

# Логические алгоритмы классификации

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)»

сентябрь 2011

## Определения и обозначения

$X$  — пространство объектов;

$Y$  — множество ответов;

$f_1(x), \dots, f_n(x)$  — признаки;

$X^\ell = (x_i, y_i)_{i=1}^\ell$  — обучающая выборка,  $y_i = y(x_i)$ .

### Определение (пока неформальное)

*Закономерность (правило, rule) — это предикат  $\varphi: X \rightarrow \{0, 1\}$ , удовлетворяющий двум требованиям:*

- 1) интерпретируемость ( $\varphi$  зависит от 1–7 признаков);*
- 2) информативность относительно класса  $c \in Y$ :*

$$p_c(\varphi) = \#\{x_i: \varphi(x_i) = 1 \text{ и } y_i = c\} \rightarrow \max;$$

$$n_c(\varphi) = \#\{x_i: \varphi(x_i) = 1 \text{ и } y_i \neq c\} \rightarrow \min;$$

Если  $\varphi(x) = 1$ , то говорят « $\varphi$  выделяет  $x$ » ( $\varphi$  covers  $x$ ).

## Содержание

- 1 Понятия закономерности и информативности**
  - Определения и обозначения
  - Интерпретируемость
  - Информативность
- 2 Методы поиска информативных закономерностей**
  - Жадный алгоритм
  - Алгоритмы на основе отбора признаков
  - Бинаризация данных
- 3 Композиции закономерностей**
  - Голосование закономерностей
  - Решающий список
  - Решающие деревья и леса

## Требование интерпретируемости

### Пример (из области медицины)

*Если возраст  $> 60$  и пациент ранее перенёс инфаркт,  
то операцию не делать, риск отрицательного исхода 60%.*

### Пример (из области кредитного скоринга)

*Если в анкете указан домашний телефон  
и зарплата  $> \$2000$  и сумма кредита  $< \$5000$   
то кредит можно выдать, риск дефолта 5%.*

### Требования интерпретируемости:

- 1)  $\varphi$  зависит от малого числа признаков;
- 2) формула  $\varphi$  выражается на естественном языке.

## Виды закономерностей

Параметрическое семейство *конъюнкций пороговых условий*:

$$\varphi(x) = \bigwedge_{j \in J} [\alpha_j \leq f_j(x) \leq \beta_j].$$

Параметрическое семейство *шаров*:

$$\varphi(x) = \left[ \sum_{j \in J} \alpha_j |f_j(x) - f_j(x_0)|^\gamma \leq R^\gamma \right].$$

Параметрическое семейство *полуплоскостей*:

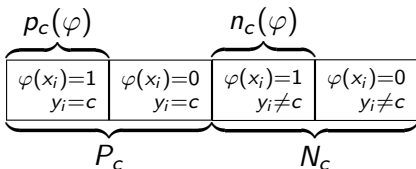
$$\varphi(x) = \left[ \sum_{j \in J} \alpha_j f_j(x) \geq \alpha_0 \right].$$

Параметрическое семейство *синдромных правил*:

$$\varphi(x) = \left[ \sum_{j \in J} [\alpha_j \leq f_j(x) \leq \beta_j] \geq K \right].$$

Основная проблема — отбор признаков  $J \subseteq \{1, \dots, n\}$ .

## Логический (эвристический) критерий закономерности



### Определение

Предикат  $\varphi(x)$  — логическая  $\varepsilon, \delta$ -закономерность класса  $c \in Y$

$$E_c(\varphi, X^\ell) = \frac{n_c(\varphi)}{p_c(\varphi) + n_c(\varphi)} \leq \varepsilon;$$

$$D_c(\varphi, X^\ell) = \frac{p_c(\varphi)}{\ell} \geq \delta.$$

**Проблема:** хотелось бы иметь один скалярный критерий.

## Статистический критерий информативности

**Точный тест Фишера.** Пусть  $X$  — в.п., выборка  $X^\ell$  — i.i.d.  
 Гипотеза  $H_0$ :  $y(x)$  и  $\varphi(x)$  — независимые случайные величины.  
 Тогда вероятность реализации пары  $(p, n)$  — ГГР:

$$P(p, n) = \frac{C_P^p C_N^n}{C_{P+N}^{p+n}}, \quad 0 \leq p \leq P, \quad 0 \leq n \leq N,$$

где  $C_N^n = \frac{N!}{n!(N-n)!}$  — биномиальные коэффициенты.

### Определение

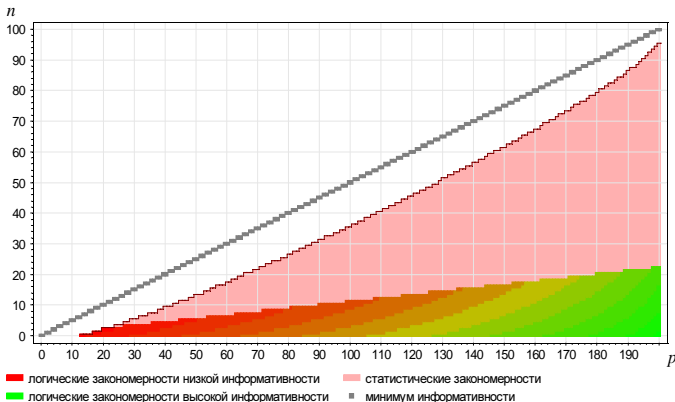
*Информативность предиката  $\varphi(x)$  относительно класса  $c \in Y$ :*

$$I_c(\varphi, X^\ell) = -\ln \frac{C_{P_c}^{p_c(\varphi)} C_{N_c}^{n_c(\varphi)}}{C_{P_c+N_c}^{p_c(\varphi)+n_c(\varphi)}},$$

$I_c(\varphi, X^\ell) \geq I_0$  — статистическая закономерность класса  $c \in Y$ .

## Соотношение логического и статистического критериев

Области логических ( $\varepsilon = 0.1$ ) и статистических ( $I_0 = 5$ ) закономерностей в координатах  $(p, n)$  при  $P = 200, N = 100$ .



Философский вопрос: закономерность == неслучайность?



## Задача перебора конъюнкций

Пусть  $\mathcal{B}$  — конечное множество *элементарных предикатов*.  
Множество конъюнкций с ограниченным числом термов из  $\mathcal{B}$ :

$$\mathcal{K}_K[\mathcal{B}] = \{ \varphi(x) = \beta_1(x) \wedge \dots \wedge \beta_k(x) \mid \beta_1, \dots, \beta_k \in \mathcal{B}, k \leq K \}.$$

Число допустимых конъюнкций:  $O(|\mathcal{B}|^K)$  — ооооочень много!

### Семейство методов локального поиска

*Окрестность*  $V(\varphi)$  — все конъюнкции, получаемые из  $V(\varphi)$  добавлением, изъятием или модификацией одного из термов.

Основная идея: на  $t$ -й итерации

$$\varphi_t := \arg \max_{\varphi \in V(\varphi_{t-1})} I_c(\varphi, X^\ell).$$

## Обобщённый алгоритм локального поиска

**Вход:** выборка  $X^\ell$ ; класс  $c \in Y$ ;  
начальное приближение  $\varphi_0$ ; параметры  $t_{\max}$ ,  $d$ ,  $\varepsilon$ ;

**Выход:** конъюнкция  $\varphi$ ;

- 
- 1:  $I^* := I_c(\varphi_0, X^\ell)$ ;  $\varphi^* := \varphi_0$ ;
  - 2: **для всех**  $t = 1, \dots, t_{\max}$
  - 3:  $\varphi_t := \arg \max_{\varphi \in V(\varphi_{t-1})} I_c(\varphi, X^\ell)$  — перспективная конъюнкция;
  - 4:  $\varphi_t^* := \arg \max_{\substack{\varphi \in V(\varphi_{t-1}) \\ E_c(\varphi) < \varepsilon}} I_c(\varphi, X^\ell)$  — лучшая конъюнкция;
  - 5: **если**  $I_c(\varphi_t^*) > I^*$  **то**  $t^* := t$ ;  $\varphi^* := \varphi_t^*$ ;  $I^* := I_c(\varphi^*)$ ;
  - 6: **если**  $t - t^* > d$  **то вернуть** ;
  - 7: **вернуть**  $\varphi^*$ ;

## Частные случаи и модификации

- **жадный алгоритм:**  
 $V(\varphi)$  — только добавления термов;  $\varphi_0 = \emptyset$ ;
- **стохастический локальный поиск (SLS):**  
 $V(\varphi)$  — случайное подмножество всевозможных добавлений, удалений, модификаций термов;  $\varphi_0 = \emptyset$ ;
- **стабилизация:**  
 $V(\varphi)$  — удаления термов или изменение параметров в термах;  $\varphi_0 \neq \emptyset$ ;
- **редукция:**  
 $V(\varphi)$  — только удаления термов;  $\varphi_0 \neq \emptyset$ ;  
 $I_c(\varphi, X^k)$  оценивается **по контрольной выборке**  $X^k$ ;
- **поиск в ширину:**  
на каждой итерации строится множество конъюнкций  
 $\Phi_t = \{\varphi_t\}$ .

## Поиск закономерностей — это отбор признаков

Отличия от методов отбора признаков:

- вместо внешнего критерия  $Q_{\text{ext}} \rightarrow \min$   
критерий информативности  $I_c \rightarrow \max$ ;
- есть ограничение на число признаков  $|J| \leq K$ .
- надо построить не одно, а *много различных* правил,  
которые должны образовать покрытие выборки.

Все механизмы отбора признаков подходят:

- добавления–удаления;
- поиск в глубину;
- поиск в ширину;
- генетические (эволюционные) алгоритмы;
- случайный поиск с адаптацией.

## Вспомогательная задача бинаризации вещественного признака

**Цель:** сократить перебор предикатов вида  $[\alpha \leq f(x) \leq \beta]$ .

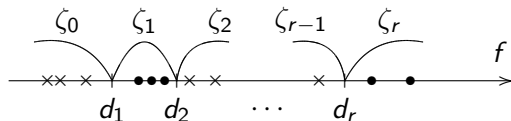
**Дано:** выборка значений вещественного признака  $f(x_i)$ ,  $x_i \in X^\ell$ .

**Найти:** наиболее информативное разбиение области значений признака на относительно небольшое число зон:

$$\zeta_0(x) = [f(x) < d_1];$$

$$\zeta_s(x) = [d_s \leq f(x) < d_{s+1}], \quad s = 1, \dots, r-1;$$

$$\zeta_r(x) = [d_r \leq f(x)].$$



## Алгоритм разбиения области значений признака на зоны

**Вход:** выборка  $X^\ell$ ; класс  $c \in Y$ ; параметры  $r$  и  $\delta_0$ .

**Выход:**  $D = \{d_1 < \dots < d_r\}$  — последовательность порогов;

- 
- 1:  $D := \emptyset$ ; упорядочить выборку  $X^\ell$  по возрастанию  $f(x_i)$ ;
  - 2: **для всех**  $i = 2, \dots, \ell$
  - 3: **если**  $f(x_{i-1}) \neq f(x_i)$  и  $[y_{i-1} = c] \neq [y_i = c]$  **то**
  - 4:     добавить порог  $\frac{1}{2}(f(x_{i-1}) + f(x_i))$  в конец  $D$ ;
  - 5: **повторять**
  - 6:     **для всех**  $d_j \in D, j = 1, \dots, |D| - 1$
  - 7:          $\delta I_j := I_c(\zeta_{i-1} \vee \zeta_i \vee \zeta_{i+1}) - \max\{I_c(\zeta_{i-1}), I_c(\zeta_i), I_c(\zeta_{i+1})\}$ ;
  - 8:          $i := \arg \max_s \delta I_s$ ;
  - 9:     **если**  $\delta I_j > \delta_0$  **то**
  - 10:         слить зоны  $\zeta_{i-1}, \zeta_i, \zeta_{i+1}$ , удалив  $d_j$  и  $d_{j+1}$  из  $D$ ;
  - 11: **пока**  $|D| > r + 1$ .

## Взвешенное голосование закономерностей

$R_c = \{\varphi_c^t(x) : t = 1, \dots, T_c\}$  — список закономерностей класса  $c$ .

Взвешенное голосование (weighted voting):

$$a(x) = \arg \max_{c \in Y} \Gamma_c(x), \quad \Gamma_c(x) = \sum_{t=1}^{T_c} \alpha_c^t \varphi_c^t(x).$$

Жадный алгоритм построения композиции:

**Вход:** выборка  $X^\ell$ , семейство правил  $\Phi$ ;

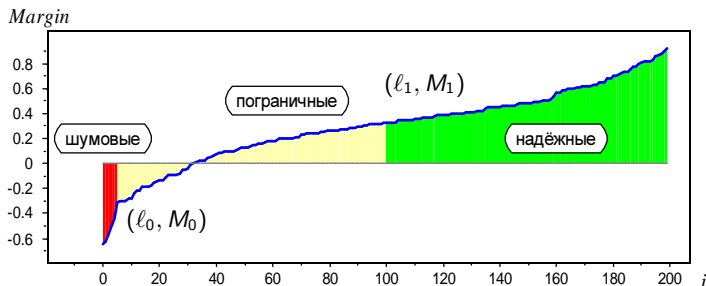
**Выход:**  $(\varphi_c^t, \alpha_c^t)$ ,  $t = 1, \dots, T_c$ ,  $\forall c \in Y$ ;

- 1: для всех правил  $\varphi \in \Phi$
- 2: если  $\varphi$  — закономерность для некоторого класса  $c$   
и  $\varphi$  существенно отличается от  $\forall \psi \in R_c$   
и  $\varphi$  существенно улучшает композицию то
- 3: добавить  $\varphi_c^t = \varphi$  в  $R_c$  и оценить  $\alpha_c^t$ ;

## Что значит «правило существенно улучшает композицию»?

Отступ (margin) объекта  $x_i$ :  $M(x_i) = \Gamma_{y_i}(x_i) - \max_{y \neq y_i} \Gamma_y(x_i)$ .

Распределение объектов выборки  $X^\ell$  по значениям отступов:



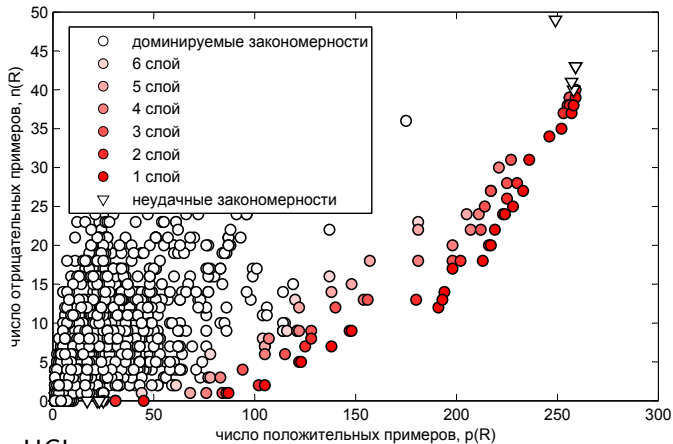
Задача максимизации и выравнивания отступов:

$$\sum_{i=l_0}^{l_1} M(x_i) \rightarrow \max_{(\varphi_c^t, \alpha_c^t)} .$$



## Что значит «правила существенно различаются»?

Парето-фронт — множество недоминируемых правил  
(правило  $R$  недоминируемо, если правее и ниже правил нет)



задача UCI:german

## Бустинг закономерностей

Два класса,  $Y = \{-1, +1\}$ .

Композиция из  $T = T_{-1} + T_{+1}$  закономерностей:

$$a_T(x) = \text{sign} \left( \underbrace{\sum_{t=1}^{T+1} \alpha_{+1}^t \varphi_{+1}^t(x)}_{\Gamma_{+1}(x)} - \underbrace{\sum_{t=1}^{T-1} \alpha_{-1}^t \varphi_{-1}^t(x)}_{\Gamma_{-1}(x)} \right), \quad \alpha_c^t > 0.$$

**Эвристика 1:** добавляем  $\alpha_c^{T+1} \varphi_c^{T+1}(x)$ , не трогая  $\alpha_c^t \varphi_c^t(x)$ ,  $t < T$ .

**Эвристика 2:** экспоненциальная аппроксимация:

$$\begin{aligned} Q_T &= \sum_{i=1}^{\ell} [\Gamma_{y_i}(x_i) - \Gamma_{-y_i}(x_i) < 0] \leq \\ &\leq \tilde{Q}_T = \sum_{i=1}^{\ell} \underbrace{\exp(\Gamma_{-y_i}(x_i) - \Gamma_{y_i}(x_i))}_{w_i} = \sum_{i=1}^{\ell} w_i. \end{aligned}$$

## Основная теорема бустинга

Добавляем ещё одну закономерность  $\alpha_c \varphi_c(x)$ :

$$\tilde{Q}_{T+1}(\varphi_c, \alpha_c) = \sum_{y_i=c} w_i e^{-\alpha_c \varphi_c(x_i)} + \sum_{y_i \neq c} w_i e^{\alpha_c \varphi_c(x_i)}$$

### Теорема

Минимум функционала  $\tilde{Q}_{T+1}(\varphi_c, \alpha_c)$  достигается при

$$\varphi_c^* = \arg \max_{\varphi} \sqrt{p_c^w(\varphi)} - \sqrt{n_c^w(\varphi)};$$

$$\alpha_c^* = \frac{1}{2} \ln \frac{p_c^w(\varphi_c^*)}{n_c^w(\varphi_c^*)};$$

$$p_c^w(\varphi) = \sum_{i=1}^{\ell} w_i [y_i = c] \varphi(x_i);$$

$$n_c^w(\varphi) = \sum_{i=1}^{\ell} w_i [y_i \neq c] \varphi(x_i);$$

## Алгоритм AdaBoost для закономерностей

**Вход:** выборка  $X^\ell$ ; семейство правил  $\Phi$ ; параметры  $T, \delta$ ;

**Выход:** закономерности и их веса  $\varphi_c^t(x), \alpha_c^t, t = 1..T_c, c \in Y$ ;

- 
- 1: инициализация:  $w_i := 1, i = 1, \dots, \ell$ ;
  - 2: **для всех**  $t = 1, \dots, T$
  - 3:  $c := c_t$  — выбрать класс закономерности;
  - 4:  $\varphi_c^t := \arg \max_{\varphi \in \Phi} \sqrt{p_c^w(\varphi)} - \sqrt{n_c^w(\varphi)}$ ;
  - 5:  $\alpha_c^t := \frac{1}{2} \ln \frac{p_c^w(\varphi) + \delta}{n_c^w(\varphi) + \delta}$ ;
  - 6: **для всех**  $i = 1, \dots, \ell$
  - 7: **если**  $\varphi_c(x_i) = 1$  **то**  $w_i := \begin{cases} w_i \exp(-\alpha_c^t), & y_i = c; \\ w_i \exp(\alpha_c^t), & y_i \neq c; \end{cases}$
  - 8: нормировка:  $Z := \frac{1}{\ell} \sum_{i=1}^{\ell} w_i$ ;  $w_i := w_i / Z, i = 1, \dots, \ell$ ;

## Алгоритм AdaBoost для закономерностей

### Достоинства:

- Бустинг одновременно
  - maximизирует информативность правил;
  - увеличивает их различность;
  - строит из них покрытие выборки.
- Низкое переобучение благодаря выравниванию отступов.
- Универсальность: можно использовать любое семейство  $\Phi$ .

### Недостатки:

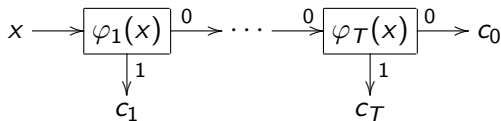
- Чувствительность к шуму (лучше использовать AnyBoost).
- В некоторых задачах переобучение всё же наблюдается, хотя, при очень больших  $T \sim 10^2 \div 10^3$ .
- Утрата интерпретируемости при  $T$ , больших  $T \sim 10 \div 20$ .

## Определение решающего списка

*Решающий список* (decision list, DL)

— алгоритм классификации  $a: X \rightarrow Y$ , который задаётся закономерностями  $\varphi_1(x), \dots, \varphi_T(x)$  классов  $c_1, \dots, c_T \in Y$ :

- 1: **для всех**  $t = 1, \dots, T$
- 2: **если**  $\varphi_t(x) = 1$  **то**
- 3: **вернуть**  $c_t$ ;
- 4: **вернуть**  $c_0$ .



«Особый ответ»  $c_0$  — отказ от классификации объекта  $x$ .

## Построение решающего списка

**Вход:** выборка  $X^\ell$ ; семейство предикатов  $\Phi$ ;

параметры:  $T_{\max}$ ,  $I_{\min}$ ,  $E_{\max}$ ,  $\ell_0$ ;

**Выход:** решающий список  $\{\varphi_t, c_t\}_{t=1}^T$ ;

- 
- 1:  $U := X^\ell$ ;
  - 2: **для всех**  $t := 1, \dots, T_{\max}$
  - 3:  $c := c_t$  — выбрать класс из  $Y$ ;
  - 4:  $\Phi' = \{\varphi \in \Phi : E_c(\varphi, U) \leq E_{\max}\}$ ;
  - 5:  $\varphi_t := \arg \max_{\varphi \in \Phi'} I_c(\varphi, U)$ ;
  - 6: **если**  $I_c(\varphi_t, U) < I_{\min}$  **то выход**;
  - 7: исключить из выборки объекты, выделенные правилом  $\varphi_t$ :  
 $U := \{x \in U : \varphi_t(x) = 0\}$ ;
  - 8: **если**  $|U| \leq \ell_0$  **то выход**;

## Замечания к алгоритму построения решающего списка

- Параметр  $E_{\max}$  позволяет управлять сложностью списка:  
 $E_{\max} \downarrow \Rightarrow p(\varphi_t) \downarrow, T \uparrow.$
- Стратегии выбора класса  $c_t$ :
  - 1) все классы по очереди;
  - 2) на каждом шаге определяется оптимальный класс:  
$$(\varphi_t, c_t) := \arg \max_{\varphi \in \Phi', c \in Y} I_c(\varphi, U);$$
- Простой обход проблемы пропусков в данных.
- Другие названия:
  - комитет с логикой старшинства;
  - голосование по старшинству;
  - машина покрывающих множеств (SCM);



## Решающие списки: достоинства и недостатки

### Достоинства:

- Интерпретируемость и простота классификации.
- Универсальность: можно использовать любое семейство  $\Phi$ .
- Допустимы разнотипные данные и данные с пропусками.
- Правила получаются различными. Можно построить несколько списков и по ним проголосовать.

### Недостатки:

- При неудачном выборе  $\Phi$  список может не построиться, будет много отказов от классификации.
- Список плохо интерпретируется, если он длинный и/или правила различных классов следуют вперемежку.
- Качество классификации обычно ниже, чем у голосования, когда правила могут компенсировать ошибки друг друга.

## Сухой остаток

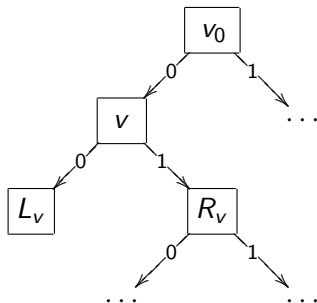
- Правило — это интерпретируемый предикат  $X \rightarrow \{0, 1\}$ .
- Закономерность — это информативное правило.
- Существует очень много критериев информативности.
- Статистический критерий — для поиска закономерностей.
- Логический  $\varepsilon, \delta$ -критерий — для отбора закономерностей.
- Механизмы Rule Induction — те же, что Features Selection.
- Принципы отбора закономерностей в композицию:
  - максимизация информативности;
  - увеличение различности правил;
  - максимизации отступов для надёжного покрытия выборки.

## Определение бинарного решающего дерева

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

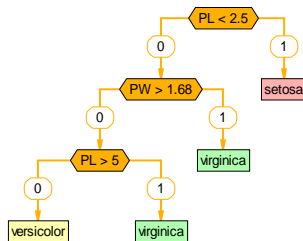
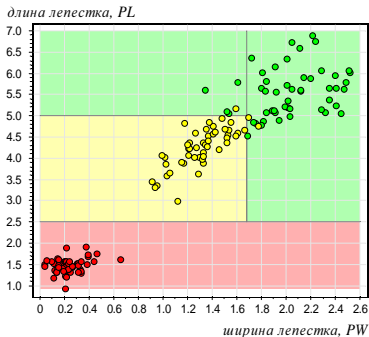
- 1)  $\forall v \in V_{\text{внутр}} \rightarrow$  предикат  $\beta_v : X \rightarrow \{0, 1\}$ ,  $\beta \in \mathcal{B}$
- 2)  $\forall v \in V_{\text{лист}} \rightarrow$  имя класса  $c_v \in Y$ .

- 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**  $c_v$ .



## 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.

## Жадный алгоритм построения дерева ID3

- 1: **ПРОЦЕДУРА** LearnID3 ( $U \subseteq X^\ell$ );
- 2: **если** все объекты из  $U$  лежат в одном классе  $c \in Y$  **то**
- 3: **вернуть** новый лист  $v$ ,  $c_v := c$ ;
- 4: найти предикат с максимальной информативностью:  
$$\beta := \arg \max_{\beta \in \mathcal{B}} I(\beta, U);$$
- 5: разбить выборку на две части  $U = U_0 \cup U_1$  по предикату  $\beta$ :  
$$U_0 := \{x \in U : \beta(x) = 0\};$$
$$U_1 := \{x \in U : \beta(x) = 1\};$$
- 6: **если**  $U_0 = \emptyset$  или  $U_1 = \emptyset$  **то**
- 7: **вернуть** новый лист  $v$ ,  $c_v := \text{Мажоритарный класс}(U)$ ;
- 8: создать новую внутреннюю вершину  $v$ :  $\beta_v := \beta$ ;  
построить левое поддерево:  $L_v := \text{LearnID3}(U_0)$ ;  
построить правое поддерево:  $R_v := \text{LearnID3}(U_1)$ ;
- 9: **вернуть**  $v$ ;

## Разновидности критериев ветвления

1. Отделение одного класса (слишком сильное ограничение):

$$I(\beta, X^\ell) = \max_{c \in Y} I_c(\beta, U).$$

2. Многоклассовый энтропийный критерий:

$$I(\beta, X^\ell) = \sum_{c \in Y} h\left(\frac{P_c}{\ell}\right) - \frac{p}{\ell} \sum_{c \in Y} h\left(\frac{p_c}{p}\right) - \frac{\ell - p}{\ell} \sum_{c \in Y} h\left(\frac{P_c - p_c}{\ell - p}\right),$$

где  $P_c = \#\{x_i: y_i = c\}$ ,  $p = \#\{x_i: \beta(x_i) = 1\}$ ,  $h(z) \equiv -z \log_2 z$ .

3. Критерий Джини:

$$I(\beta, X^\ell) = \#\{(x_i, x_j): \beta(x_i) = \beta(x_j) \text{ и } y_i \neq y_j\}.$$

4.  $D$ -критерий В.И.Донского:

$$I(\beta, X^\ell) = \#\{(x_i, x_j): \beta(x_i) \neq \beta(x_j) \text{ и } y_i = y_j\}.$$

## Обработка пропусков

### На стадии обучения:

- Если  $\beta(x_i)$  не определено, то при вычислении  $I(\beta, U)$  объект  $x_i$  исключается из выборки  $U$ .
- Для  $\forall v \in V_{\text{внутр}}$  оценивается:  
 $\hat{p}_L = |U_0|/|U|$  — вероятность левой ветви;  
 $\hat{p}_R = |U_1|/|U|$  — вероятность правой ветви.

### На стадии классификации:

- $\beta_v(x)$  не определено  $\Rightarrow$  пропорциональное распределение:

$$\hat{P}_v(y|x) = \begin{cases} [y = c_v], & v \in V_{\text{лист}}; \\ \hat{p}_L \hat{P}_{L_v}(y|x) + \hat{p}_R \hat{P}_{R_v}(y|x), & v \in V_{\text{внутр}}. \end{cases}$$

- Окончательное решение — байесовское правило:

$$y = \arg \max_{y \in Y} \hat{P}_v(y|x).$$

## Решающие деревья: достоинства и недостатки

### Достоинства:

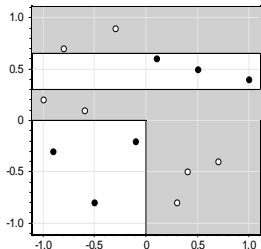
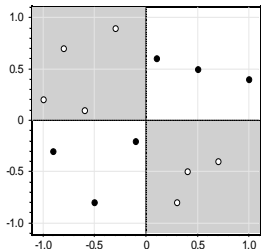
- Интерпретируемость и простота классификации.
- Гибкость: можно варьировать множество  $\mathcal{B}$ .
- Допустимы разнотипные данные и данные с пропусками.
- Трудоёмкость линейна по длине выборки  $O(|\mathcal{B}|hl)$ .
- Не бывает отказов от классификации.

### Недостатки:

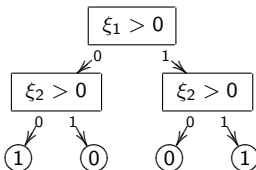
- Жадный ID3 переусложняет структуру дерева, и, как следствие, сильно переобучается.
- Фрагментация выборки: чем дальше  $v$  от корня, тем меньше статистическая надёжность выбора  $\beta_v, c_v$ .
- Высокая чувствительность к шуму, к составу выборки, к критерию информативности.



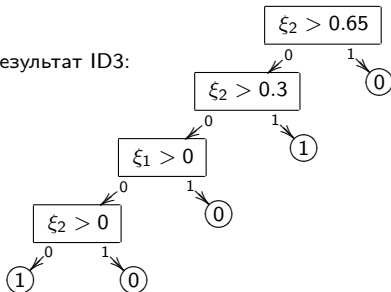
## Жадный ID3 переусложняет структуру дерева



Оптимальное дерево для задачи XOR:



Результат ID3:



## Стратегия пред-просмотра (look ahead)

Шаг 6:

**если**  $U_0 = \emptyset$  или  $U_1 = \emptyset$  **то**

**вернуть** новый лист  $v$ ,  $c_v :=$  Мажоритарный класс( $U$ );

Шаг 6 заменяется на более ресурсоёмкую процедуру:

**для всех** деревьев  $T$  глубины  $h$

$r(U)$  = число ошибок дерева  $T$  на выборке  $U$ ;

**вернуть** новый лист  $v$ ,  $\beta_v :=$  корень лучшего поддеревя;

Достоинства:

- Задача XOR уже решается почти идеально.

Недостатки:

- При  $h > 2$  оооооочень долго.
- На реальных данных улучшение незначительно.

## Стратегия пред-редукции (pre-pruning)

Шаг 6:

**если**  $U_0 = \emptyset$  или  $U_1 = \emptyset$  **то**  
**вернуть** новый лист  $v$ ;

Шаг 6 заменяется на более мягкое условие:

**если**  $I(\beta, U) \leq I_0$  **то**  
**вернуть** новый лист  $v$ ;

**Достоинства:**

- Сразу строится более простое дерево.

**Недостатки:**

- Качество дерева может и не улучшиться.

## Стратегия пост-редукции (post-pruning: C4.5, CART)

$X^k$  — независимая контрольная выборка,  $k \approx 0.5l$ .

- 1: **для всех**  $v \in V_{\text{внутр}}$
- 2:  $S_v :=$  подмножество объектов  $X^k$ , дошедших до  $v$ ;
- 3: **если**  $S_v = \emptyset$  **то**
- 4: **вернуть** новый лист  $v$ ,  $c_v := \text{Мажоритарный класс}(U)$ ;
- 5: число ошибок при классификации  $S_v$  четырьмя способами:
  - $r(v)$  — поддеревом, растущим из вершины  $v$ ;
  - $r_L(v)$  — поддеревом левой дочерней вершины  $L_v$ ;
  - $r_R(v)$  — поддеревом правой дочерней вершины  $R_v$ ;
  - $r_c(v)$  — к классу  $c \in Y$ .
- 6: в зависимости от того, какое из них минимально:
  - сохранить поддерево  $v$ ;
  - заменить поддерево  $v$  поддеревом  $L_v$ ;
  - заменить поддерево  $v$  поддеревом  $R_v$ ;
  - заменить поддерево  $v$  листом,  $c_v := \arg \min_{c \in Y} r_c(v)$ .

## Преобразование решающего дерева в решающий список (C4.5-rules)

- Для любого бинарного решающего дерева

$$a(x) = \arg \max_{y \in Y} \sum_{v \in V_{\text{лист}}} [c_v = y] K_v(x),$$

где  $K_v(x)$  — конъюнкция по всем рёбрам пути  $[v_0, v]$ :

$$K_v(x) = \bigwedge_{(u, R_u)} \beta_u(x) \bigwedge_{(u, L_u)} \bar{\beta}_u(x).$$

- Редукция  $K_v(x)$ ,  $\forall v \in V_{\text{лист}}$  по контрольной выборке  $X^k$ .

**Достоинства:**

- Переобучение, как правило, уменьшается.

**Недостатки:**

- Преобразование в список необратимо: это уже не дерево.

## Решающие леса

### Бустинг решающих деревьев:

- базовые классификаторы — деревья ограниченной глубины.

### Случайные леса (Random Forests):

- $T$  деревьев обучаются по случайным подвыборкам (bagging);
- $\mathcal{B}$  — случайное множество гиперплоскостей;
- $I(\beta, U)$  — энтропийный критерий информативности;
- Простое голосование по  $T$  деревьям.

### Решающий список из решающих деревьев:

- При образовании статистически ненадёжного листа этот лист заменяется переходом к следующему дереву.
- Следующее дерево строится по всем объектам, прошедшим через ненадёжные листы предыдущего дерева.

## Чередующееся решающее дерево — ADT (Alternating Decision Tree) [Freund, Mason, 1999]

Пусть  $Y = \{-1, +1\}$ .

Для каждой вершины  $v \in V$  задаётся:

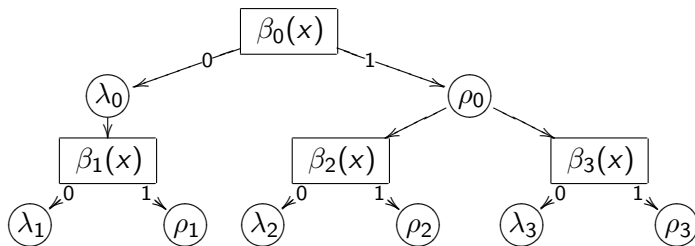
- 1) предикат  $\beta_v : X \rightarrow \{0, 1\}$ ,  $\beta \in \mathcal{B}$
- 2)  $L_v, R_v$  — множества левых и правых дочерних вершин;
- 3)  $\lambda_v, \rho_v$  — левый и правый коэффициент.

Классификация вычисляется рекурсивно:

$a(x) = \text{sign ADT}(x, v_0)$ ,  $v_0$  — корень дерева;

$$\text{ADT}(x, v) = \begin{cases} \lambda_v + \sum_{t \in L_v} \text{ADT}(x, t), & \beta_v(x) = 0; \\ \rho_v + \sum_{t \in R_v} \text{ADT}(x, t), & \beta_v(x) = 1; \end{cases}$$

## Пример



Многократное применение рекуррентной формулы приводит к взвешенному голосованию конъюнкций:

$$\begin{aligned}
 ADT &= \bar{\beta}_0(\lambda_0 + \bar{\beta}_1\lambda_1 + \beta_1\rho_1) + \\
 &+ \beta_0(\rho_0 + \bar{\beta}_2\lambda_2 + \beta_2\rho_2 + \bar{\beta}_3\lambda_3 + \beta_3\rho_3) = \\
 &= (\bar{\beta}_0\lambda_0 + \beta_0\rho_0) + (\bar{\beta}_0\bar{\beta}_1\lambda_1 + \bar{\beta}_0\beta_1\rho_1) + \\
 &+ (\beta_0\bar{\beta}_2\lambda_2 + \beta_0\beta_2\rho_2) + (\beta_0\bar{\beta}_3\lambda_3 + \beta_0\beta_3\rho_3).
 \end{aligned}$$



## Бинарное решающее дерево является частным случаем ADT

... если наложить на ADT ограничение  $|L_v| = |R_v| \leq 1$   
и положить  $\lambda_v, \rho_v$  равными либо 0, либо меткам классов:

$$\lambda_v = \begin{cases} 0, & |L_v| = 1; \\ c_v \in Y, & |L_v| = 0; \end{cases}$$

$$\rho_v = \begin{cases} 0, & |R_v| = 1; \\ c_v \in Y, & |R_v| = 0. \end{cases}$$

Коэффициенты  $\lambda_v = c_v$  и  $\rho_v = c_v$  играют роль терминальных вершин.

В ADT терминальных вершин нет, все вершины равноправны.

## Обучение ADT методом бустинга

В общем случае

$$a(x) = \text{sign} \sum_{v \in V} \underbrace{\lambda_v K_v(x) \bar{\beta}_v(x) + \rho_v K_v(x) \beta_v(x)}_{\Gamma(x)},$$

где  $K_v(x)$  — конъюнкция всех предикатов на пути  $[v_0, v]$ .

**Идея наращивания дерева.** Создав вершину  $u$ , надо решить:

- 1) к какой вершине  $v$  её прицепить;
- 2) слева (тогда  $K_u = K_v \bar{\beta}_v$ ) или справа (тогда  $K_u = K_v \beta_v$ );
- 3) какой взять предикат  $\beta_u$  из семейства  $\mathcal{B}$ ;
- 4) как оптимизировать  $\lambda_u, \rho_u$ .

$$a(x) = \text{sign} \left( \Gamma(x) + K_u(x) (\lambda_u \bar{\beta}_u(x) + \rho_u \beta_u(x)) \right).$$

## Обучение ADT методом бустинга

Упростим обозначения:

$$K_i \equiv K_u(x_i), \quad \beta_i \equiv \beta_u(x_i), \quad \bar{\beta}_i \equiv \bar{\beta}_u(x_i), \quad \lambda \equiv \lambda_u, \quad \rho \equiv \rho_u.$$

Функционал числа ошибок после добавления вершины  $u$ :

$$Q = \sum_{i=1}^{\ell} \left[ \Gamma(x_i) y_i + \lambda K_i \bar{\beta}_i y_i + \rho K_i \beta_i y_i < 0 \right].$$

**AbaBoost:** экспоненциальная аппроксимация  $[M < 0] \leq e^{-M}$

$$Q \leq \tilde{Q} = \sum_{i=1}^{\ell} \underbrace{\exp(-\Gamma(x_i) y_i)}_{w_i} \exp(-K_i y_i (\lambda \bar{\beta}_i + \rho \beta_i)) \rightarrow \min_{\lambda, \rho, v, \beta}.$$

В отличие от AdaBoost, требуется найти 2 коэффициента  $\lambda, \rho$ .

## Основная теорема бустинга

### Теорема

Минимум функционала  $\tilde{Q}$  достигается при

$$\lambda^* = \frac{1}{2} \ln \frac{p(K\bar{\beta}^*)}{n(K\bar{\beta}^*)}; \quad \rho^* = \frac{1}{2} \ln \frac{p(K\beta^*)}{n(K\beta^*)};$$

$$\beta^* = \arg \min_{K, \beta} \left( W(\bar{K}) + 2\sqrt{p(K\bar{\beta})n(K\bar{\beta})} + 2\sqrt{p(K\beta)n(K\beta)} \right);$$

$$p(K\beta) = \sum_{i=1}^{\ell} w_i K_i \beta_i [y_i = +1];$$

$$n(K\beta) = \sum_{i=1}^{\ell} w_i K_i \beta_i [y_i = -1];$$

$$W(\bar{K}) = \sum_{i=1}^{\ell} w_i \bar{K}_i.$$

$p(K\beta)n(K\beta)$  — «число пар объектов» разных классов справа;  
 $W(\bar{K})$  — «число объектов», не дошедших до вершины  $v$ .

## Алгоритм AdaBoost для обучения ADT

**Вход:** выборка  $X^\ell$ ; семейство правил  $\mathcal{B}$ ; параметры  $T, \delta$ ;

**Выход:**  $ADT = \langle \beta_v, L_v, R_v, \lambda_v, \rho_v \rangle_{v \in V}$ ;

- 
- 1: инициализация:  $V := \emptyset$ ;  $w_i := 1, i = 1, \dots, \ell$ ;
  - 2: **для всех**  $t = 1, \dots, T$
  - 3: создать вершину  $u$ ;  $L_u := R_u := \emptyset$ ;
  - 4: минимизируя  $\tilde{Q}$ , найти: (1) предикат  $\beta_u \in \mathcal{B}$ , (2) к какой вершине  $v \in V$  присоединить  $u$ , (3) слева или справа;
  - 5: **если**  $u$  присоединяется слева **то**  $L_v := L_v \cup \{u\}$ ;
  - 6: **если**  $u$  присоединяется справа **то**  $R_v := R_v \cup \{u\}$ ;
  - 7:  $\lambda_u := \frac{1}{2} \ln \frac{p(K\bar{\beta}_u) + \delta}{n(K\bar{\beta}_u) + \delta}$ ;  $\rho_u := \frac{1}{2} \ln \frac{p(K\beta_u) + \delta}{n(K\beta_u) + \delta}$ ;
  - 8: **для всех**  $i = 1, \dots, \ell$  таких, что  $K_u(x_i) = 1$
  - 9:  $w_i := w_i \exp(-\lambda_u \bar{\beta}(x_i) y_i - \rho_u \beta(x_i) y_i)$ ;
  - 10: нормировка:  $Z := \frac{1}{\ell} \sum_{i=1}^{\ell} w_i$ ;  $w_i := w_i / Z, i = 1, \dots, \ell$ ;

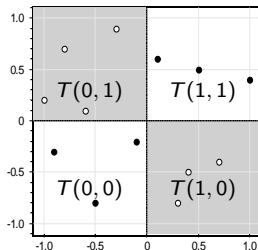
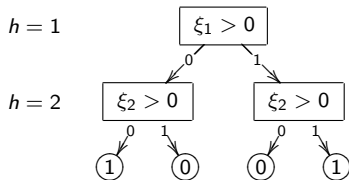
## Невнимательные решающие деревья — ODT (Oblivious Decision Tree) [1993]

**Решение проблемы фрагментации:** строится сбалансированное дерево высоты  $H$ ; для всех узлов уровня  $h$  условие ветвления  $\beta_h(x)$  одинаково; на уровне  $h$  ровно  $2^{h-1}$  вершин; пространство  $X$  делится на  $2^H$  ортантов.

Классификатор задаётся *таблицей решений*  $T: \{0, 1\}^H \rightarrow Y$ :

$$a(x) = T(\beta_1(x), \dots, \beta_H(x)).$$

**Пример:** задача XOR,  $H = 2$ .



## Алгоритм обучения ODT

**Вход:** выборка  $X^\ell$ ; семейство правил  $\mathcal{B}$ ; глубина дерева  $H$ ;

**Выход:** условия  $\beta_h$ ,  $h = 1, \dots, H$ ; таблица  $T: \{0, 1\}^H \rightarrow Y$ ;

- 1: для всех  $h = 1, \dots, H$
- 2: найти предикат с максимальной информативностью:  

$$\beta_h := \arg \max_{\beta \in \mathcal{B}} I(\beta_1, \dots, \beta_{h-1}, \beta);$$
- 3: для всех  $b \equiv (b_1, \dots, b_H) \in \{0, 1\}^H$
- 4: классификация по мажоритарному правилу:

$$T(b_1, \dots, b_H) := \arg \max_{c \in Y} \sum_{i=1}^{\ell} [y_i = c] \prod_{h=1}^H [\beta_h(x_i) = b_h];$$

$$I(\beta_1, \dots, \beta_h) = \sum_{c \in Y} h \left( \frac{P_c}{\ell} \right) - \sum_{b \in \{0,1\}^h} \frac{|X_b|}{\ell} \sum_{c \in Y} h \left( \frac{|X_b \cap X_c|}{|X_b|} \right);$$

$$X_b = \{x_i: \beta_s(x_i) = b_s, s = 1, \dots, h\}, \quad X^\ell = \bigsqcup_{b \in \{0,1\}^h} X_b.$$

## Резюме в конце лекции

- Преимущества решающих деревьев:
  - интерпретируемость,
  - допускаются разнотипные данные,
  - возможность обхода пропусков;
- Недостатки решающих деревьев:
  - переобучение,
  - фрагментация,
  - неустойчивость к шуму, составу выборки, критерию;
- Post-pruning уменьшает переобучение;
- ADT, ODT снижают фрагментацию;
- При наличии пропусков в данных оценивают апостериорные вероятности классов.
- Голосование деревьев — один из лучших методов обучения.