

Задача классификации (обучение с учителем)

Задача восстановления зависимости $y: X \rightarrow Y$, $|Y| < \infty$
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* $\tilde{x}_i = (\tilde{x}_i^1, \dots, \tilde{x}_i^n)$, $i = 1, \dots, k$:

$$\begin{pmatrix} \tilde{x}_1^1 & \dots & \tilde{x}_1^n \\ \dots & \dots & \dots \\ \tilde{x}_k^1 & \dots & \tilde{x}_k^n \end{pmatrix} \xrightarrow{a?} \begin{pmatrix} a(\tilde{x}_1) \\ \dots \\ a(\tilde{x}_k) \end{pmatrix}$$

Определение бинарного решающего дерева

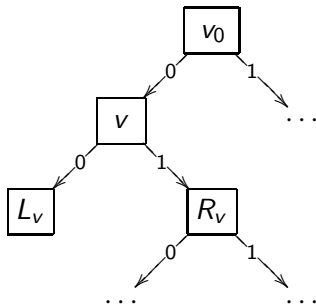
Бинарное решающее дерево — алгоритм классификации $a(x)$, задающийся бинарным деревом:

1) $\forall v \in V_{\text{внутр}} \rightarrow$ предикат $\beta_v : X \rightarrow \{0, 1\}$, $\beta_v \in \mathcal{B}$,

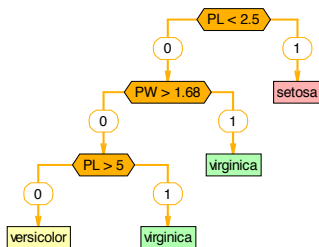
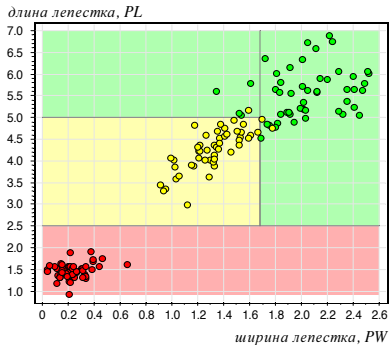
2) $\forall v \in V_{\text{лист}} \rightarrow$ имя класса $c_v \in Y$,

где \mathcal{B} — множество бинарных признаков или предикатов (например, вида $\beta(x) = [x^j \geq \theta_j]$, $x^j \in \mathbb{R}$)

```
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  $c_v$ .
```

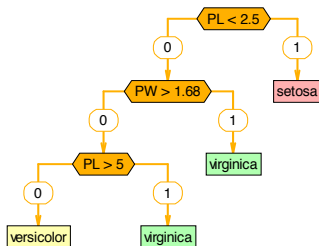
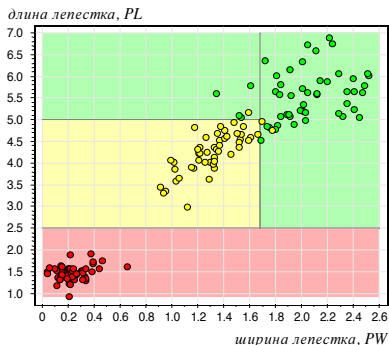


Задача Фишера о классификации цветков ириса на 3 класса, в выборке по 50 объектов каждого класса, 4 признака.



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

Решающее дерево → покрывающий набор конъюнкций



setosa

$$r_1(x) = [PL \leq 2.5]$$

virginica

$$r_2(x) = [PL > 2.5] \wedge [PW > 1.68]$$

virginica

$$r_3(x) = [PL > 5] \wedge [PW \leq 1.68]$$

versicolor

$$r_4(x) = [PL > 2.5] \wedge [PL \leq 5] \wedge [PW < 1.68]$$