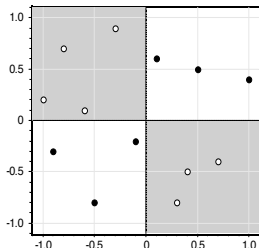
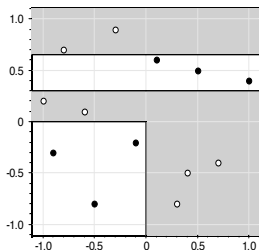
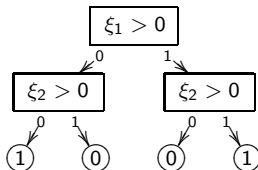


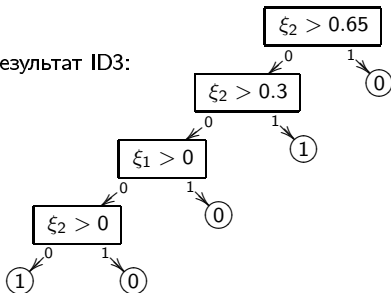
Жадный ID3 переусложняет структуру дерева



Оптимальное дерево для задачи XOR:



Результат ID3:



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

- 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, $c_v := \arg \min_{c \in Y} r_c(v)$.

Обобщение на случай регрессии: $Y = \mathbb{R}$, $c_v \in \mathbb{R}$

Пусть U_v — множество объектов x_i , дошедших до вершины v

Значения в терминальных вершинах — МНК-решение:

$$c_v := \hat{y}(U_v) = \frac{1}{|U_v|} \sum_{x_i \in U_v} y_i$$

Критерий информативности — среднеквадратичная ошибка

$$I(\beta, U_v) = \sum_{x_i \in U_v} (\hat{y}_i(\beta) - y_i)^2,$$

где $\hat{y}_i(\beta) = \beta(x_i)\hat{y}(U_{v1}) + (1 - \beta(x_i))\hat{y}(U_{v0})$

— прогноз после ветвления β и разбиения $U_v = U_{v0} \sqcup U_{v1}$

Среднеквадратичная ошибка со штрафом за сложность дерева

$$C_{\alpha} = \sum_{x_j=1}^{\ell} (\hat{y}_i - y_i)^2 + \alpha |V_{\text{лист}}| \rightarrow \min$$

При увеличении α дерево последовательно упрощается.

Причём последовательность вложенных деревьев единственна.

Из этой последовательности выбирается дерево с минимальной ошибкой на тестовой выборке (Hold-Out).

Для случая классификации используется аналогичная стратегия усечения, с критерием Джини.

- Преимущества решающих деревьев:
 - интерпретируемость,
 - допускаются разнотипные данные,
 - возможность обхода пропусков;
- Недостатки решающих деревьев:
 - переобучение,
 - фрагментация,
 - неустойчивость к шуму, составу выборки, критерию;
- Способы устранения этих недостатков:
 - редукция,
 - композиции (леса) деревьев.