

# Композиции классификаторов (часть 2)

К. В. Воронцов  
vokov@forecsys.ru

Этот курс доступен на странице вики-ресурса  
<http://www.MachineLearning.ru/wiki>  
«Машинное обучение (курс лекций, К.В.Воронцов)»

ШАД Яндекс • 1 октября 2019

## 1 Простое голосование

- Бэггинг и метод случайных подпространств
- Комитетный бустинг
- Случайные леса

## 2 Разложение ошибки на смещение и разброс

- Теория: общая формула разложения
- Частные случаи
- Композиции

## 3 Смеси алгоритмов

- Идея областей компетентности
- Итерационный метод обучения смеси
- Последовательное наращивание смеси

## Определение композиции (напоминание)

$X^\ell = (x_i, y_i)_{i=1}^\ell \subset X \times Y$  — обучающая выборка,  $y_i = y^*(x_i)$ ;

$a(x) = C(b(x))$  — алгоритм, где

$b: X \rightarrow R$  — базовый алгоритм (алгоритмический оператор),

$C: R \rightarrow Y$  — решающее правило,

$R$  — пространство оценок;

### Определение

Композиция базовых алгоритмов  $b_1, \dots, b_T$

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

где  $F: R^T \rightarrow R$  — корректирующая операция.

Зачем вводится  $R$ ?

В задачах классификации множество отображений

$\{F: R^T \rightarrow R\}$  существенно шире, чем  $\{F: Y^T \rightarrow Y\}$ .

## Примеры корректирующих операций (напоминание)

- **Пример 1:** Простое голосование (Simple Voting):

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

- **Пример 2:** Взвешенное голосование (Weighted Voting):

$$F(b_1(x), \dots, b_T(x)) = \sum_{t=1}^T \alpha_t b_t(x), \quad x \in X, \quad \alpha_t \in \mathbb{R}.$$

- **Пример 3:** Смесь алгоритмов (Mixture of Experts)

$$F(b_1(x), \dots, b_T(x)) = \sum_{t=1}^T g_t(x) b_t(x), \quad x \in X, \quad g_t: X \rightarrow \mathbb{R}.$$

## Стохастические методы построения композиций

Чтобы алгоритмы в композиции были различными

- их обучают по (случайным) подвыборкам,
- либо по (случайным) подмножествам признаков.

**Первую идею** реализует bagging (bootstrap aggregation) [Breiman, 1996]: подвыборки длины  $\ell$  с повторениями, доля объектов, попадающих в выборку:  $(1 - \frac{1}{e}) \approx 0.632$

**Вторую идею** реализует RSM (random subspace method) [Ho, 1998].

**Совместим обе идеи в одном алгоритме.**

$\mathcal{F} = \{f_1, \dots, f_n\}$  — признаки,

$\mu(\mathcal{G}, U)$  — метод обучения алгоритма по подвыборке  $U \subseteq X^\ell$ , использующий только признаки из  $\mathcal{G} \subseteq \mathcal{F}$ .

## Бэггинг и метод случайных подпространств

- Вход:** обучающая выборка  $X^\ell$ ; параметры:  $T$ ;  
 $\ell'$  — длина обучающих подвыборок;  
 $n'$  — длина признакового подописания;  
 $\varepsilon_1$  — порог качества базовых алгоритмов на обучении;  
 $\varepsilon_2$  — порог качества базовых алгоритмов на контроле;  
**Выход:** базовые алгоритмы  $b_t$ ,  $t = 1, \dots, T$ ;

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

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

## Сравнение: boosting — bagging — RSM

- 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

### И ещё несколько эмпирических наблюдений:

- Веса алгоритмов не столь важны для выравнивания отступов
- Веса объектов не столь важны для обеспечения различности
- Короткие композиции из «сильных» алгоритмов типа SVM строить труднее, чем длинные из слабых

## Оптимизация распределения отступов на каждом шаге

Возьмём  $Y = \{\pm 1\}$ ,  $F(b_1, \dots, b_T) = \frac{1}{T} \sum_{t=1}^T b_t$ ,  $C(b) = \text{sign}(b)$ .

Функционал качества композиции — число ошибок на обучении:

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

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

**Эвристика:** чтобы  $b_t$  компенсировал ошибки композиции,

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

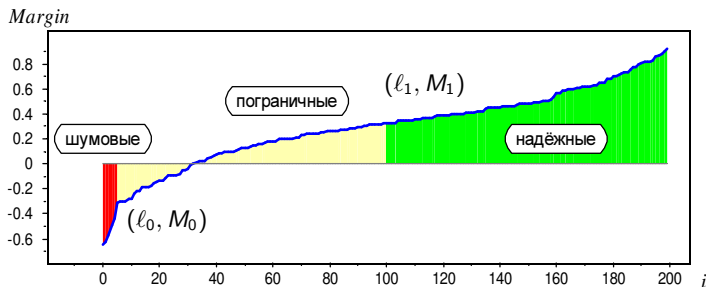
где  $U_t = \{x_i : M_0 < M_{i,t-1} \leq M_1\}$ ,

$M_0, M_1$  — параметры метода обучения.



## Подбор параметров $M_0$ и $M_1$

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



**Принцип максимизации и выравнивания 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  $x_i$  надёжно классифицируется.

## Алгоритм ComBoost (Committee Boosting)

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

**Выход:**  $b_1, \dots, b_T$

- 1:  $b_1 := \arg \min_b Q(b, X^\ell)$ ;  
упорядочить  $X^\ell$  по возрастанию  $M_i = y_i b_1(x_i)$ ,  $i = 1, \dots, \ell$ ;
- 2: **для всех**  $t = 1, \dots, T$
- 3:   **для всех**  $k = \ell_1, \dots, \ell_2$  с шагом  $\Delta\ell$
- 4:      $U_t = \{x_i \in X^\ell : \ell_0 \leq i \leq k\}$ ;
- 5:      $b_{tk} := \arg \min_b Q(b, U_t)$ ;
- 6:   выбрать наилучший  $b_t \in \{b_{tk}\}$  по критерию  $Q$ ;
- 7:   обновить отступы:  $M_i := M_i + y_i b_t(x_i)$ ,  $i = 1, \dots, \ell$ ;
- 8:   упорядочить выборку  $X^\ell$  по возрастанию отступов  $M_i$ ;
- 9:   опция: скорректировать значения параметров  $\ell_0, \ell_1, \Delta\ell$ ;
- 10: **пока**  $Q$  существенно улучшается.

## Результаты эксперимента на 4 задачах из репозитория UCI

Средняя частота ошибок на контроле по 50 случайным разбиениям в отношении «обучение : контроль» = 4 : 1.

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

ComBoost<sub>0</sub> — с подбором  $l_0$  и  $l_1 = l_2$  по скользящему контролю;

ComBoost — с подбором длины подвыборки  $U_t$ ;

Parzen — окно Парзена с евклидовой метрикой и подбором ширины окна скользящим контролем LOO.

## Результаты эксперимента на 4 задачах из репозитория UCI

### Мощность композиций:

Число базовых алгоритмов	ionosphere	pima	bupa	votes
ComBoost <sub>0</sub> над SVM	4	2	5	2
ComBoost над SVM	5	2	5	3
AdaBoost над SVM	65	18	15	8

**Критерий останова:** отсутствие существенного улучшения качества классификации обучающей выборки.

---

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

## Обобщение для задач с произвольным числом классов

Пусть теперь  $Y = \{1, \dots, M\}$ .

Композиция — 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 только два изменения:

— изменится определение отступа  $M_i$ :

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

— в алгоритме ComBoost на шаге 3 придётся решать, за какой класс строить очередной базовый алгоритм, кроме того, немного изменится шаг 7 (пересчёт отступов).

## Преобразование простого голосования во взвешенное

Линейный классификатор над признаками  $b_t(x)$ :

$$a(x) = \text{sign} \sum_{t=1}^T \alpha_t b_t(x),$$

1. Метод обучения: SVM, логистическая регрессия, и т.п.:

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

2. Регуляризация:  $\alpha_t \geq 0$  либо LASSO:  $\sum_{t=1}^T |\alpha_t| \leq \varkappa$ .

3. Наивный байесовский классификатор

приводит к простому аналитическому решению:

$$\alpha_t = \ln \frac{1 - p_t}{p_t}, \quad t = 1, \dots, T,$$

где  $p_t$  — оценка вероятности ошибки базового алгоритма  $b_t$ .

## Случайный лес (Random Forest)

### Обучение случайного леса:

- бэггинг над решающими деревьями, без pruning
- признак в каждой вершине дерева выбирается из случайного подмножества  $k$  из  $n$  признаков
- регрессия:  $k = \lfloor n/3 \rfloor$ ; классификация:  $k = \lfloor \sqrt{n} \rfloor$

### Подбор числа деревьев $T$ по критерию *out-of-bag*:

число ошибок на объектах  $x_i$ , если не учитывать голоса деревьев, для которых  $x_i$  был обучающим:

$$\text{out-of-bag}(a) = \sum_{i=1}^{\ell} \left[ \text{sign} \left( \sum_{t=1}^T [x_i \notin U_t] b_t(x_i) \right) \neq y_i \right] \rightarrow \min$$

Это несмещённая оценка обобщающей способности.

*Breiman L.* Random Forests. Machine Learning, 2001.

## Задача минимизации среднеквадратичного риска

Задача регрессии:  $Y = \mathbb{R}$

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

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

Среднеквадратичный риск:

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

Идеальный минимизатор среднеквадратичного риска:

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

Основная мера качества метода обучения  $\mu$ :

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



# Разложение ошибки на шум, вариацию и смещение

## Теорема

В случае квадратичной функции потерь для любого  $\mu$

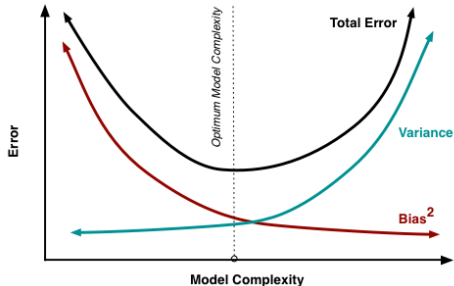
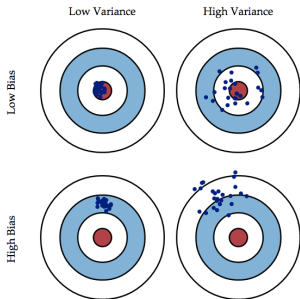
$$\begin{aligned}
 Q(\mu) = & \underbrace{E_{x,y} (a^*(x) - y)^2}_{\text{шум (noise)}} + \\
 & + \underbrace{E_{x,y} (\bar{a}(x) - a^*(x))^2}_{\text{смещение (bias)}} + \\
 & + \underbrace{E_{x,y} E_{X^\ell} (\mu(X^\ell)(x) - \bar{a}(x))^2}_{\text{разброс (variance)}},
 \end{aligned}$$

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

# Разложение ошибки на шум, вариацию и смещение

Качественное понимание: по мере роста сложности модели

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



## Упражнение. Метод $k$ ближайших соседей

Вероятностная модель данных:  $p(y|x) = f(x) + \mathcal{N}(0, \sigma^2)$

Метод  $k$  ближайших соседей:

$$a(x) = \frac{1}{k} \sum_{j=1}^k y(x^{(j)}),$$

где  $x^{(j)}$  —  $j$ -й сосед объекта  $x$

$a^*(x) = f(x)$  — истинная зависимость

$\bar{a}(x) = E_{x^\ell} \frac{1}{k} \sum_{j=1}^k f(x^{(j)})$  — средний ответ

Разложение bias-variance (чем меньше  $k$ , тем сложнее модель):

$$Q(\mu) = \underbrace{\sigma^2}_{\text{шум}} + \underbrace{E_{x,y} \left( \bar{a}(x) - f(x) \right)^2}_{\text{смещение}} + \underbrace{\frac{1}{k} \sigma^2}_{\text{разброс}}$$

## Простое голосование

Обучение базовых алгоритмов по случайным подвыборкам:

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

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

Смещение композиции совпадает со смещением отдельного базового алгоритма:

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

Разброс состоит из дисперсии и ковариации:

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

## Резюме. Почему бустинг и бэггинг работают?

- бэггинг уменьшает разброс
- бустинг уменьшает и смещение, и разброс
- корреляция базовых алгоритмов увеличивает разброс

Случайный лес (Random Forest) — разновидность бэггинга:

- RF уменьшает корреляции базовых алгоритмов
- RF — один из самых сильных методов машинного обучения.
- RF обычно лишь немного уступает градиентному бустингу.
- Бэггинг позволяет вычислять оценки out-of-bag:
  - для оптимизации числа базовых алгоритмов  $T$ ,
  - для оценивания важности признаков,
  - для выбора параметра  $k$  в RF.

## Квазилинейная композиция (смесь алгоритмов)

Смесь алгоритмов (Mixture of Experts)

$$a(x) = C \left( \sum_{t=1}^T g_t(x) b_t(x) \right),$$

$b_t: X \rightarrow \mathbb{R}$  — базовый алгоритм,

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

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.

## Вид функций компетентности

Функции компетентности выбираются из содержательных соображений и могут определяться:

- признаком  $f(x)$ :

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

- неизвестным направлением  $\alpha \in \mathbb{R}^n$ :

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

- расстоянием 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  $\alpha, \beta \in \mathbb{R}$  — параметры, *частично* обучаемые 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 функционалов  $Q_t$ :

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

Итерационный процесс, аналогичный EM-алгоритму:

- 1: начальное приближение функций компетентности  $g_t$ ;
- 2: **повторять**
- 3: **M-шаг**: при фиксированных  $g_t$  обучить все  $b_t$ ;
- 4: **E-шаг**: при фиксированных  $b_t$  оценить все  $g_t$ ;
- 5: **пока** значения компетентностей  $g_t(x_i)$  не стабилизируются.

## Алгоритм ME (Mixture of Experts): обучение смеси алгоритмов

Итерационный процесс, аналогичный EM-алгоритму:

**Вход:** выборка  $X^\ell$ , начальные  $(g_t)_{t=1}^T$ , параметры  $T, \delta, \gamma$ ;

**Выход:**  $g_t(x), b_t(x), t = 1, \dots, T$ ;

1: **повторять**

2:  $(\tilde{g}_1(x_i), \dots, \tilde{g}_T(x_i)) := \text{SoftMax}(g_1(x_i), \dots, g_T(x_i); \gamma)$ ;

3:  $\tilde{g}_t^0 := \tilde{g}_t$  для всех  $t = 1, \dots, T$ ;

4: **M-шаг:** при фиксированных  $g_t$  обучить все  $b_t$ :

$$b_t := \arg \min_b \sum_{i=1}^{\ell} \tilde{g}_t(x_i) \mathcal{L}(b(x_i), y_i), \quad t = 1, \dots, T;$$

5: **E-шаг:** при фиксированных  $b_t$  оценить все  $g_t$ :

$$g_t := \arg \min_{g_t} \sum_{i=1}^{\ell} \mathcal{L}\left(\sum_{s=1}^T \tilde{g}_s(x_i) b_s(x_i), y_i\right), \quad t = 1, \dots, T;$$

6: **пока**  $\max_{t,i} |\tilde{g}_t(x_i) - \tilde{g}_t^0(x_i)| > \delta$ .

Обучение смеси с автоматическим определением числа  $T$ 

**Вход:** выборка  $X^\ell$ , параметры  $\ell_0$ ,  $\mathcal{L}_0$ ,  $\delta$ ,  $\gamma$ ;

**Выход:**  $T$ ,  $g_t(x)$ ,  $b_t(x)$ ,  $t = 1, \dots, T$ ;

1: начальное приближение:

$$b_1 := \arg \min_b \sum_{i=1}^{\ell} \mathcal{L}(b(x_i), y_i), \quad g_1(x_i) := 1, \quad i = 1, \dots, \ell;$$

2: **для всех**  $t = 2, \dots, T$

3: множество трудных объектов:

$$X_t := \{x_i : \mathcal{L}(a_{t-1}(x_i), y_i) > \mathcal{L}_0\};$$

4: **если**  $|X_t| \leq \ell_0$  **то выход**;

5:  $b_t := \arg \min_b \sum_{x_i \in X_t} \mathcal{L}(b(x_i), y_i)$ ;

6:  $g_t := \arg \min_{g_t} \sum_{i=1}^{\ell} \mathcal{L}\left(\sum_{s=1}^t g_s(x_i) b_s(x_i), y_i\right)$ ;

7:  $(g_s, b_s)_{s=1}^t := \text{ME}(X^\ell, (g_s)_{s=1}^t, t, \delta, \gamma)$ ;

## Резюме

- Простое голосование может работать не хуже бустинга
- Причины эффективности композиций: уменьшается разброс, корреляция и смещение базовых алгоритмов
- Смеси алгоритмов эффективны в тех задачах, где есть содержательные модели областей компетентности.