

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 • 15 сентября 2015

Содержание

- 1 Стохастические методы построения композиций**
 - Композиции классификаторов
 - Бэггинг и метод случайных подпространств
 - 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], причём подвыборки берутся длины ℓ с повторениями, 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
- Короткие композиции из «сильных» алгоритмов типа SVM строить труднее, чем длинные из слабых

Случайный лес (Random Forest)

Основные идеи: [Breiman, 2001]

- бэггинг над решающими деревьями
- признак в каждой вершине дерева выбирается из случайного подмножества k из n признаков
- рекомендации:
 - для регрессии $k = \lfloor n/3 \rfloor$
 - для классификации $k = \lfloor \sqrt{n} \rfloor$

Основные свойства:

- случайный лес — один из самых сильных методов машинного обучения
- обычно лишь немного уступает градиентному бустингу
- но намного проще реализуется и распараллеливается

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:

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

Резюме. Почему бустинг и бэггинг работают?

- бэггинг уменьшает разброс
- бустинг уменьшает и смещение, и разброс
- композиции тем менее эффективны, чем сильнее коррелируют базовые алгоритмы
- случайный лес уменьшает корреляции базовых алгоритмов

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 $\alpha \in \mathbb{R}^n$:

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

- 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: обучение смеси алгоритмов

Итерационный процесс, аналогичный EM-алгоритму:

Вход: выборка X^ℓ , нормированные $(g_t)_{t=1}^T$, **параметры** T, δ, γ ;

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

1: **повторять**

2: $g_t^0 := g_t$ для всех $t = 1, \dots, T$;

3: **M-шаг:** при фиксированных g_t обучить все b_t :

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

4: **E-шаг:** при фиксированных b_t оценить все g_t :

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

5: нормировать компетентности:

$$(g_1(x_i), \dots, g_T(x_i)) := \text{SoftMax}(g_1(x_i), \dots, g_T(x_i); \gamma);$$

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

Обучение смеси с автоматическим определением числа T

Вход: выборка X^ℓ , 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

- Обучение смесей алгоритмов основано на принципе «разделяй и властвуй».
- Смеси алгоритмов имеет смысл строить в тех задачах, где есть априорные соображения о виде областей компетентности.
- Кроме последовательного метода построения смесей алгоритмов, известен ещё иерархический.