

Композиции классификаторов (часть 1)

К. В. Воронцов, А. В. Зухба
vokov@forecsys.ru
a__l@mail.ru

сентябрь 2016

Содержание

1 Напоминание

- Задачи обучения композиций
- Бэггинг и метод случайных подпространств
- Простое и взвешенное голосование

2 Разложения ошибки на смещение и разброс

- Основные понятия и определения
- Частные случаи
- Композиции

Определение композиции

$X^\ell = (x_i, y_i)_{i=1}^\ell \subset X \times Y$ — 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 пространств оценок и решающих правил

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

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

где $R = \mathbb{R}$, $b: X \rightarrow \mathbb{R}$, $C(b) \equiv \text{sign}(b)$.

- **Пример 2:** классификация на M классов $Y = \{1, \dots, M\}$:

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

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

- **Пример 3:** регрессия, $Y = R = \mathbb{R}$:
 $C(b) \equiv b$ — решающее правило не нужно.

Примеры корректирующих операций

- **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)

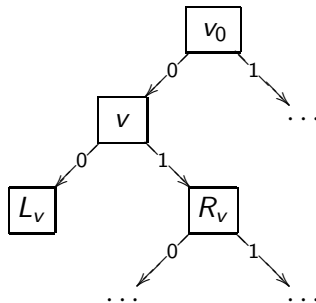
$$F(b_1(x), \dots, b_T(x)) = \sum_{t=1}^T g_t(x) b_t(x), \quad x \in X, \quad g_t: X \rightarrow \mathbb{R}.$$

Определение бинарного решающего дерева

Бинарное решающее дерево — алгоритм классификации $a(x)$, задающийся бинарным деревом:

- 1) $\forall v \in V_{\text{внутр}} \rightarrow$ предикат $\beta_v : X \rightarrow \{0, 1\}$, $\beta_v \in \mathcal{B}$
- 2) $\forall v \in V_{\text{лист}} \rightarrow$ имя класса $c_v \in Y$.

- 1: $v := v_0$;
- 2: **пока** $v \in V_{\text{внутр}}$
- 3: **если** $\beta_v(x) = 1$ **то**
- 4: переход вправо:
 $v := R_v$;
- 5: **иначе**
- 6: переход влево:
 $v := L_v$;
- 7: **вернуть** c_v .



Линейные алгоритмы

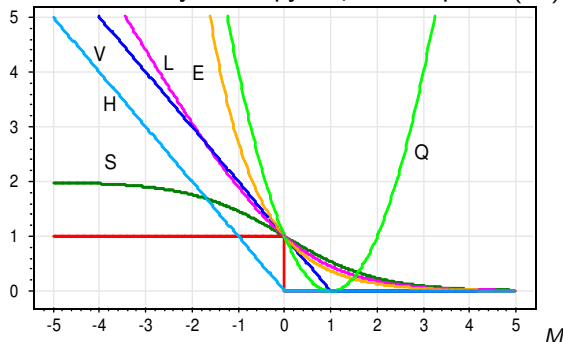
- Задача классификации с двумя классами, $Y = \{-1, +1\}$;
обучающая выборка $X^\ell = (x_i, y_i)_{i=1}^\ell$;
алгоритм классификации $a(x, w) = \text{sign } f(x, w)$,
 $f(x, w)$ — дискриминантная функция,
 w — вектор параметров.
- $f(x, w) = 0$ — разделяющая поверхность;
 $M_i(w) = y_i f(x_i, w)$ — отступ (margin) объекта x_i ;
 $M_i(w) < 0 \iff$ алгоритм $a(x, w)$ ошибается на x_i .
- Минимизация эмпирического риска:

$$Q(w) = \sum_{i=1}^{\ell} [M_i(w) < 0] \leq \tilde{Q}(w) = \sum_{i=1}^{\ell} \mathcal{L}(M_i(w)) \rightarrow \min_w;$$

функция потерь $\mathcal{L}(M)$ невозрастающая, неотрицательная.

Непрерывные аппроксимации пороговой функции потерь

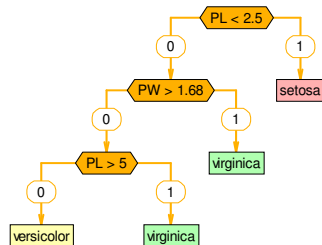
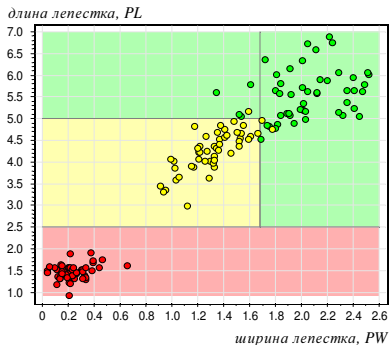
Часто используемые функции потерь $\mathcal{L}(M)$:



- | | |
|-----------------------------|--------------------------------|
| $V(M) = (1 - M)_+$ | — кусочно-линейная (SVM); |
| $S(M) = 2(1 + e^M)^{-1}$ | — сигмоидная (ANN); |
| $L(M) = \log_2(1 + e^{-M})$ | — логарифмическая (LR); |
| $E(M) = e^{-M}$ | — экспоненциальная (AdaBoost); |
| $Q(M) = (1 - M)^2$ | — квадратичная (FLD). |

Пример решающего дерева

Задача Фишера о классификации цветков ириса на 3 класса, в выборке по 50 объектов каждого класса, 4 признака.



На графике: в осях двух самых информативных признаков (из 4) два класса разделились без ошибок, на третьем 3 ошибки.

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], причём подвыборки берутся длины ℓ с возвращениями, как в методе bootstrap.

Вторую идею реализует RSM (random subspace method) [Duin, 2002].

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^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: 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.
- Не удаётся строить короткие композиции из «сильных» алгоритмов типа SVM (только длинные из слабых).

Возможно ли строить композиции проще и аккуратнее?

Простое голосование в задаче классификации

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:

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

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

Эвристика: чтобы b_{t+1} компенсировал ошибки композиции,

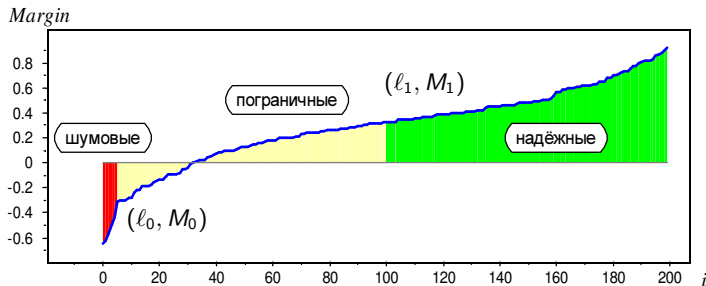
$$Q(b, U) = \sum_{x_i \in U} [y_i b(x_i) < 0] \rightarrow \min_b,$$

где $U = \{x_i : M_0 < M_{it} \leq M_1\}$,

M_0, M_1 — параметры метода обучения.

Подбор параметров M_0 и M_1

Упорядочим объекты по возрастанию отступов M_{it} :



Принцип максимизации и выравнивания отступов.

Два случая, когда b_{t+1} на объекте x_i обучать не надо:

$M_{it} < M_0$, $i < \ell_0$ — объект x_i шумовой;

$M_{it} > M_1$, $i > \ell_1$ — объект x_i уже надёжно классифицируется.

Основные понятия и определения

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 \rightarrow алгоритм

Среднеквадратичный риск:

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

Минимум среднеквадратичного риска, «недостижимый идеал»:

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

Основная мера качества метода обучения μ :

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

Разложение ошибки на шум, вариацию и смещение

Теорема

В случае квадратичной функции потерь для любого μ

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

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

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

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

Разложение bias-variance:

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

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:

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