

Решающие деревья

Лекция 4

Тест



Повторение

X – множество объектов (их признаковое описание)

Y – множество истинных ответов

\hat{Y} – множество предсказанных ответов. Получаем по формуле:

$$\hat{Y} = f(X),$$

где $f(X)$ – модель машинного обучения

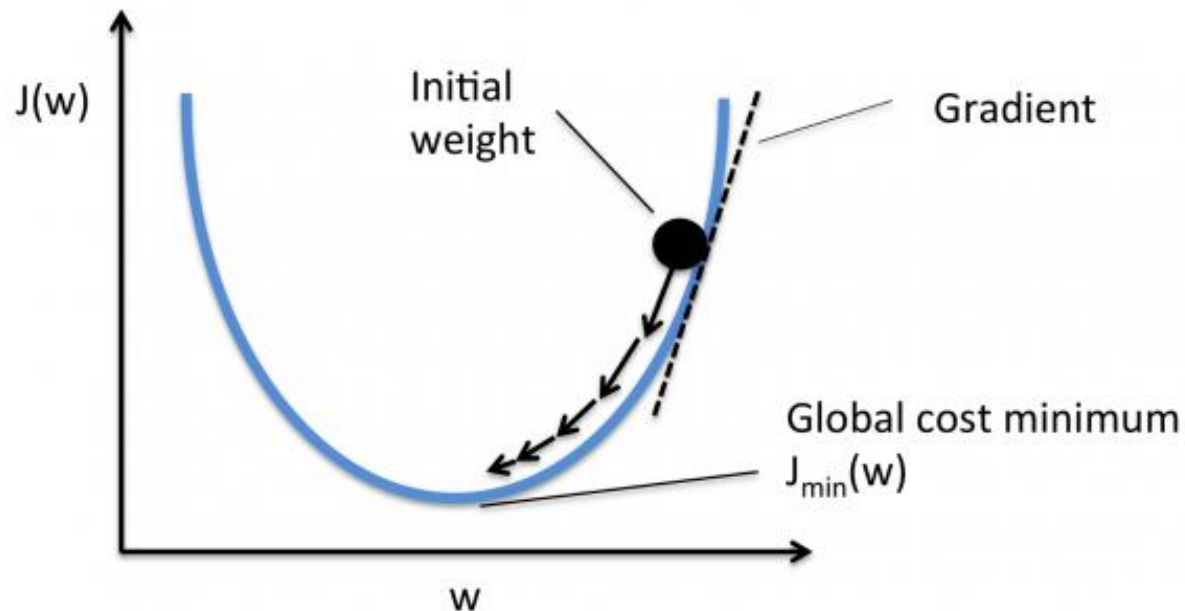
Если модель линейная, то

$$\hat{Y} = X^T W$$

Повторение

Обновление весов модели происходит с учетом **антиградиента** функции потерь.

$$W = W - \eta \cdot \nabla_W Q(W)$$

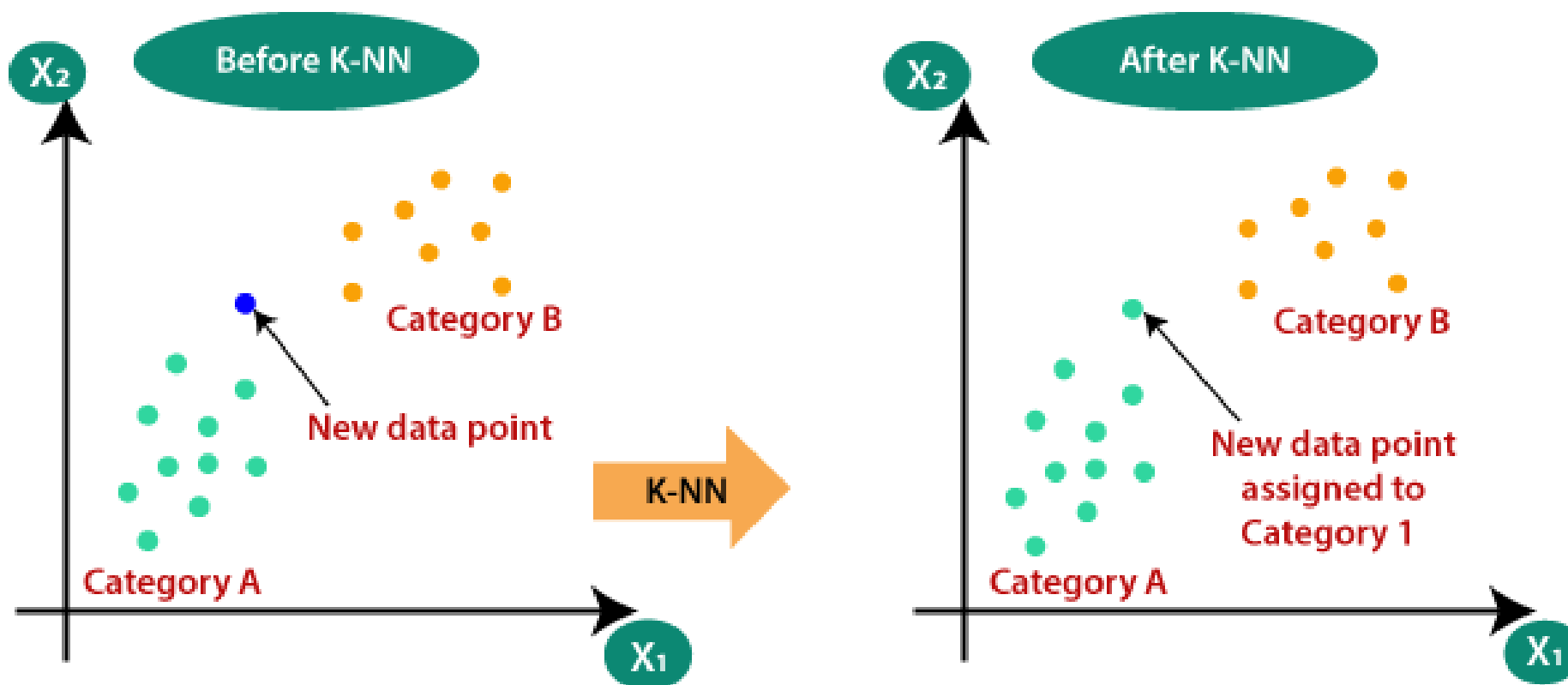


Новый мир без градиента

Есть модели, где градиент не нужен:

- К-ближайших соседей
- Наивный байесовский классификатор (регрессор)
- Решающее дерево

К-ближайших соседей



К-ближайших соседей. Расстояние

Варианты расстояний:

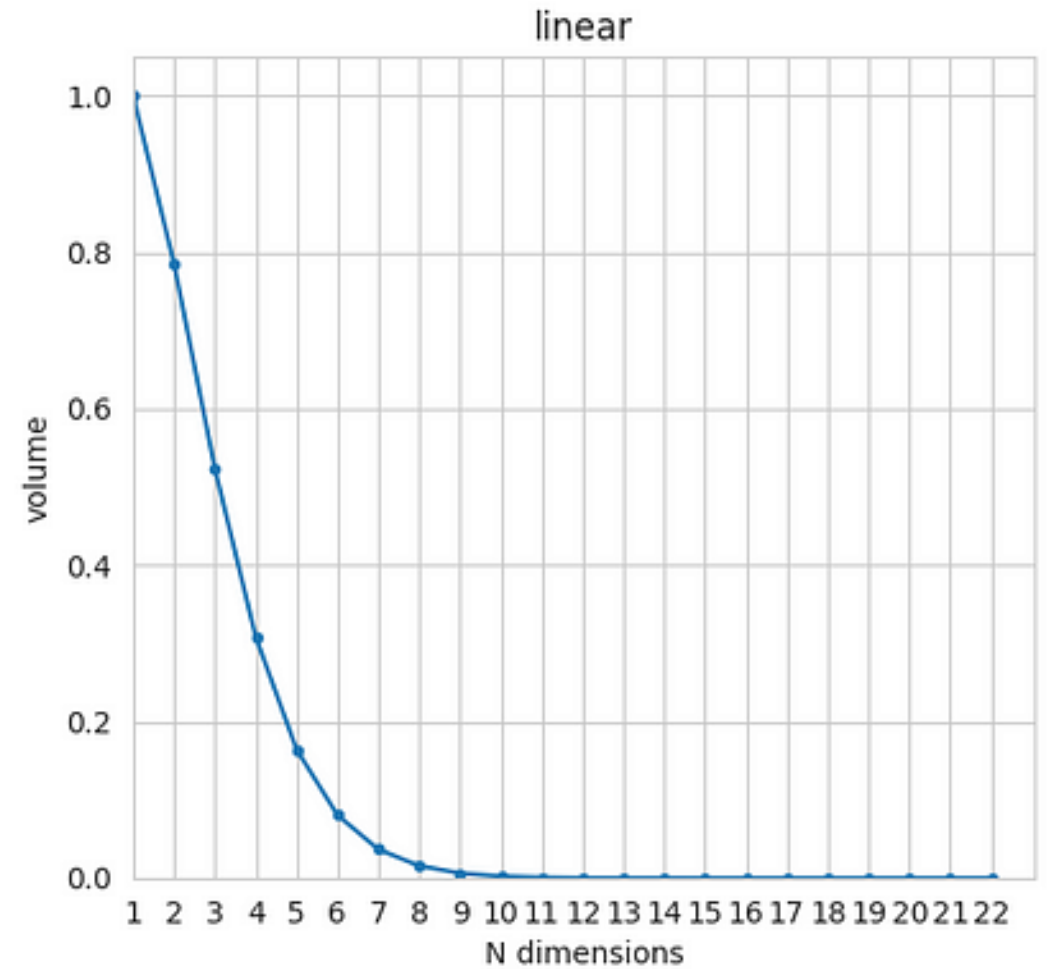
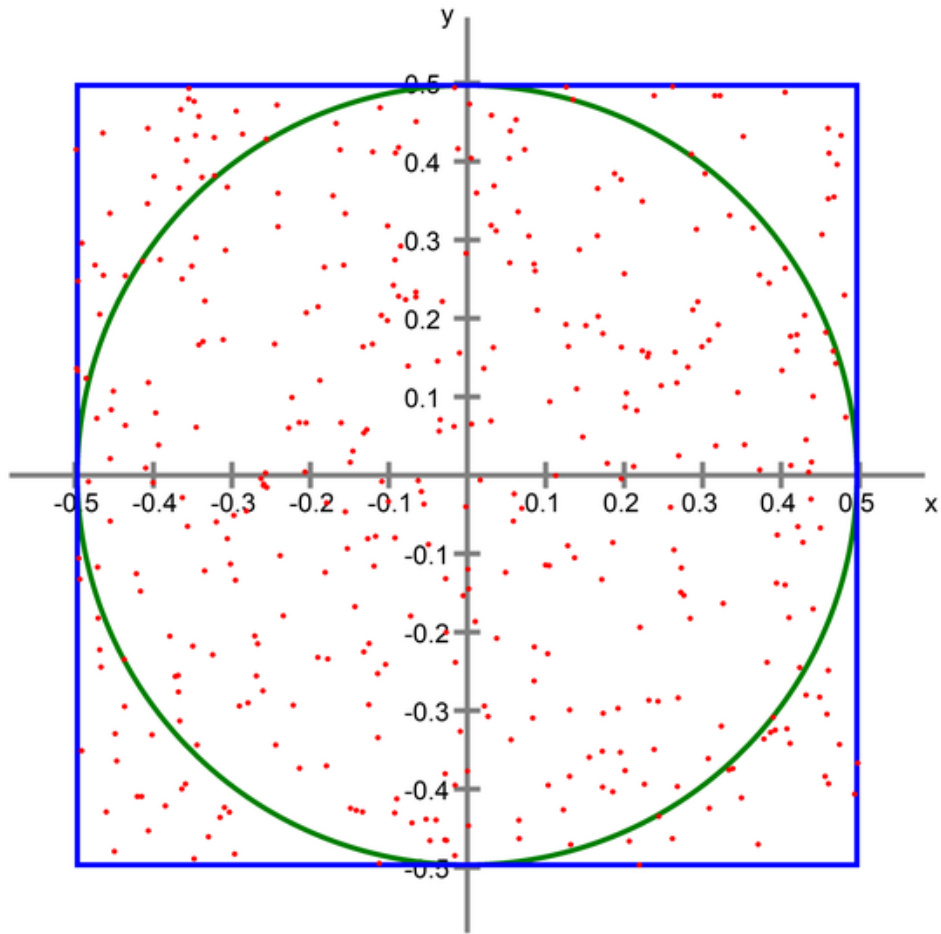
Евклидово (L2): $Q(x^k, x^m) = \sum_{i=1}^n \sqrt{(x_i^k - x_i^m)^2}$

Манхэттенское расстояние (L1): $Q(x^k, x^m) = \sum_{i=1}^n |x_i^k - x_i^m|$

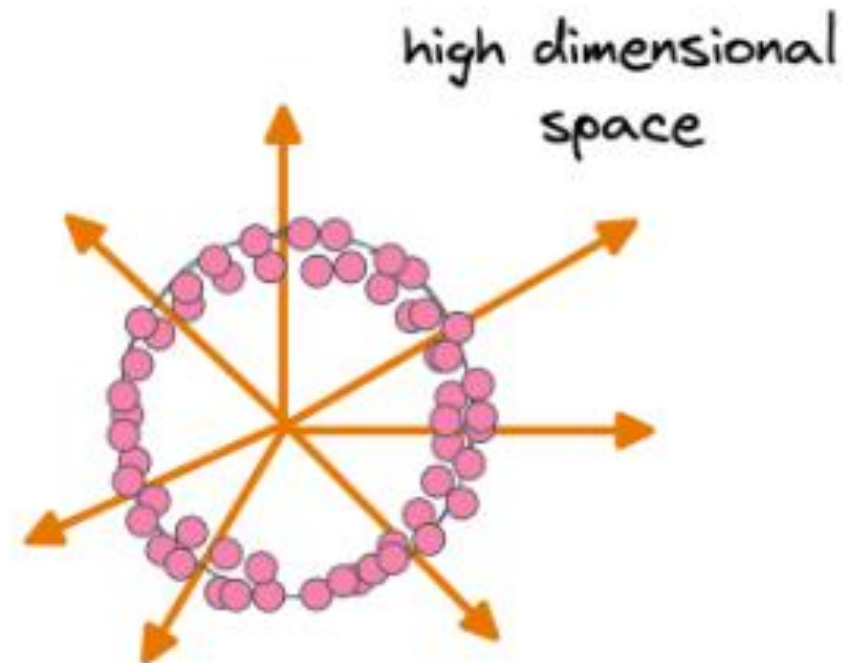
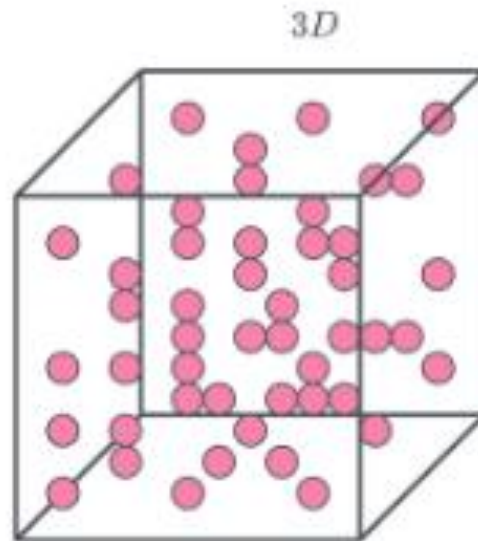
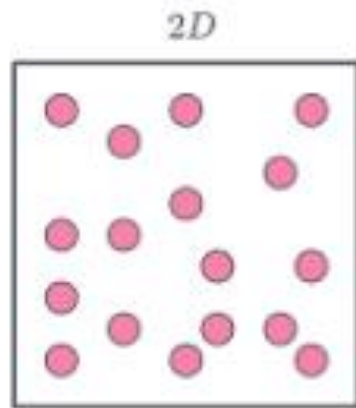
Косинусное расстояние: $Q(x^k, x^m) = \frac{\sum_{i=1}^n x_i^k \cdot x_i^m}{\sum_{i=1}^n \sqrt{(x_i^k)^2} \cdot \sum_{i=1}^n \sqrt{(x_i^m)^2}}$

Необходимо масштабировать признаки при использовании KNN.

Проклятие размерности



Проклятие размерности



Вопрос на засыпку

В выборке 1000 человек, из них 400 мужчин, 600 женщин.

Целевых событий ($y = 1$) 100 штук, 80 из которых у мужчин.

Какова вероятность целевого событий при условии, что перед нами мужчина?

Наивный байесовский классификатор

Теорема Байеса:

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$$

где $P(A)$ – априорная вероятность события A ,

$P(B)$ – априорная вероятность события B ,

$P(B|A)$ – условная вероятность наступления события B при наступлении события A (правдоподобие)

Наивный байесовский классификатор

Например:

В выборке 1000 человек, из них 400 мужчин, 600 женщин.

Целевых событий ($y = 1$) 100 штук, 80 из которых у мужчин.

Какова вероятность целевого событий при условии, что перед нами мужчина?

Перепишем в терминах машинного обучения:

$$P(y = 1 | X = \text{«Мужчина»}) = \frac{P(y = 1) \cdot P(X = \text{«Мужчина»} | y = 1)}{P(X = \text{«Мужчина»})}$$

Наивный байесовский классификатор

1000 человек, 400 М, 600 Ж. 100 «1», из них 80 М, 20 Ж.

$$P(y = 1) = \frac{100}{1000} = 0.1$$

$$P(X = \text{«Мужчина»}) = \frac{400}{1000} = 0.4$$

$$P(X = \text{«Мужчина»} | y = 1) = \frac{80}{100} = 0.8$$

Наивный байесовский классификатор

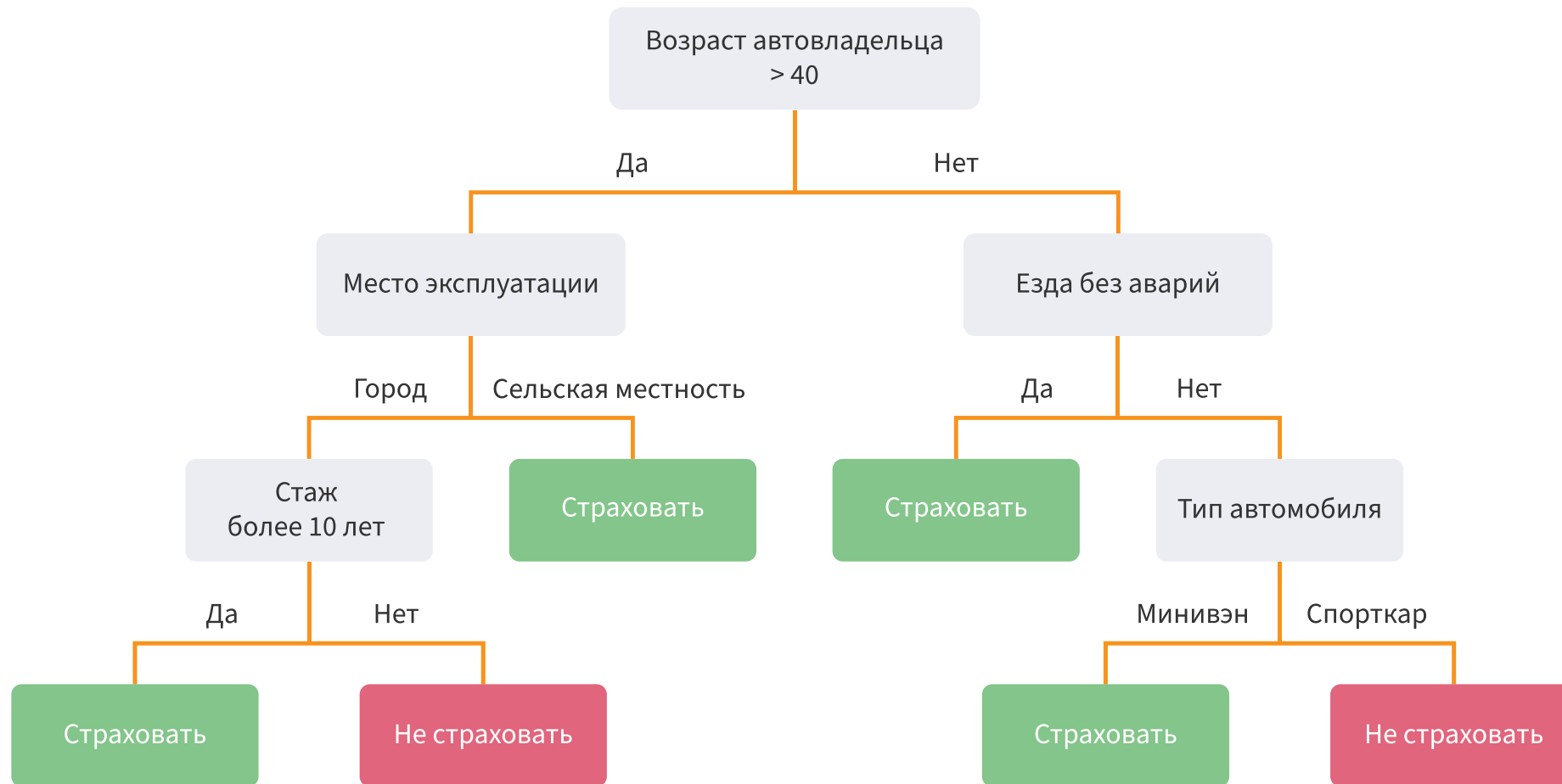
Тогда

$$P(y = 1|X = \text{«Мужчина»}) = \frac{0.1 \cdot 0.8}{0.4} = 0.2 = 20\%$$

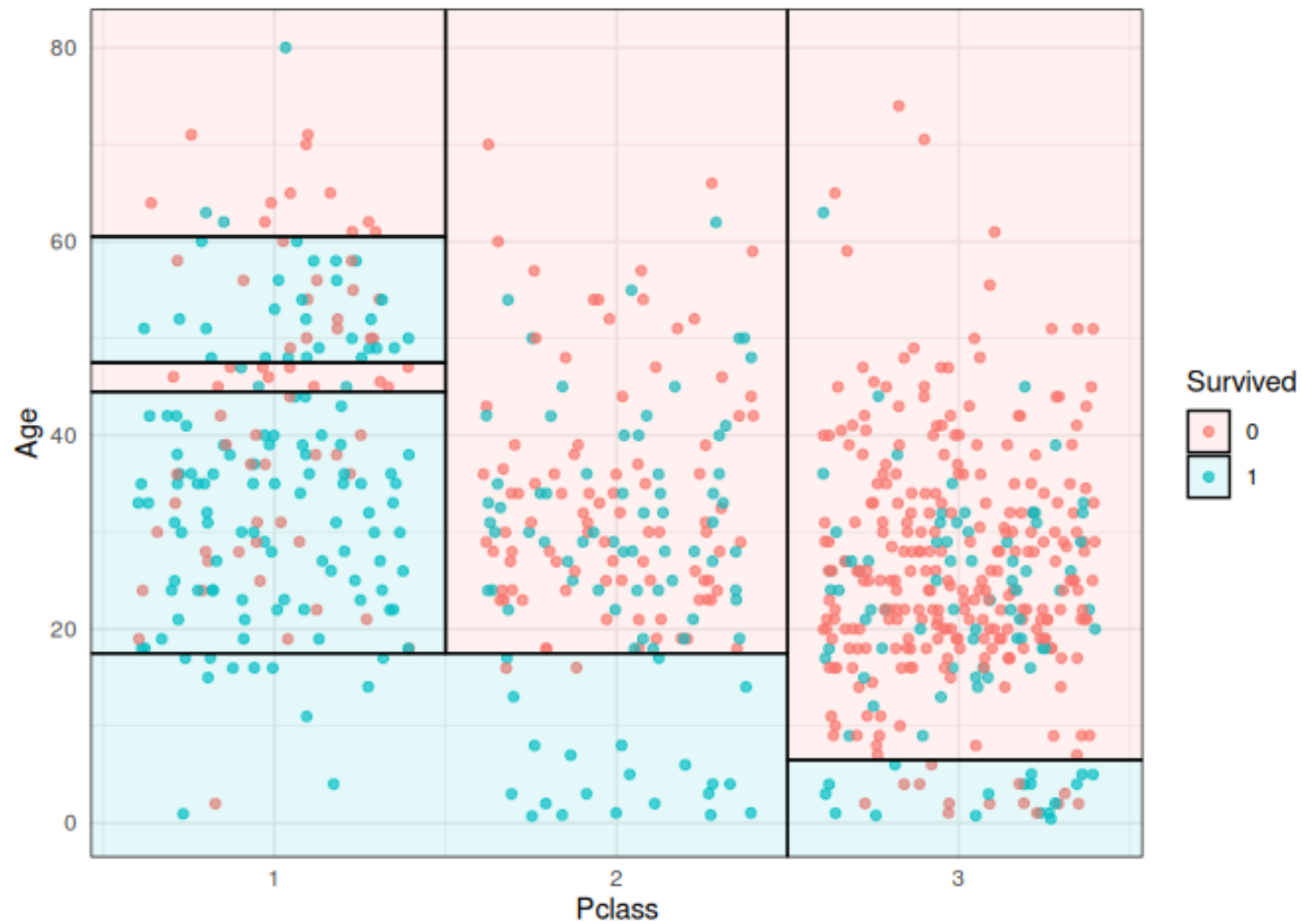
Если признаков несколько, то:

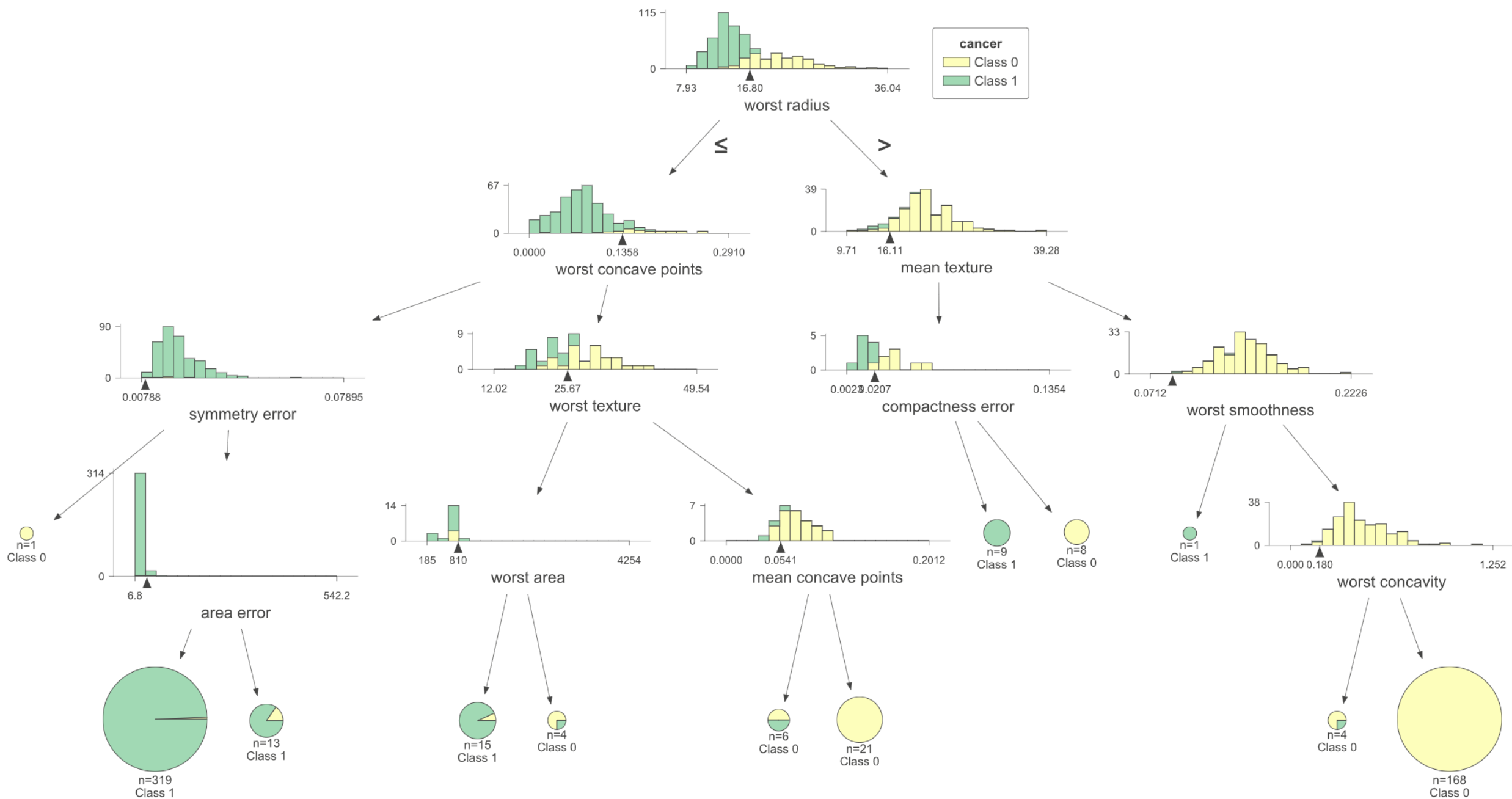
$$P(y = 1|X_1, X_2) = \frac{P(y = 1) \cdot P(X_1|y = 1) \cdot P(X_2|y = 1)}{P(X_1, X_2)}$$

Решающее дерево



Решающее дерево





Решающее дерево. Критерий информативности

Хотим для каждого потенциального разбиения знать, насколько оно хорошо.

Опишем функционал качества $Q(R)$ по формуле:

$$Q(R, X_i, t) = H(R_{head}) - \frac{size_{left}}{size_{head}} H(R_{left}) - \frac{size_{right}}{size_{head}} H(R_{right})$$

где $H(R)$ -- критерий информативности (impurity criterion)

Решающее дерево. Регрессия

Для регрессии введем критерий информативности:

$$H(R) = \frac{1}{N} \sum_{i=0}^N (y_i - c)^2 \rightarrow \min$$

$$H(R) = \frac{1}{N} \sum_{i=0}^N \left(y_i - \frac{1}{N} \sum_{i=0}^N y_i \right)^2 \rightarrow \min$$

Решающее дерево. Классификация

1. Оцениваем ошибку классификации

$$H(R) = \frac{1}{N} \sum_{i=0}^N [y_i \neq c] \rightarrow \min$$

В листе дерево всегда предсказывает наиболее часто встречающийся класс, поэтому:

$$H(R) = \frac{1}{N} \sum_{i=0}^N [y_i \neq c] = 1 - p_k$$

Решающее дерево. Классификация

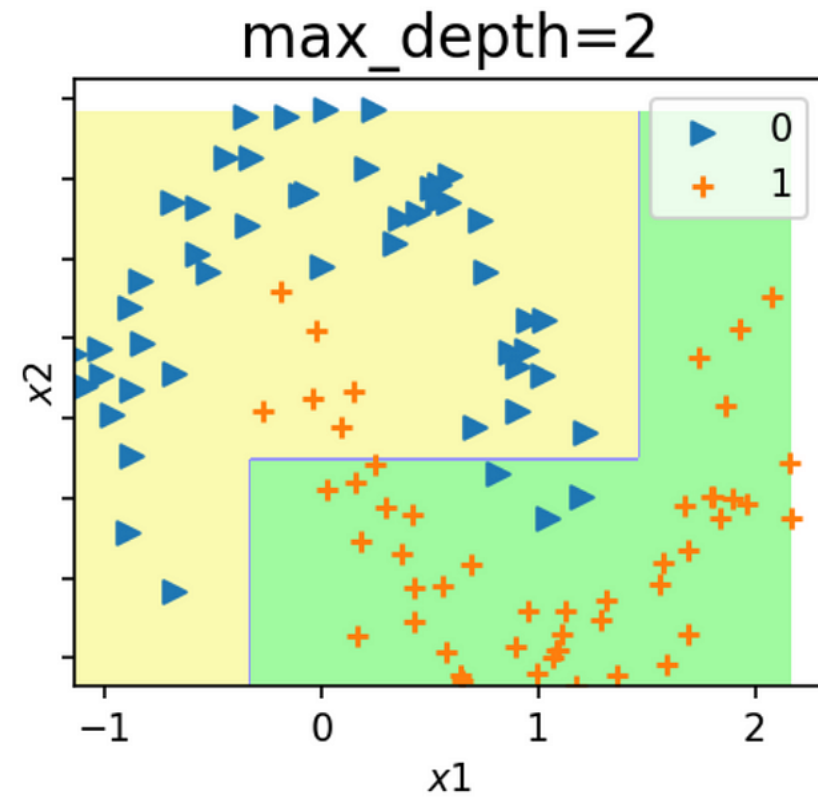
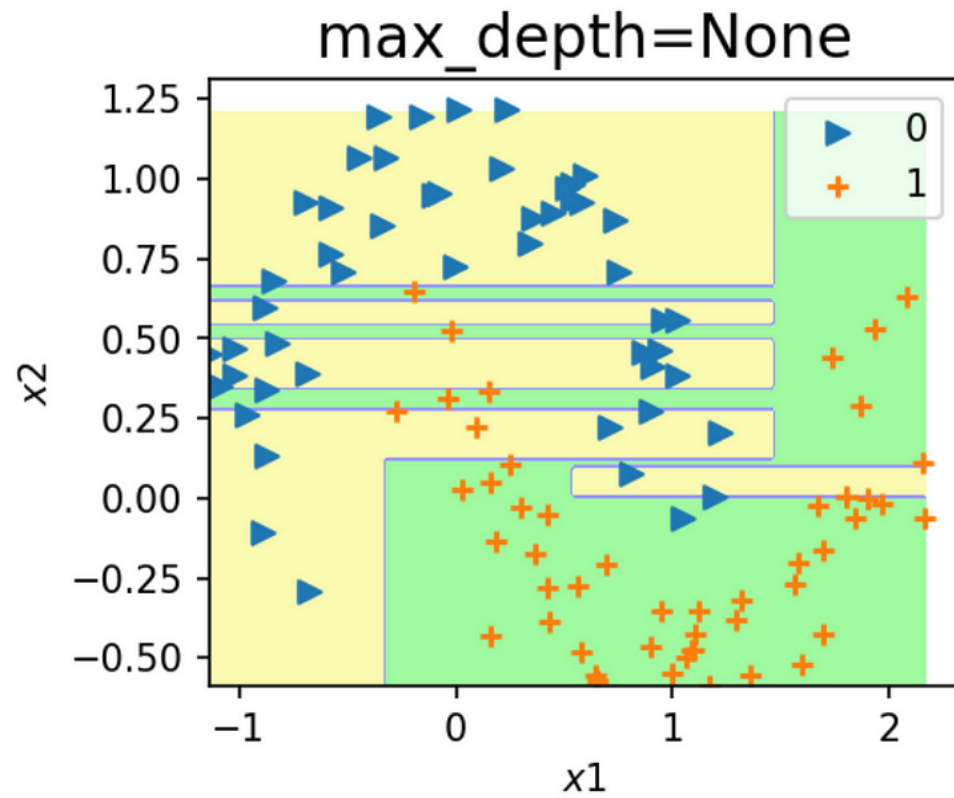
2. Критерий Джини

$$H(R) = - \sum_{k=1}^K p_k (1 - p_k) \rightarrow \min$$

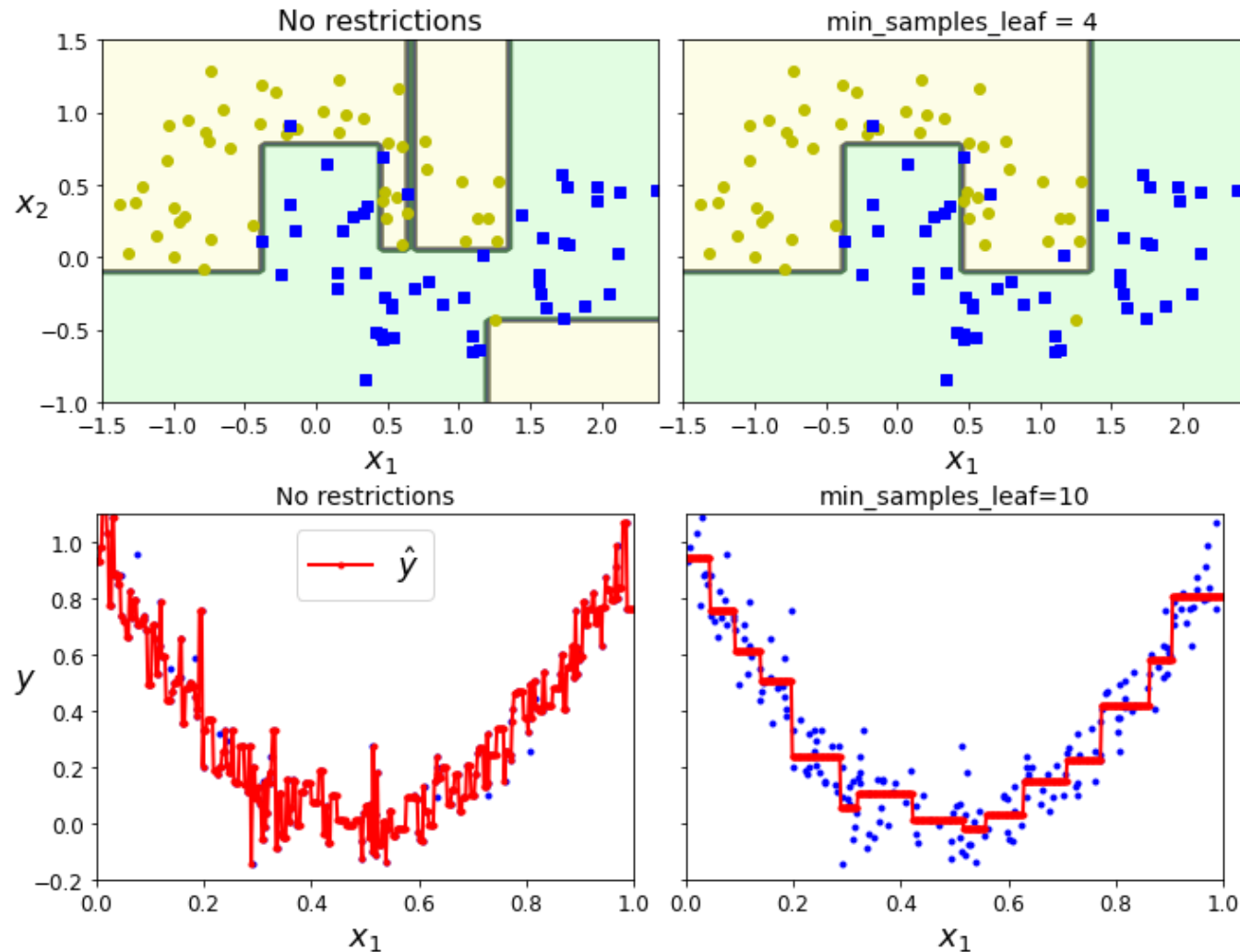
3. Энтропийный критерий

$$H(R) = - \sum_{k=1}^K p_k \log p_k \rightarrow \min$$

Решающее дерево. Глубина



Решающее дерево. Количество наблюдений в листьях



Решающее дерево. Параметры

criterion – критерий разбиения. Для классификации Джини, энтропия или логлосс.

splitter – стратегия разбиение. Best или random.

max_depth – максимальная глубина дерева. От 1 до бесконечности. По умолчанию *None*.

min_samples_split – минимальное количество наблюдений в листе, который будут делить дальше. От 1 до бесконечности. По умолчанию 2.

min_samples_leaf – минимальное количество наблюдений в дочернем листе. От 1 до бесконечности. По умолчанию 1.

Решающие деревья. Обрезка

Tree Pruning Example

