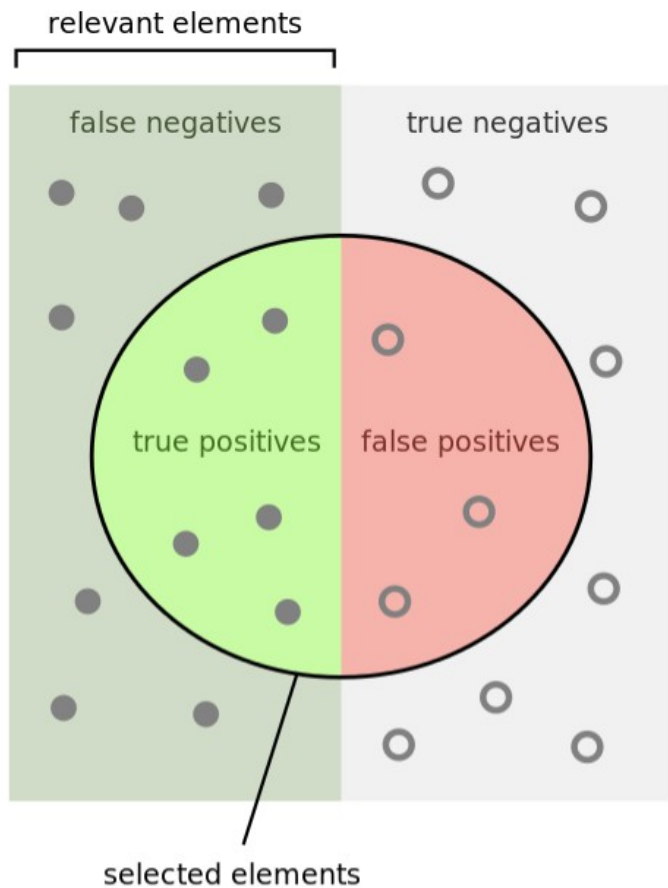


Деревья

Рыжиков Артём, НИУ ВШЭ,
17.11.2018

Offtop. Метрики качества



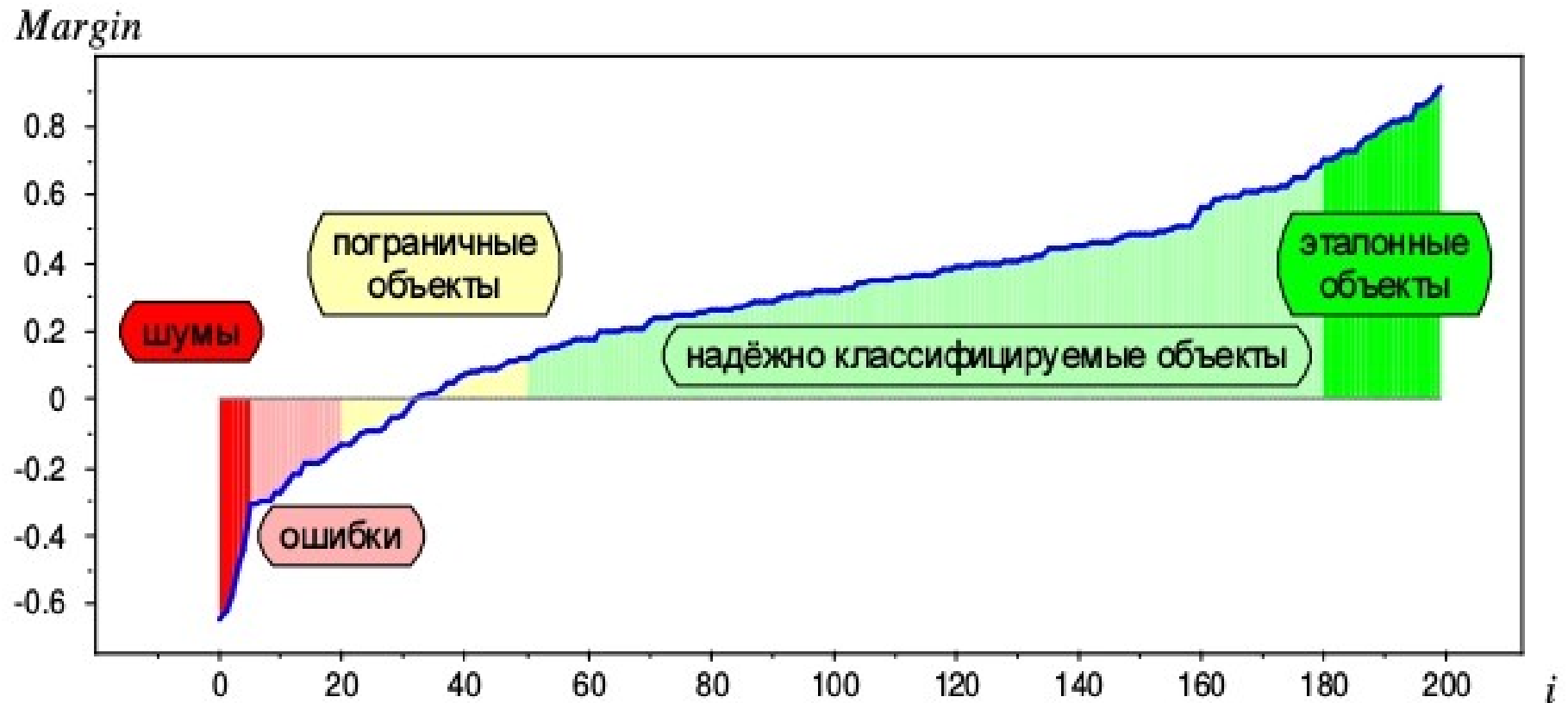
How many selected
items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

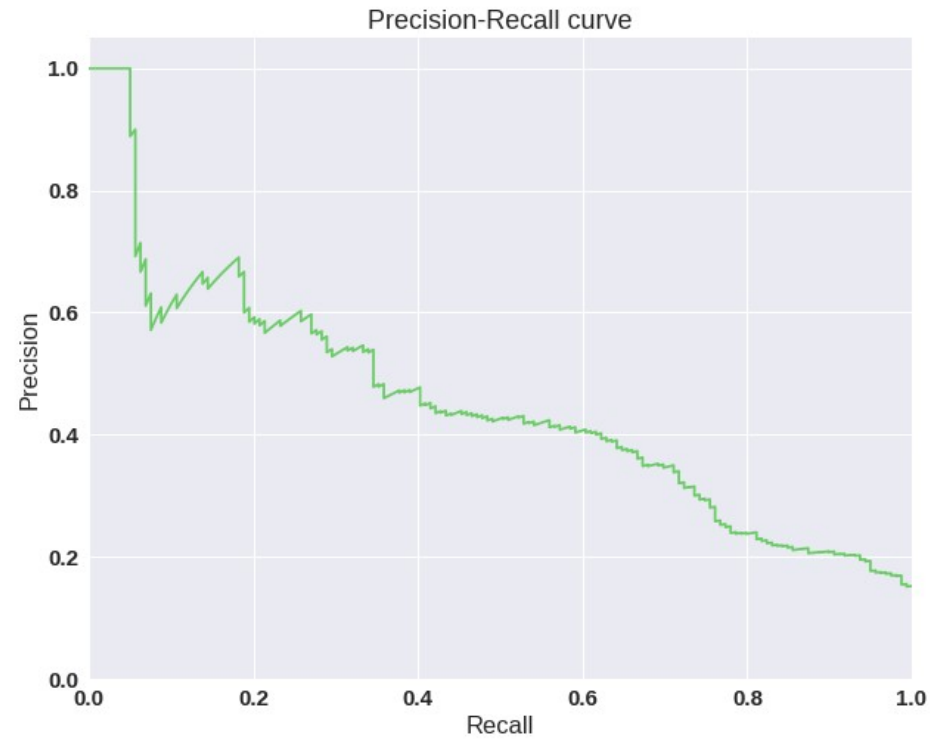
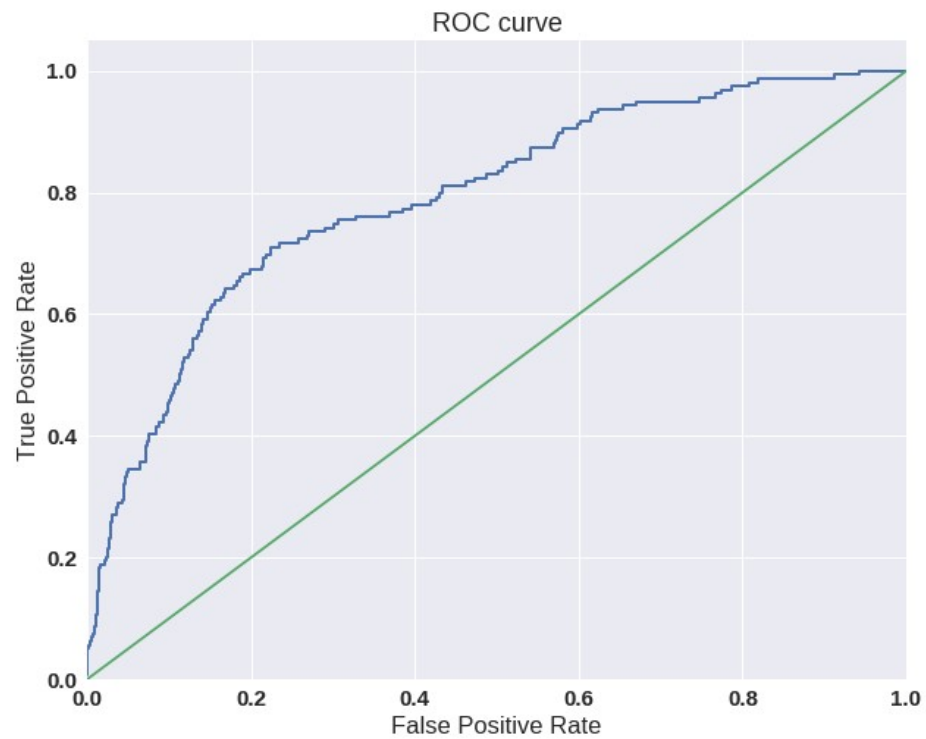
How many relevant
items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Threshold



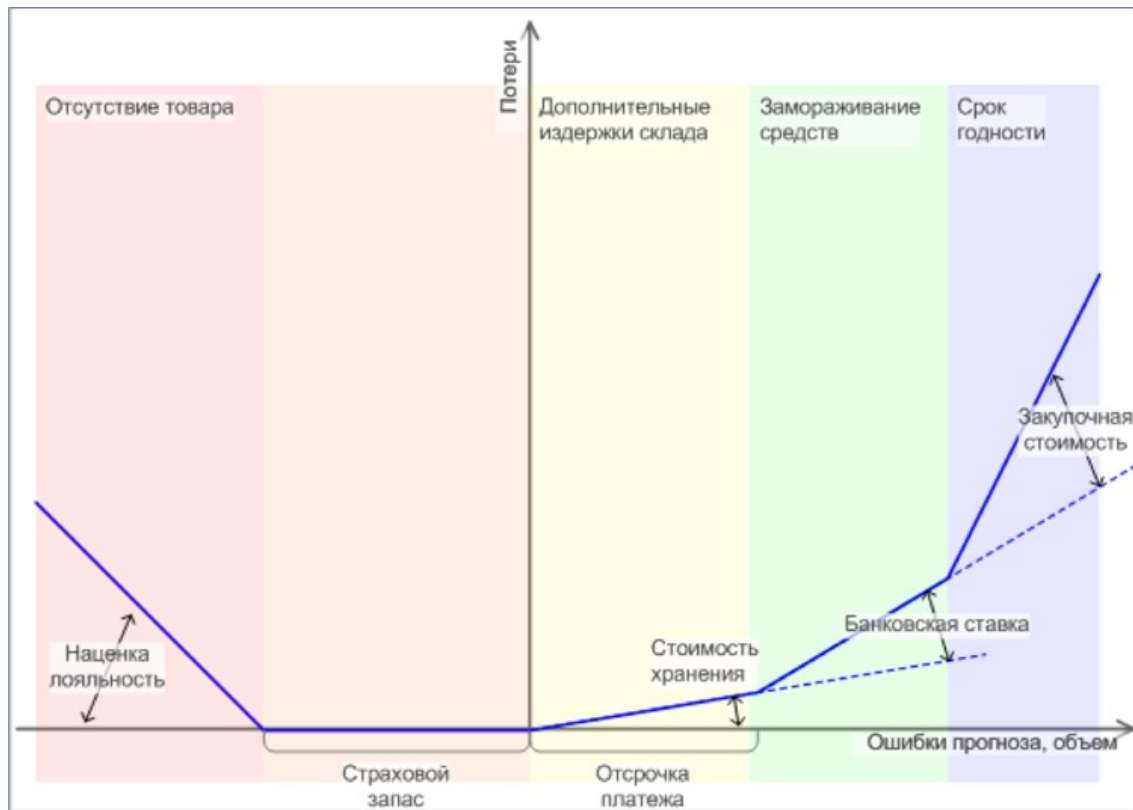
Метрики



Ссылки

- [Habr](#)
- [Подробнее](#) про различие между ROC AUC и PR AUC

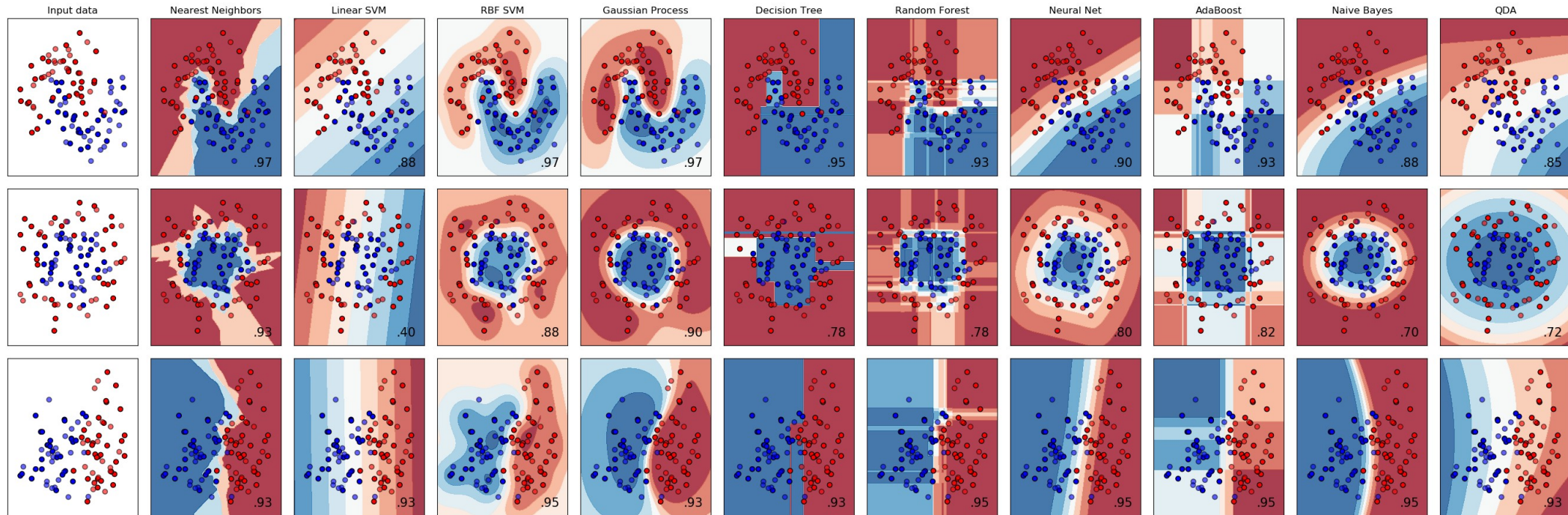
Метрики в задаче регрессии



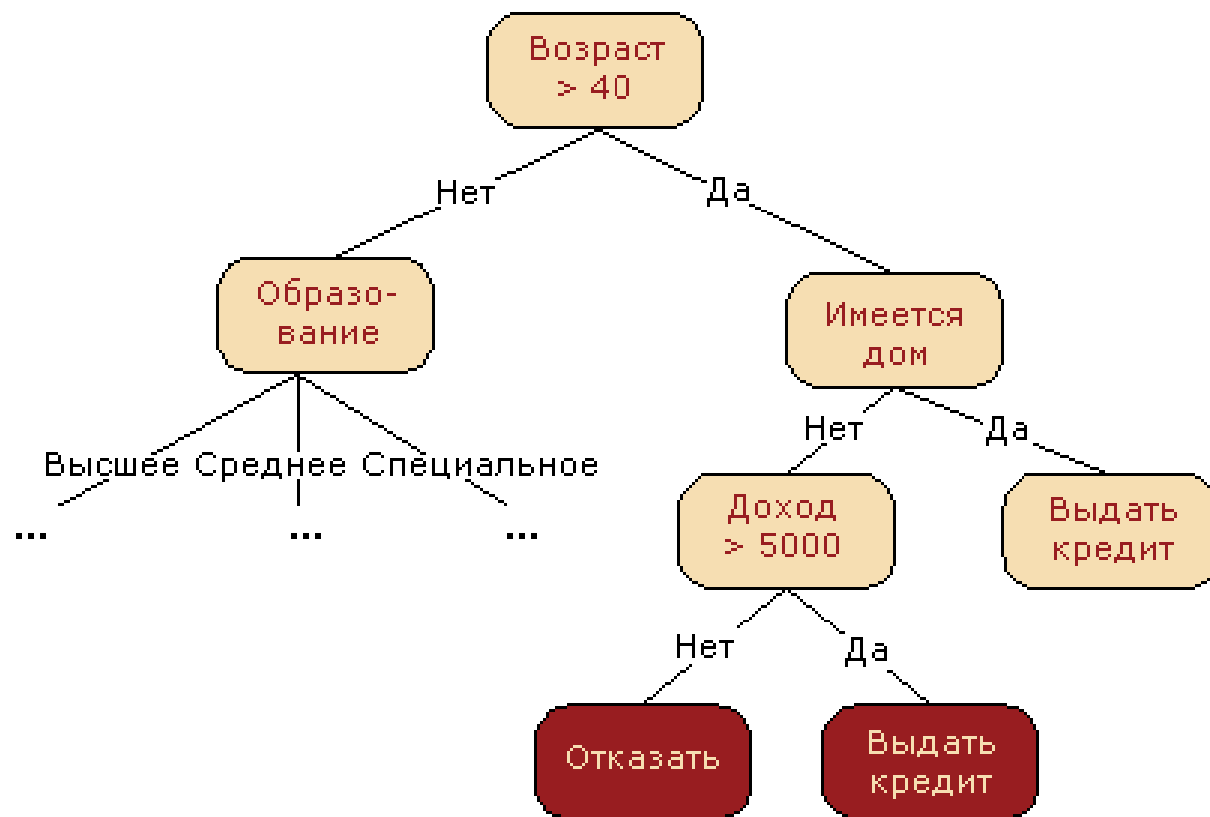
Ссылки

- Лекции [Евгения Соколова](#)

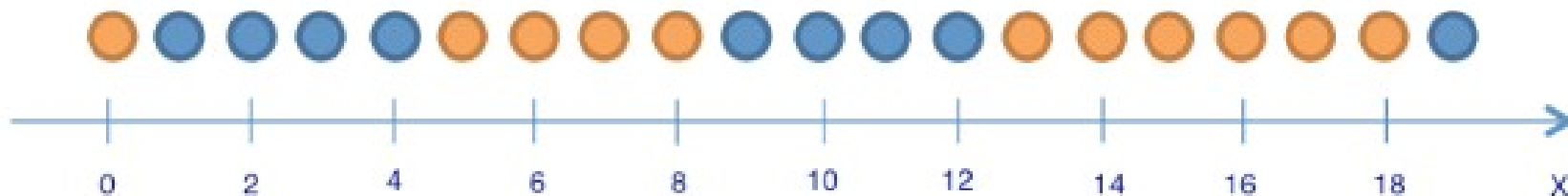
Логические методы классификации



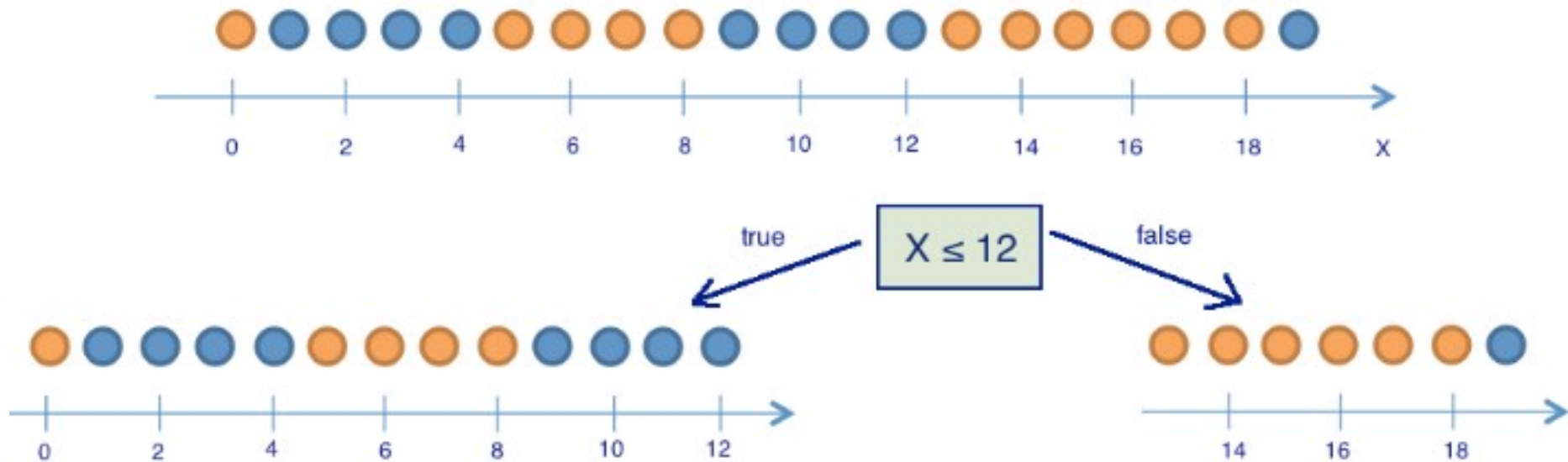
Дерево принятия решений



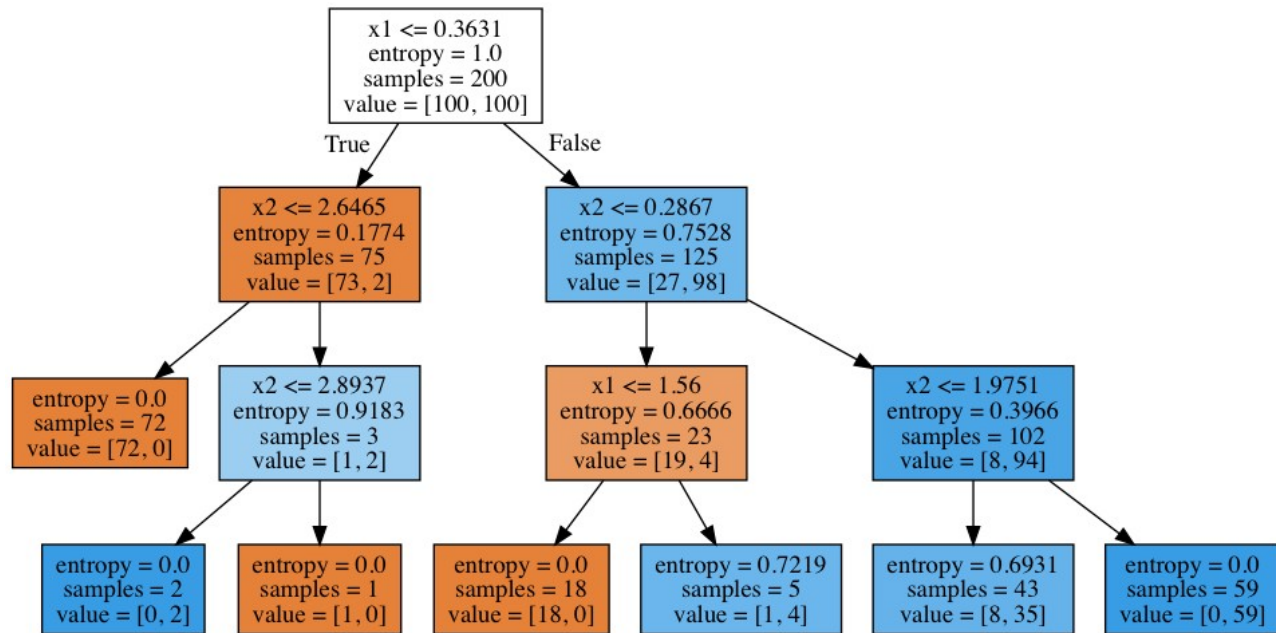
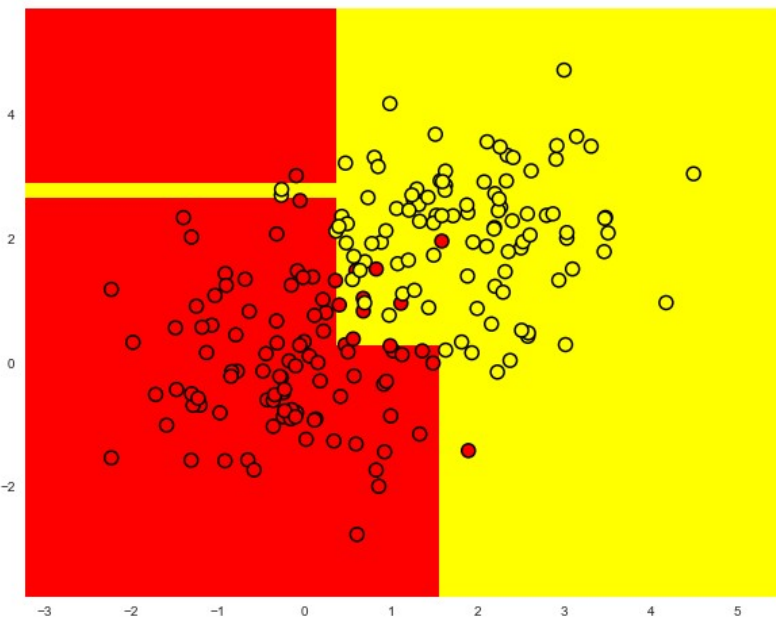
Информационный прирост



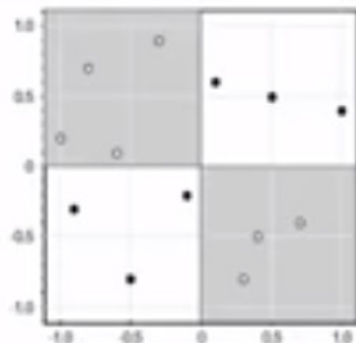
Информационный прирост



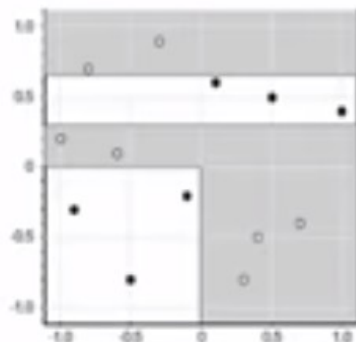
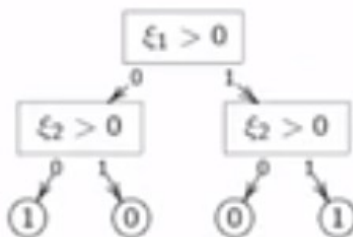
Пример дерева



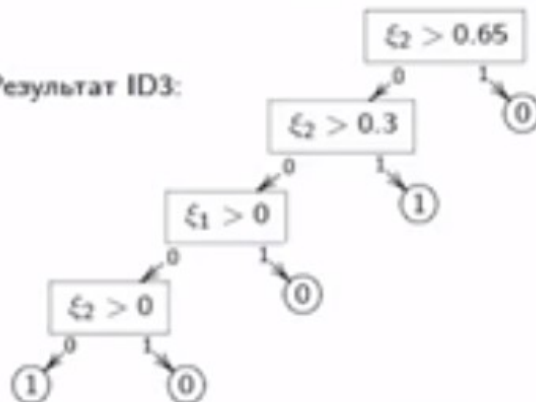
Переобучаемость ID3



Оптимальное дерево для задачи XOR:



Результат ID3:



C4.5, CART

Редукция дерева («стрижка», pruning: C4.5, CART)

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

Ссылки

- [Видеолекции](#)
- [Habr](#)