

# Линейные классификаторы и случайные блуждания

Евгений Соколов  
ВМК МГУ

13 апреля 2013 г.

- $\mathbb{X} = \{x_1, \dots, x_L\}$  — конечное генеральное множество объектов;
- $A = \{a_1, \dots, a_D\}$  — конечное множество алгоритмов;
- $I(a, x) = [\text{алгоритм } a \text{ ошибается на объекте } x]$  — индикатор ошибки;
- $t(a, X)$  — число ошибок алгоритма  $a$  на выборке  $X$ .
- $\nu(a, X) = t(a, X)/|X|$  — частота ошибок алгоритма  $a$  на выборке  $X$ .
- Метод обучения  $\mu : 2^{\mathbb{X}} \rightarrow A$  по произвольной выборке  $X \subset \mathbb{X}$  выбирает алгоритм  $a \in A$ .
- Вероятность переобучения:

$$Q_\varepsilon = \mathbf{P}[\nu(\mu X, \bar{X}) - \nu(\mu X, X) \geq \varepsilon] = \frac{1}{C_L^\ell} \sum_{X \sqcup \bar{X}} [\nu(\mu X, \bar{X}) - \nu(\mu X, X) \geq \varepsilon]$$

Комбинаторные оценки вероятности переобучения имеют вид

$$B_\varepsilon = \sum_{a \in \mathbb{A}} b(a),$$

где  $a$  — вершины графа расслоения-связности.

Рассмотрим более общую задачу.

- $G = (V, E)$  — граф.
- $f : V \rightarrow \mathbb{R}$  — функция на его вершинах.
- $\mathbf{p} = (p(v))_{v \in V}$  — распределение на вершинах.

Требуется оценить величину

$$F = \mathbb{E}_{\mathbf{p}} f = \sum_{v \in V} f(v)p(v).$$

Требуется оценить величину

$$F = \mathbb{E}_{\mathbf{p}} f = \sum_{v \in V} f(v)p(v).$$

- ❶ Допустим, мы умеем выбирать вершину  $v$  из распределения  $\mathbf{q} = (q(v))_{v \in V}$  на  $V$ .

- Можно использовать оценку

$$\hat{F} = \frac{1}{n} \sum_{i=1}^n f(v_i) \frac{p(v_i)}{q(v_i)}.$$

- Оценка является несмещенной и состоятельной.

- ❷ Зачастую такое сэмплирование невозможно (наш случай!).

- Стартуем из одной из вершин графа, ходим по нему, набирая вершины.
  - Как организовать ходьбу по графу так, чтобы были известны вероятности попадания в каждую вершину?
  - Случайные блуждания!
  - Формализуются с помощью аппарата марковских цепей.

- Последовательность случайных величин  $\{V_t \in V \mid t = 0, 1, 2, \dots\}$  называется марковской цепью, если

$$\mathbb{P}(V_{t+1} = v \mid V_1 = v_1, \dots, V_t = v_t) = \mathbb{P}(V_{t+1} = v \mid V_t = v_t).$$

- $V_t$  — вершина, в которой блуждание находится в момент  $t$ .
- Переходы возможны только по ребрам.
- Цепь называется *неприводимой*, если из любой ее вершины можно попасть в любую с ненулевой вероятностью:

$$\forall v, u \in V \exists s : \mathbb{P}(V_s = v \mid V_0 = u) > 0.$$

- Цепь называется *апериодической*, если для всех вершин  $v$  выполнено

$$k_v = \gcd\{t \mid \mathbb{P}(V_t = v \mid V_0 = v) > 0\} = 1.$$

- Вероятности перехода:

$$p_{vu} = \mathbb{P}(V_{t+1} = v \mid V_t = u).$$

- Матрица перехода:  $P = (p_{vu})_{v,u \in V}$ .
- Стационарное распределение:

$$P^T \pi = \pi.$$

- Если цепь неприводимая и апериодическая, то у нее гарантированно есть стационарное распределение.
- Если цепь неприводимая, но имеет период, то у нее есть стационарное распределение в обобщенном смысле, для которого все дальнейшие результаты также верны.

Требуется оценить величину

$$F = \mathbb{E}_p f = \sum_{v \in V} f(v)p(v).$$

Подойдет ли следующая оценка?

$$\mu_n(f) = \frac{1}{n} \sum_{i=1}^n f(V_i).$$

### Теорема (Закон больших чисел для марковских цепей)

Пусть  $\{V_t\}$  — неприводимая марковская цепь со стационарным распределением  $\pi$ . Тогда для любого начального распределения  $\mathbb{P}(V_0 = v)$  и для любой функции  $f : V \rightarrow \mathbb{R}$ , такой что  $\mathbb{E}_\pi |f| < \infty$ , выполнено

$$\mu_n(f) \xrightarrow[n \rightarrow \infty]{\text{a. s.}} \mathbb{E}_\pi f.$$

- Итак, оценка

$$\mu_n(f) = \frac{1}{n} \sum_{i=1}^n f(V_i)$$

является состоятельной для  $\mathbb{E}_\pi f$ .

- Как перейти от  $\mathbb{E}_\pi$  к  $\mathbb{E}_p$ ?
- Зададим корректирующую функцию

$$w(v) = \frac{p(v)}{\pi(v)}$$

и построим оценку

$$\mu_n(wf) = \frac{1}{n} \sum_{i=1}^n w(V_i) f(V_i).$$

- Можно показать, что

$$\mu_n(wf) \xrightarrow[n \rightarrow \infty]{\text{a. s.}} \mathbb{E}_p f.$$



- Конкретный случай:

$$p_{vu} = \begin{cases} \frac{1}{d(v)}, & \text{если } (v, u) \in E; \\ 0, & \text{иначе.} \end{cases}$$

- Такое блуждание называется *простым*.
- Известно стационарное распределение:  $\pi(v) = \frac{\deg(v)}{2|E|}$ .
- Весовая функция имеет вид

$$w(v) = \frac{p(v)}{\pi(v)} = \frac{2|E|}{|V|} \frac{1}{\deg(v)}$$

- Не знаем  $|V|$  и  $|E|$ !
- Рассмотрим оценку

$$\hat{F}_{\text{RW}} = \frac{\mu_n(wf)}{\mu_n(w)} = \frac{\sum_{i=1}^n w(V_i) f(V_i)}{\sum_{i=1}^n w(V_i)} = \frac{\sum_{i=1}^n f(V_i) / \deg(V_i)}{\sum_{i=1}^n 1 / \deg(V_i)}.$$

- Оценка является состоятельной и несмещенной.

- Вклады алгоритмов из разных слоев графа расслоения-связности отличаются на порядки;
- Дисперсия функции  $f$  большая, и для оценивания  $\mathbb{E}_p f$  требуется большая выборка;
- При этом в пределах одного слоя дисперсия мала!
- Новая оценка:

$$\hat{F}'_{RW} = \frac{1}{|V|} \sum_{m=1}^L \hat{f}_m \hat{c}_m,$$

где  $\hat{f}_m$  — среднее значение функции  $f$  по  $m$ -му слою,  
 $\hat{c}_m$  — мощность  $m$ -го слоя.

Подробности:

- Пусть  $C_1, \dots, C_L$  — слои графа,  $c_1, \dots, c_L$  — их мощности.
- Преобразуем матожидание:

$$F = \mathbb{E}_{\mathbf{p}} f = \frac{1}{|V|} \sum_{v \in V} f(v) = \frac{1}{|V|} \sum_{m=1}^L \left( \frac{1}{c_m} \sum_{v \in C_m} f(v) \right) c_m.$$

- Новая оценка:

$$\hat{F}'_{RW} = \frac{1}{|V|} \sum_{m=1}^L \hat{f}_m \hat{c}_m;$$

$$\hat{f}_m = \frac{\sum_{i=1}^n \frac{f(V_i)}{\deg(V_i)} [m(V_i) = m]}{\sum_{i=1}^n \frac{1}{\deg(V_i)} [m(V_i) = m]};$$

$$\hat{c}_m = \frac{|V|}{n} \sum_{i=1}^n \left( \frac{1}{\deg(V_i)} \sum_{u \in N(V_i)} [m(u) = m] \right).$$

**Вход:** Граф  $G = (V, E)$ ; число итераций  $N$ ;  
набор стартовых вершин  $P = (v^1, \dots, v^s)$ ;

**Выход:** Выборка вершин графа  $v_1, v_2, \dots, v_N$ ;

1: для  $i = 1, \dots, N$

2: выбрать  $v \in P$  с вероятностью  $\frac{\deg(v)}{\sum_{u \in P} \deg(u)}$ ;

3: с вероятностью  $\frac{1}{2}$

4: выбрать вершину  $v'$  из равномерного

5: распределения на  $\{v' \in V \mid (v, v') \in E\}$ ;

6:  $v_i := v'$ ;

7: иначе

8:  $v' := v$ ;  $v_i := v$ ;

9: Заменить в  $P$  вершину  $v$  на  $v'$ ;

**Основная задача:** построение непереобученных композиций линейных классификаторов вида

$$a(x) = \text{sign} \sum_{i=1}^p \text{th} \langle w_i, x \rangle.$$

*Каждый базовый классификатор строится в подпространстве малой размерности.*

Пусть известна некоторая оценка вероятности переобучения

$$P[\nu(\mu X, \bar{X}) - \nu(\mu X, X) \geq \varepsilon] \leq \eta(\varepsilon).$$

Обратив ее, можно оценить частоту ошибок на контроле: с вероятностью не менее  $1 - \eta$

$$\nu(\mu X, \bar{X}) \leq \nu(\mu X, X) + \varepsilon(\eta).$$

Величину в правой части неравенства будем использовать в качестве критерия при отборе признаков.

- Базовый алгоритм — линейный классификатор, настроенный методом SVM;
- Классификатор строится по отобранным признакам;
- Отбор признаков: жадное наращивание с минимизацией комбинаторного критерия:

$$\nu(\mu X, X) + \varepsilon(1/2) \rightarrow \min .$$

- Для вычисления комбинаторной оценки на каждом наборе признаков используем случайные блуждания.

**Вход:** Выборка  $X$ ; параметры  $T, \ell_0, \ell_1$ ;

**Выход:** Базовые линейные классификаторы  $b_1, \dots, b_T$

1: Инициализировать веса и отступы:

$$w_i = 1, M_i = 0 \text{ для всех } i = 1, \dots, \ell;$$

2: для  $t = 1, \dots, T$ , пока не выполнен критерий останова

3: Обучить базовый алгоритм (SVM по отобраннным признакам):

$$b_t := \underset{b}{\operatorname{argmin}} Q(b, X, W);$$

4: Обновить значения отступов:

$$M_i = M_i + y_i b_t(x_i) \text{ для всех } i = 1, \dots, L;$$

5: упорядочить выборку  $X$  по возрастанию отступов  $M_i$ ;

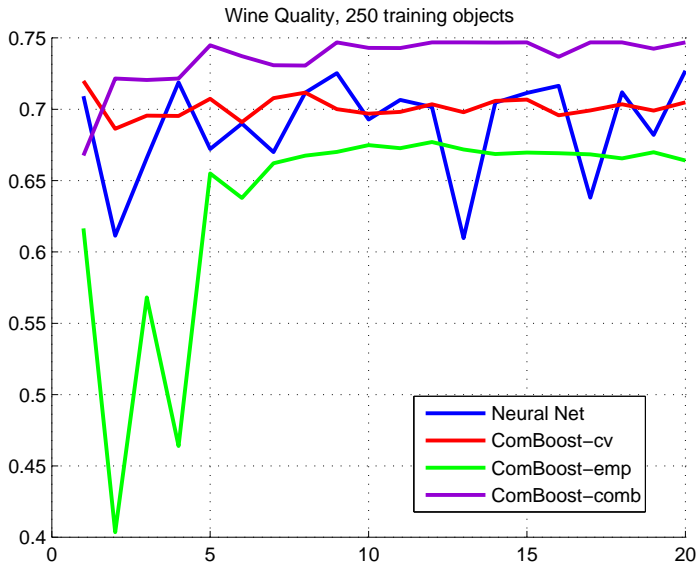
6: Отобрать объекты для обучения следующего базового алгоритма:

$$w_i = [\ell_0 \leq i \leq \ell_1];$$

	Wine Quality	Statlog	Waveform	Faults
Оценка, основанная на эмпирическом риске $\nu(\mu, \mathbb{X})$	64,70	84,92	84,78	73,39
Двухслойная нейронная сеть, настроенная методом backpropagation	72,06	85,41	86,79	74,51
Оценка (1), в которой обращается комбинаторная оценка [Воронцов, 2010], вычисленная с помощью случайных блужданий	69,48	86,26	85,77	<b>77,81</b>
Оценка скользящего контроля, вычисленная по 100 случайным разбиениям	71,06	85,26	86,38	75,76
Оценка (1), в которой обращается комбинаторная оценка с истоками, вычисленная с помощью случайных блужданий	<b>74,68</b>	<b>86,75</b>	<b>86,91</b>	74,03

$$Q_c = \nu(a_0, X) + \varepsilon \left(\frac{1}{2}\right) \quad (1)$$





Вопросы?