

При  $\gamma \rightarrow \infty$  SoftMax выделяет максимальную из  $T$  величин.

---

*Растригин Л. А., Эренштейн Р. Х.* Коллективные правила распознавания. 1981.

*Hien D. Nguyen, Faïcel Chamroukhi.* Practical and theoretical aspects of mixture-of-experts modeling: An overview. 2018

## Вид функций компетентности

Функции компетентности определяются из практических соображений, в зависимости от особенностей задачи, например:

- по признаку  $f(x)$ :

$$g(x; \alpha, \beta) = \sigma(\alpha f(x) + \beta), \quad \alpha, \beta \in \mathbb{R};$$

- по неизвестному направлению  $\alpha \in \mathbb{R}^n$ :

$$g(x; \alpha, \beta) = \sigma(x^T \alpha + \beta), \quad \alpha \in \mathbb{R}^n, \beta \in \mathbb{R};$$

- по расстоянию до неизвестной точки  $\alpha \in \mathbb{R}^n$ :

$$g(x; \alpha, \beta) = \exp(-\beta \|x - \alpha\|^2), \quad \alpha \in \mathbb{R}^n, \beta \in \mathbb{R};$$

где параметры  $\alpha, \beta$  фиксируются или обучаются по выборке,

$\sigma(z) = \frac{1}{1+e^{-z}}$  — сигмоидная функция.

## Выпуклые функции потерь

Функция потерь  $\mathcal{L}(b, y)$  называется *выпуклой* по  $b$ , если  $\forall y \in Y, \forall b_1, b_2 \in R, \forall g_1, g_2 \geq 0: g_1 + g_2 = 1$ , выполняется

$$\mathcal{L}(g_1 b_1 + g_2 b_2, y) \leq g_1 \mathcal{L}(b_1, y) + g_2 \mathcal{L}(b_2, y).$$

**Интерпретация:** потери растут не медленнее, чем величина отклонения от правильного ответа  $y$ .

**Примеры** выпуклых функций потерь:

$$\mathcal{L}(b, y) = \begin{cases} (b - y)^2 & \text{— квадратичная (МНК-регрессия);} \\ e^{-by} & \text{— экспоненциальная (AdaBoost);} \\ \log_2(1 + e^{-by}) & \text{— логарифмическая (LR);} \\ (1 - by)_+ & \text{— кусочно-линейная (SVM).} \end{cases}$$

**Пример** невыпуклой функции потерь:  $\mathcal{L}(b, y) = [by < 0]$ .

## Основная идея применения выпуклых функций потерь

Пусть  $\forall x \sum_{t=1}^T g_t(x) = 1$  и функция потерь  $\mathcal{L}$  выпукла.

Тогда  $Q(a)$  распадается на  $T$  независимых критериев  $Q_t$ :

$$Q(a) = \sum_{i=1}^{\ell} \mathcal{L} \left( \sum_{t=1}^T g_t(x_i) b_t(x_i), y_i \right) \leq \sum_{t=1}^T \underbrace{\sum_{i=1}^{\ell} g_t(x_i) \mathcal{L}(b_t(x_i), y_i)}_{Q_t(g_t, b_t)}$$

Итерационный процесс, два шага на каждой итерации:

начальное приближение функций компетентности  $g_t$ ;

**повторять**

- обучить все  $b_t := \arg \min_b Q_t(g_t, b)$  при фиксированных  $g_t$ ;
- обучить все  $g_t$  при фиксированных  $b_t$ ;

**пока** значения компетентностей  $g_t(x_i)$  не стабилизируются;

## Философия ансамблирования

Ансамблировать можно только нечто гомогенное.

- 1 **Декомпозиция** — разделение модели алгоритма  $a_t$  на алгоритмический оператор  $b_t$  и решающее правило  $C$ :

$$a_t = C \circ b_t$$

- 2 **Гомогенизация** — разнородные модели имеют общее пространство оценок  $R$  и общую структуру алгоритмического оператора  $b_t$  как отображения

$$b_t: X \rightarrow R$$

- 3 **Ансамблирование** — совместное обучение базовых алгоритмических операторов для решения общей задачи:

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

---

Ю.И. Журавлёв. Об алгебраическом подходе к решению задач распознавания или классификации. Проблемы кибернетики, 1978.

## Философия многозадачного обучения

- ① **Декомпозиция** — разделение моделей  $y_t: X \rightarrow Y_t$  на векторизатор  $z = f(x, \alpha)$  и предиктор  $y_t = g_t(z, \beta)$ :

$$y_t(x) = g_t(f(x, \alpha), \beta_t)$$

- ② **Гомогенизация** — разнородные модели имеют общий векторизатор  $z = f(x, \alpha)$  и общее векторное пространство представлений (эмбедингов)  $Z$ :

$$f: X \rightarrow Z$$

- ③ **Ансамблирование** — совместное обучение общего векторизатора для решения разнородных задач:

$$\sum_{t \in T} \sum_{i \in X^t} \mathcal{L}_{ti}(g_t(f(x_{ti}, \alpha), \beta_t)) \rightarrow \min_{\alpha, \{\beta_t\}}$$

---

Yu Zhang, Qiang Yang. A survey on multi-task learning. 2021

M. Crawshaw. Multi-task learning with deep neural networks: a survey. 2020

## Философия фундаментальных моделей

- ① **Декомпозиция** — разделение моделей  $y_t: X_t \rightarrow Y_t$  на векторизатор  $z = f(x_t, \alpha_t)$  и предиктор  $y_t = g(z_t, \beta_t)$ :

$$y_t(x) = g_t(f(x_t, \alpha_t), \beta_t)$$

- ② **Гомогенизация** — разнородные модели в разнородных задачах имеют общее пространство эмбедингов  $Z$ :

$$f_t: X_t \rightarrow Z$$

- ③ **Ансамблирование** — совместное обучение эмбедингов в едином семантическом пространстве для решения разнородных задач:

$$\sum_{t \in T} \sum_{i \in X^t} \mathcal{L}_{ti}(g_t(f(x_{ti}, \alpha_t), \beta_t)) \rightarrow \min_{\{\alpha_t, \beta_t\}}$$

---

*R. Bommasani et al. (Center for Research on Foundation Models, Stanford University)*  
On the opportunities and risks of foundation models // CoRR, 20 August 2021.

## Философия аддитивной регуляризации (ARTM)

- 1 **Декомпозиция** — разделение критерия обучения модели на основной (log-правдоподобие) и регуляризатор  $R$ :

$$\sum_i \ln p(x_i | \Phi, \Theta) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

- 2 **Гомогенизация** — разнородные модели имеют общую структуру векторизатора (матричное разложение) и общий основной критерий (log-правдоподобие):

$$f_{\Phi}: X \rightarrow \Theta, \quad \sum_i \ln p(x_i | \Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

- 3 **Ансамблирование** — совместное использование регуляризаторов  $R_k$ , взятых от разнородных моделей:

$$\sum_i \ln p(x_i | \Phi, \Theta) + \sum_k \lambda_k R_k(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

---

Vorontsov K. V. Rethinking probabilistic topic modeling from the point of view of classical non-Bayesian regularization. 2023.

- Ансамбли позволяют решать сложные задачи, которые плохо решаются отдельными базовыми алгоритмами
- Обычно ансамбль строится *алгоритмом-обёрткой* (envelop): базовые алгоритмы обучаются готовыми методами
- Важное открытие середины 90-х: обобщающая способность бустинга не ухудшается с ростом сложности  $T$
- Градиентный бустинг — наиболее общий из всех бустингов:
  - произвольная функция потерь
  - произвольное пространство оценок  $R$
  - подходит для регрессии, классификации, ранжирования
- Базовые алгоритмы: компромисс качество/различность
- Чаще всего GB применяется к решающим деревьям
- Для смешивания нужна адекватная модель компетентности