

# ДРЕВОВИДНЫЕ МОДЕЛИ

---

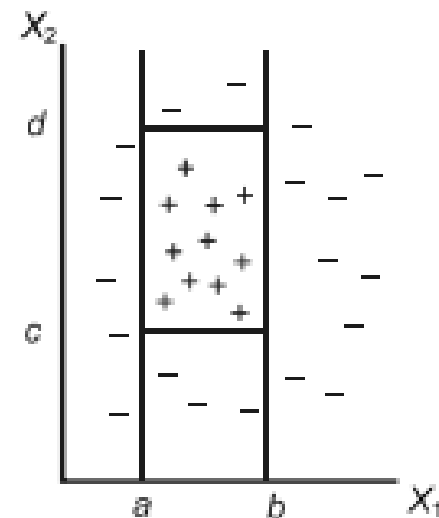
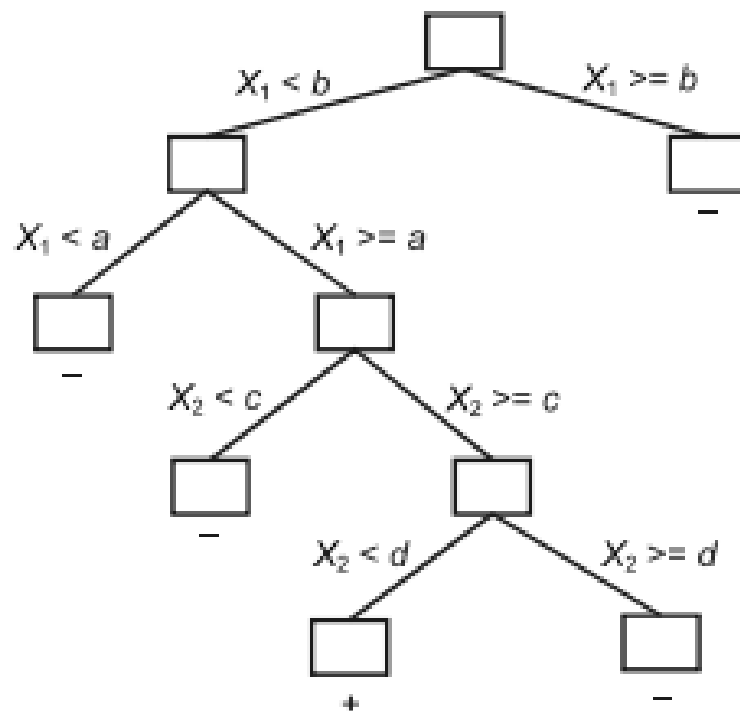
Деревья решений,  
случайный лес

# Дерево решений

*Деревья решений* - это метод, позволяющий предсказывать значения зависимой переменной в зависимости от соответствующих значений одной или нескольких предикторных (независимых) переменных. Применяется в задачах классификации и (реже) регрессии.

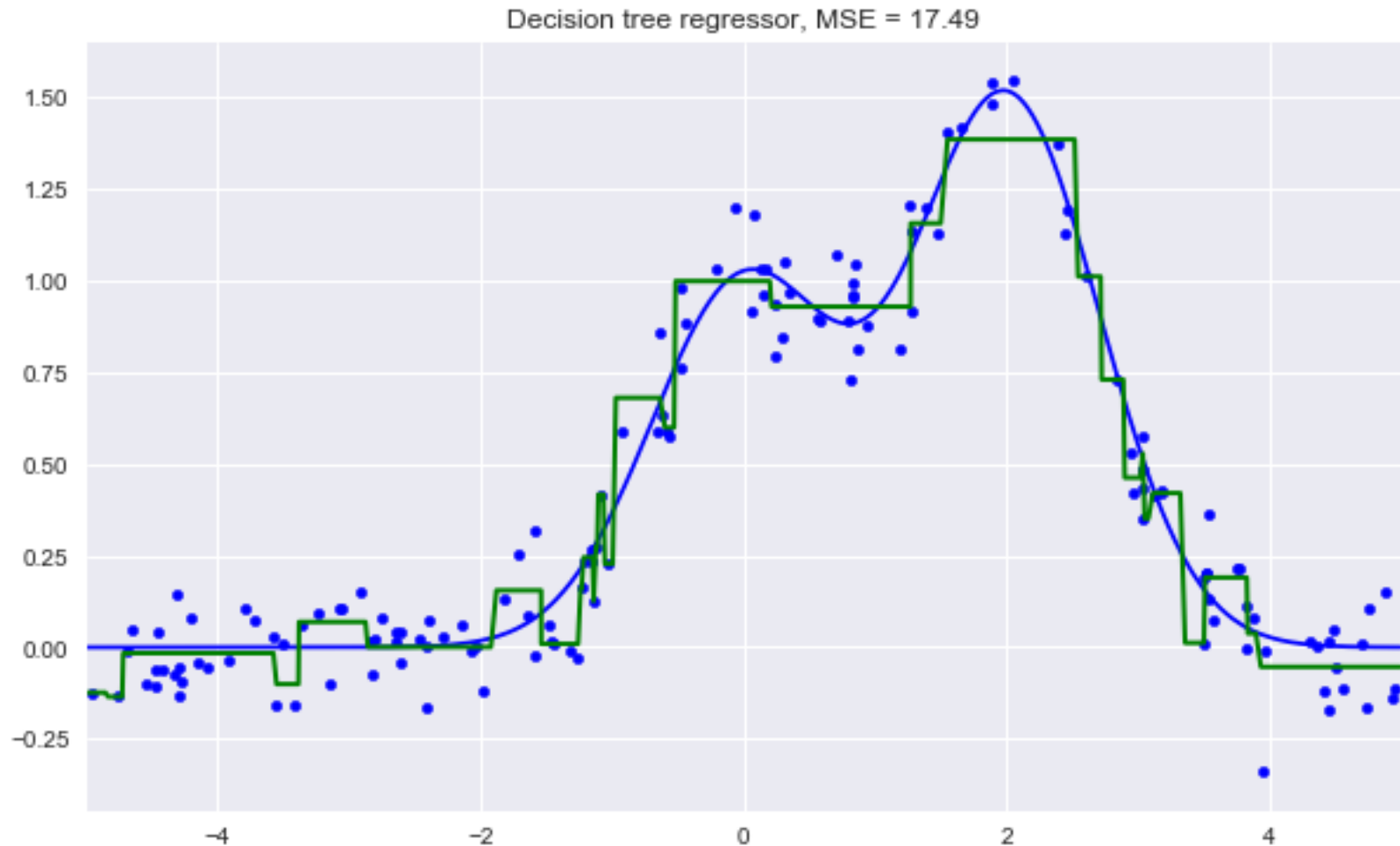


# Графическая иллюстрация нелинейного разделения классов



На рисунки приведен пример классификации объектов по двум непрерывным признакам. Объекты, относящиеся к разным классам, отмечены знаками "+" и "-".

# Использование деревьев решений в задачах регрессии



# Этапы построения дерева решений

- 1. Выбор критерия точности прогноза
- 2. Выбор типа ветвления
- 3. Определение момента прекращения ветвлений
- 4. Определение "подходящих" размеров дерева

## Выбор критерия точности прогноза

Accuracy, precision, recall – в задачах классификации

MSE, MAE – в задачах регрессии

# Выбор типа ветвления (criterion)

- Есть различные способы выбирать очередной признак для текущего ветвления:
- Алгоритм ID3, где выбор атрибута происходит на основании прироста информации ( [\*Gain\*](#) ).
- Алгоритм C4.5 (улучшенная версия ID3), где выбор атрибута происходит на основании нормализованного прироста информации ( [\*Gain Ratio\*](#) ).
- Алгоритм CART где выбор атрибута происходит на основании индекса Джини.

# Энтропия

Энтропия Шеннона для системы с  $N$  возможными состояниями:

- $H = - \sum_{i=1}^s p_i \log_2 p_i$

$p_i$  — вероятности нахождения системы в  $i$  — м состоянии

В нашем случае:

Предположим, что имеется множество  $A$ , состоящее из  $n$  элементов, обладающих свойством  $S$ , которое может принимать  $s$  различных значений,  $m_i$  — количество объектов множества  $A$ , имеющих  $i$ -е значение свойства  $S$ . Тогда

$$p_i = \frac{m_i}{n},$$

$$H(A, S) = - \sum_{i=1}^s \frac{m_i}{n} \log \frac{m_i}{n}.$$

# Прирост информации (ID3)

Предположим, что множество  $A$  элементов, характеризующихся свойством  $S$ , классифицировано посредством атрибута  $Q$ , имеющего  $q$  возможных значений. Тогда прирост информации (*information gain*) определяется как

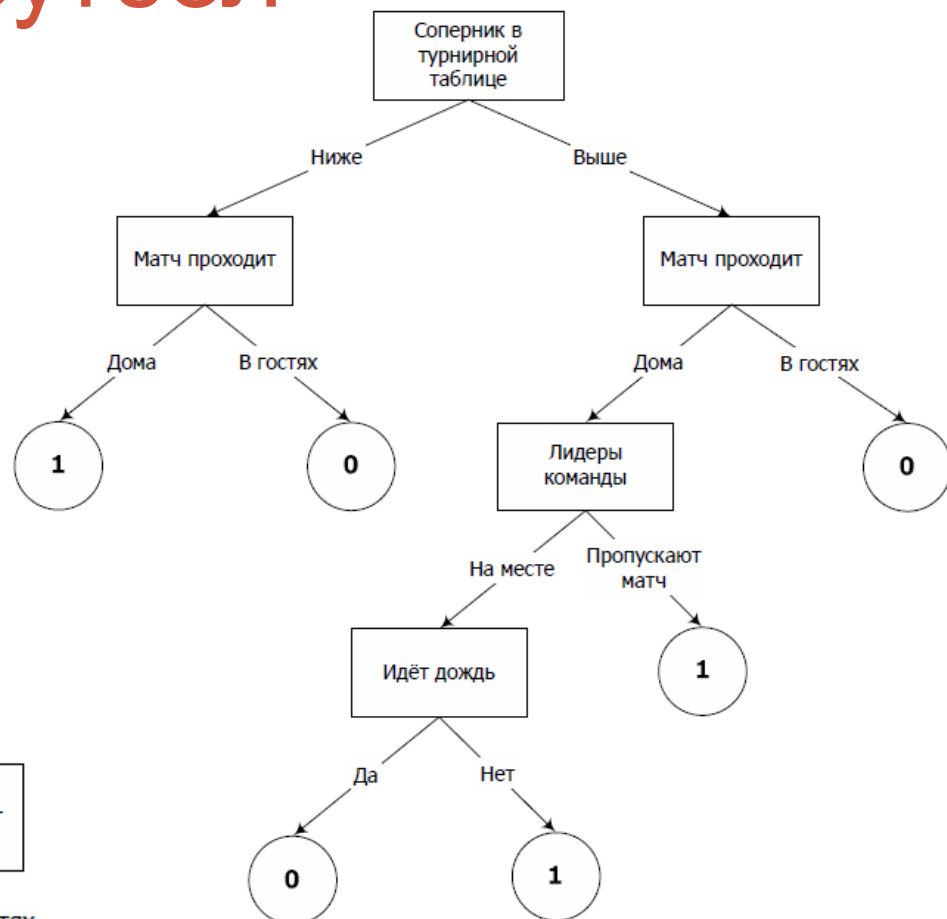
$$\text{Gain}(A, Q) = H(A, S) - \sum_{i=1}^q \frac{|A_i|}{|A|} H(A_i, S),$$

где  $A_i$  — множество элементов  $A$ , на которых атрибут  $Q$  имеет значение  $i$ .

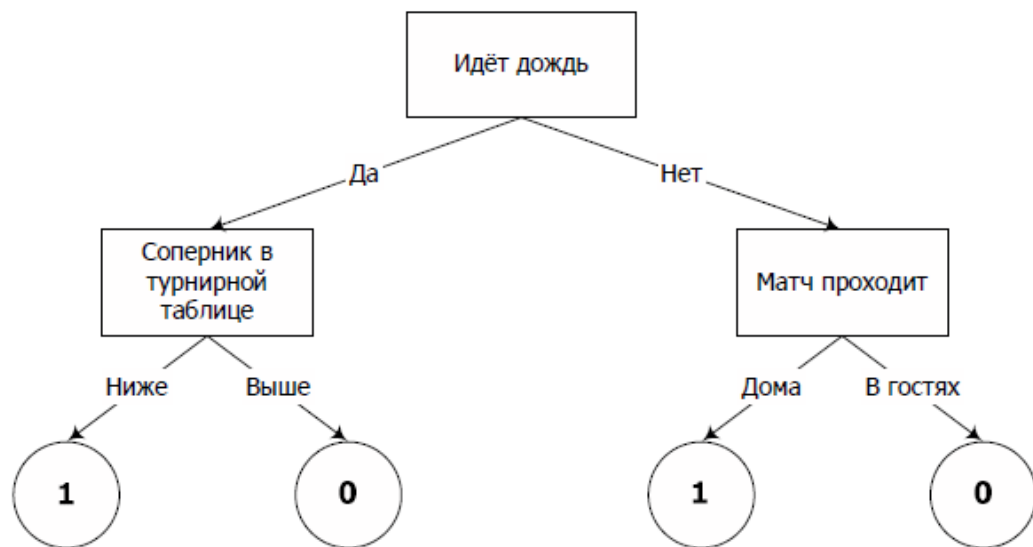


# Прогноз игры в футбол

Соперник	Играем	Лидеры	Дождь	Победа
Выше	Дома	На месте	Да	Нет
Выше	Дома	На месте	Нет	Да
Выше	Дома	Пропускают	Нет	Да
Ниже	Дома	Пропускают	Нет	Да
Ниже	В гостях	Пропускают	Нет	Нет
Ниже	Дома	Пропускают	Да	Да
Выше	В гостях	На месте	Да	Нет
Ниже	В гостях	На месте	Нет	???



Первый вариант дерева



Второй вариант дерева

# Вычисление энтропии и прироста информации

$$H(A, \text{Победа}) = -\frac{4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7} \approx 0.9852.$$

$$\begin{aligned} \text{Gain}(A, \text{Соперник}) &= H(A, \text{Победа}) - \frac{4}{7} H(A_{\text{выигрывает}}, \text{Победа}) - \frac{3}{7} H(A_{\text{проигрывает}}, \text{Победа}) \approx \\ &\approx 0.9852 - \frac{4}{7} \left( -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) - \frac{3}{7} \left( -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) \approx 0.0202. \end{aligned}$$

$$\text{Gain}(A, \text{Играем}) = H(A, \text{Победа}) - \frac{5}{7} H(A_{\text{дома}}, \text{Победа}) - \frac{2}{7} H(A_{\text{в гостях}}, \text{Победа}) \approx 0.4696.$$

$$\text{Gain}(A, \text{Лидеры}) = H(A, \text{Победа}) - \frac{3}{7} H(A_{\text{на месте}}, \text{Победа}) - \frac{4}{7} H(A_{\text{пропускают}}, \text{Победа}) \approx 0.1281.$$

$$\text{Gain}(A, \text{Дождь}) = H(A, \text{Победа}) - \frac{3}{7} H(A_{\text{да}}, \text{Победа}) - \frac{4}{7} H(A_{\text{нет}}, \text{Победа}) \approx 0.1281.$$

# Нормализованный прирост информации (C4.5)

Проблема: прирост информации выбирает атрибуты, у которых

Gain Ratio учитывает не только количество информации, требуемое для записи результата, но и количество информации, требуемое для разделения по текущему атрибуту.

Поправка:

$$\text{SplitInfo}(A, Q) = - \sum_{i=1}^q \frac{|A_q|}{|A|} \log_2 \frac{|A_q|}{|A|},$$

Сам критерий — максимизация величины

$$\text{GainRatio}(A, Q) = \frac{\text{Gain}(A, Q)}{\text{SplitInfo}(A, Q)}.$$

# Индекс Gini (CART)

Для набора тестов  $A$  и свойства  $S$ , имеющего  $s$  значений, этот индекс вычисляется как

$$\text{Gini}(A, S) = 1 - \sum_{i=1}^s \left( \frac{|A_i|}{|A|} \right)^2.$$

Соответственно, для набора тестов  $A$ , атрибута  $Q$ , имеющего  $q$  значений, и целевого свойства  $S$ , имеющего  $s$  значений, индекс вычисляется следующим образом:

$$\text{Gini}(A, Q, S) = \text{Gini}(A, S) - \sum_{j=1}^q \frac{|A_j|}{|A|} \text{Gini}(A_j, S).$$

# Правила разбиения (CART)

- 1) Вектор, подаваемый на вход дерева может содержать как порядковые так и категориальные переменные.
- 2) В каждом узле разбиение идет только по *одной переменной*.

2.1) Если переменная числового типа, то в узле формируется правило вида  $x_i \leq c$ . Где  $c$  – некоторый порог, который чаще всего выбирается как среднее арифметическое двух соседних *упорядоченных* значений переменной  $x_i$  обучающей выборки.

2.2) Если переменная категориального типа, то в узле формируется правило  $x_i \in V(x_i)$ , где  $V(x_i)$  – некоторое непустое подмножество множества значений переменной  $x_i$  в обучающей выборке.

Следовательно, для  $n$  значений числового атрибута алгоритм сравнивает  $n-1$  разбиений, а для категориального  $(2^{n-1} - 1)$ .

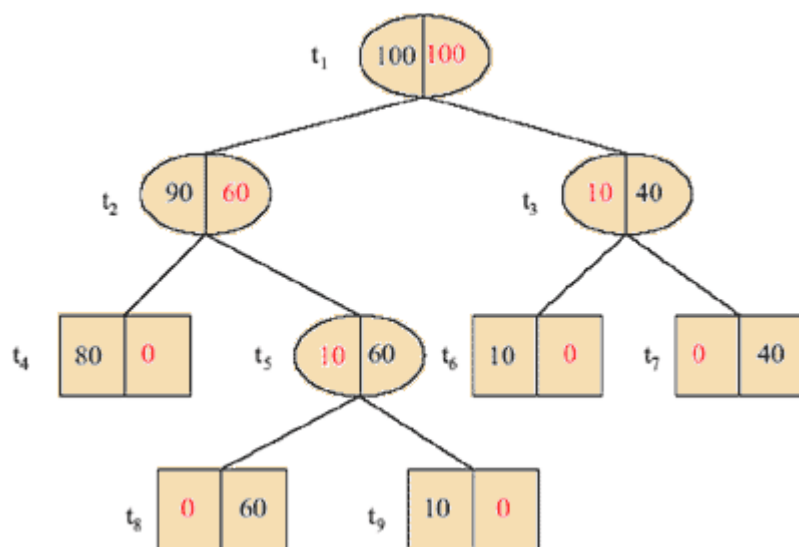
# Правила остановки

- **Минимальное число объектов, при котором выполняется расщепление (`min_samples_split`).** В этом варианте ветвление прекращается, когда все терминальные вершины, содержащие более одного класса, содержат не более чем заданное число объектов (наблюдений).
- **Минимальное число объектов в листьях (`min_samples_leaf`)**
- **Доля неклассифицированных.** В этом варианте ветвление прекращается, когда все терминальные вершины, содержащие более одного класса, содержат не более чем заданную долю неправильно классифицированных объектов (наблюдений).
- **Максимальная глубина деревьев (`max_depth`)**

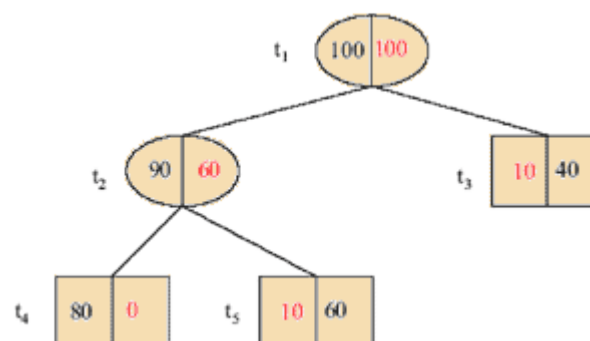
# Механизм отсечения дерева (CART)

Обозначим  $|T|$  – число листов дерева,  $R(T)$  – ошибка классификации дерева, равная отношению числа неправильно классифицированных примеров к числу примеров в обучающей выборке. Определим  $C_\alpha(T)$  – полную стоимость (оценку/показатель затраты-сложность) дерева  $T$  как:

$C_\alpha(T) = R(T) + \alpha * |T|$ , где  $|T|$  – число листов (терминальных узлов) дерева, – некоторый параметр, изменяющийся от 0 до  $+\infty$ . Полная стоимость дерева состоит из двух компонент – ошибки классификации дерева и штрафа за его сложность.

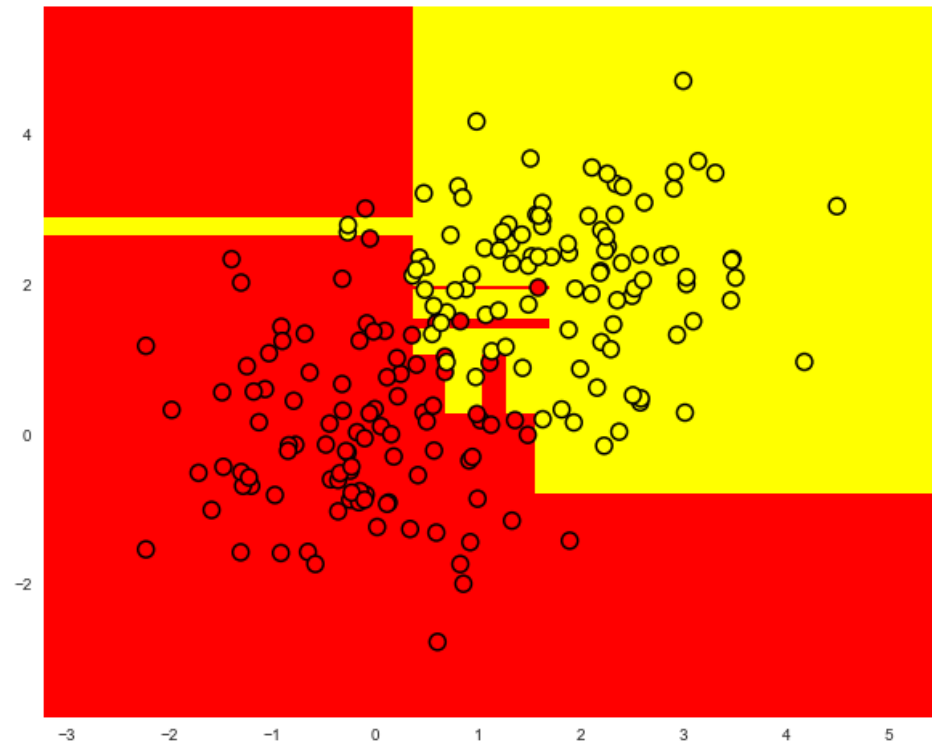
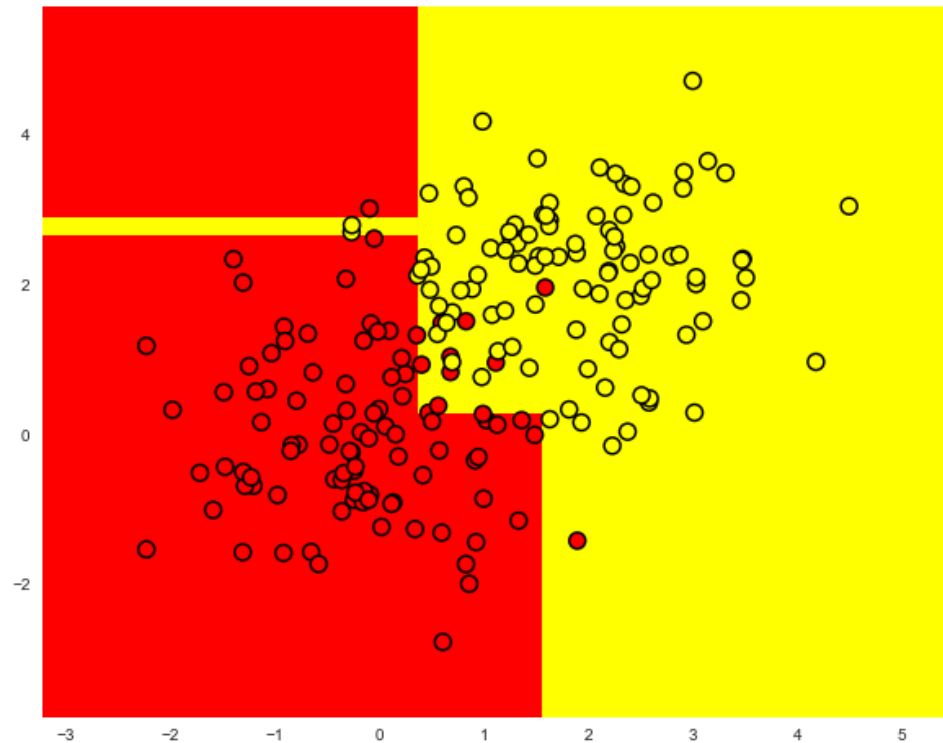


$\alpha_1 = 0,$



$\alpha_2 = 1/20,$

# Иллюстрация переобучения



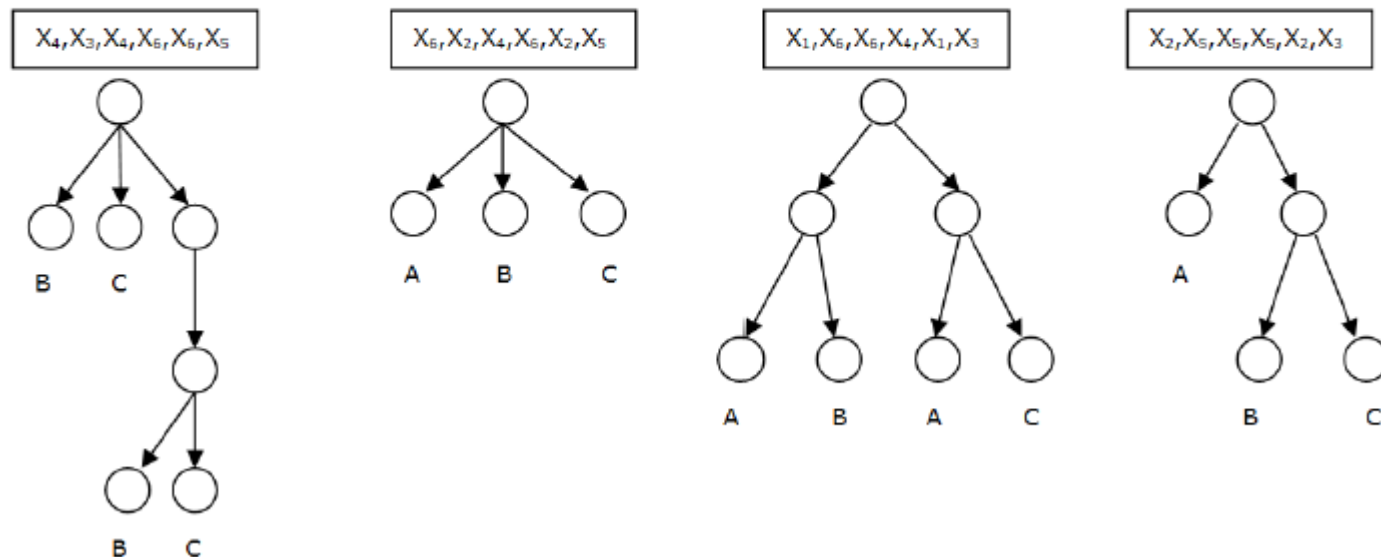


# Случайный лес (Random forest)

- Случайный лес — алгоритм машинного обучения, заключающийся в использовании комитета (ансамбля) деревьев решений.

Тренировочный набор:

$\{(X_1, A), (X_2, A), (X_3, B), (X_4, B), (X_5, C), (X_6, C)\}$



# Обучение случайного леса

- Пусть обучающая выборка состоит из  $N$  примеров, размерность пространства признаков равна  $M$ , и задан параметр  $m$  (в задачах классификации обычно  $m \approx \sqrt{M}$ ).
- Все деревья комитета строятся независимо друг от друга по следующей процедуре:
- Сгенерируем случайную подвыборку **с повторением** размером  $N$  из обучающей выборки. (Таким образом, некоторые примеры попадут в неё несколько раз, а в среднем  $N \left(1 - \frac{1}{N}\right)^N$ , т.е. примерно  $N/e$  примеров не войдут в неё вообще)
- Построим дерево, классифицирующее примеры данной подвыборки, причём в ходе создания очередного узла дерева будем выбирать признак, на основе которого производится разбиение, не из всех  $M$  признаков, а лишь из  $m$  случайно выбранных.
- Дерево строится до полного исчерпания подвыборки и не подвергается процедуре отсечения.
- Классификация объектов проводится путём голосования: каждое дерево комитета относит классифицируемый объект к одному из классов, и побеждает класс, за который проголосовало наибольшее число деревьев.
- Оптимальное число деревьев (**n\_estimators**) подбирается таким образом, чтобы минимизировать ошибку классификатора на валидационной выборке.

# Достоинства и недостатки

- Достоинства:

- Способность эффективно обрабатывать данные с большим числом признаков и классов.
- Нечувствительность к масштабированию значений признаков.
- Одинаково хорошо обрабатываются как непрерывные, так и дискретные признаки. Существуют методы построения деревьев по данным с пропущенными значениями признаков.
- Существуют методы оценивания значимости отдельных признаков в модели.
- Высокая параллелизуемость и масштабируемость.

- Недостатки:

Большой размер получающихся моделей.