

Московский энергетический институт

Кафедра Вычислительных машин систем и сетей

КУРС ПРОБЛЕМЫ ОРГАНИЗАЦИИ ВЫЧИСЛЕНИЙ

Лекция №2 на тему :

«Введение в машинное обучение»

Задача классификации (обучение с у

Задача восстановления зависимости $y: X \to Y, |Y| < \infty$ по точкам обучающей выборки $(x_i, y_i), i = 1, \dots, \ell$.

Дано: векторы $x_i = (x_i^1, \dots, x_i^n)$ — объекты обучающей выборки, $y_i = y(x_i)$ — классификации, ответы учителя, $i = 1, \dots, \ell$:

$$\begin{pmatrix} x_1^1 & \dots & x_1^n \\ \dots & \dots & \dots \\ x_\ell^1 & \dots & x_\ell^n \end{pmatrix} \xrightarrow{y^*} \begin{pmatrix} y_1 \\ \dots \\ y_\ell \end{pmatrix}$$

Найти: функцию a(x), способную классифицировать объекты произвольной тестовой выборки $\tilde{x}_i = (\tilde{x}_i^1, \dots, \tilde{x}_i^n), i = 1, \dots, k$:

$$\begin{pmatrix} \tilde{x}_1^1 & \dots & \tilde{x}_1^n \\ \dots & \dots & \dots \\ \tilde{x}_k^1 & \dots & \tilde{x}_k^n \end{pmatrix} \stackrel{\mathbf{a}^7}{\longrightarrow} \begin{pmatrix} \mathbf{a}(\tilde{x}_1) \\ \dots \\ \mathbf{a}(\tilde{x}_k) \end{pmatrix}$$

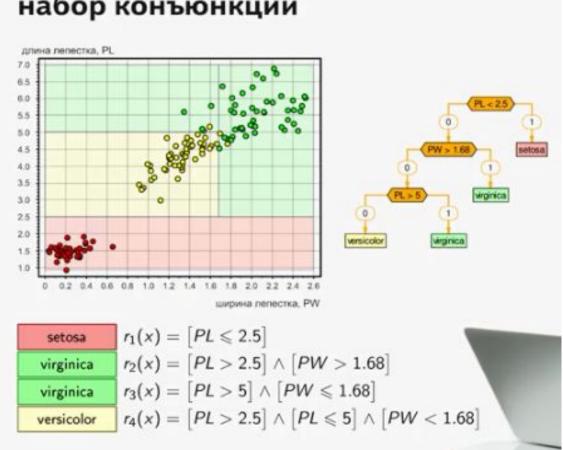
Определение бинарного решающего дерева

```
Бинарное решающее дерево — алгоритм классификации a(x),
  задающийся бинарным деревом:
 1) \forall v \in V_{\mathsf{внутр}} \to \mathsf{предикат} \ \beta_v : X \to \{0,1\}, \ \beta_v \in \mathscr{B},
  2) \forall v \in V_{\mathsf{лист}} \to \mathsf{имя} \mathsf{класса} c_v \in Y,
  где \mathscr{B} — множество бинарных признаков или предикатов
  (например, вида \beta(x) = [x^j \ge \theta_i], x^j \in \mathbb{R})
1: v := v_0;
2: пока v \in V_{\text{внутр}}
3: если \beta_{\nu}(x) = 1 то
    переход вправо: v := R_v;
   иначе
         переход влево: v := L_v;
7: вернуть c<sub>v</sub>.
```

Определение бинарного решающего дерева

```
Бинарное решающее дерево — алгоритм классификации a(x),
  задающийся бинарным деревом:
 1) \forall v \in V_{\mathtt{внутр}} \to \mathsf{предикат} \ \beta_v : X \to \{0,1\}, \ \ \beta_v \in \mathscr{B},
  2) \forall v \in V_{\mathsf{лист}} \to \mathsf{имя} \mathsf{класса} c_v \in Y,
  где \mathscr{B} — множество бинарных признаков или предикатов
  (например, вида \beta(x) = [x^j \ge \theta_i], x^j \in \mathbb{R})
1: v := v_0;
2: пока v \in V_{\text{внутр}}
3: если \beta_{\nu}(x) = 1 то
    переход вправо: v := R_v;
   иначе
         переход влево: v := L_v;
7: вернуть c<sub>v</sub>.
```

Решающее дерево — покрывающий набор конъюнкций



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
     вернуть новый лист v, c_v := c;
4: найти предикат с максимальной информативностью:
  \beta := \arg \max_{\beta \in \mathscr{B}} I(\beta, U);
5: разбить выборку на две части U=U_0\sqcup 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 := Мажоритарный класс(<math>U);
8: создать новую внутреннюю вершину v: \beta_v := \beta;
   построить левое поддерево: L_{\nu} := \text{LearnID3 } (U_0);
   построить правое поддерево: R_{\nu} := \text{LearnID3 } (U_1);
9: вернуть v;
```

Варианты критериев ветвления

1. Критерий Джини:

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

2. D-критерий В.И.Донского:

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

3. Энтропийный критерий:

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

где
$$h(z) \equiv -z \log_2 z$$
, $P_c(X^\ell) = \#\{x_i \colon y_i = c\}$, $p_c(X^\ell) = \#\{x_i \colon y_i = c \text{ и } \beta(x_i) = 1\}$, $p(X^\ell) = \#\{x_i \colon \beta(x_i) = 1\}$.

Решающие деревья ID3: достоинства

Достоинства:

- Интерпретируемость и простота классификации.
- У Гибкость: можно варьировать множество В.
- Допустимы разнотипные данные и данные с пропусками.
-) Трудоёмкость линейна по длине выборки $O(|\mathscr{B}|h\ell)$.
- 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 β_v , c_v .
- Высокая чувствительность к шуму, к составу выборки, к критерию информативности.

Усечение дерева (pruning). Алгоритм C4.5

```
X^k — независимая контрольная выборка, k \approx 0.5\ell.
 1: для всех v ∈ V_{внутр}
      S_v := подмножество объектов X^k, дошедших до v;
    если S_{\nu} = \emptyset то
        вернуть новый лист v, c_v := Мажоритарный класс(U);
    число ошибок при классификации S_{\nu} 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.
      в зависимости от того, какое из них минимально:
        сохранить поддерево v;
        заменить поддерево \nu поддеревом L_{\nu};
        заменить поддерево \nu поддеревом R_{\nu};
        заменить поддерево v листом, c_v := \arg\min_{c \in V} r_c(v).
```