

## Лекция 7

### Решающие деревья

#### 1. Введение в решающие деревья

Решающие деревья — это мощный инструмент машинного обучения, применяемый как для классификации, так и для регрессии. Структура дерева позволяет моделям принимать решения на основе последовательных вопросов, что делает их интерпретируемыми и легко объяснимыми для пользователей. Каждая ветвь дерева представляет возможные решения, а каждый лист — конечный вывод или класс. Решающие деревья особенно полезны в ситуациях, когда важно понимать процесс принятия решений, так как они наглядно показывают, как изменяются результаты в зависимости от различных признаков.

Данный метод активно используется в различных областях, таких как медицина, маркетинг, банковское дело и биоинформатика, где необходимо объяснимое и логически структурированное принятие решений.

#### 2. Основные компоненты решающих деревьев

Решающие деревья состоят из нескольких ключевых элементов:

- **Корневой узел (Root Node):** Это начальная точка дерева, которая содержит все данные и выполняет первое разбиение на основе одного из признаков.
- **Внутренние узлы (Internal Nodes):** Узлы, которые содержат вопросы или условия для разбиения данных. Каждый внутренний узел разделяет данные на подмножества, в зависимости от значений признаков.
- **Листовые узлы (Leaf Nodes):** Конечные узлы дерева, в которых содержится предсказанный класс (для классификации) или значение (для регрессии).
- **Ветви (Branches):** Связи между узлами, представляющие пути, по которым данные могут следовать от корневого узла до листового.

Процесс разбиения данных основан на выборе признака, который максимизирует качество разделения. Важнейшими показателями качества являются прирост информации и критерий Джини, о которых подробнее рассказывается далее.

#### 3. Алгоритмы построения решающих деревьев

Существует несколько алгоритмов для построения решающих деревьев, которые различаются подходами к выбору разбиений и критериям остановки. Наиболее известные из них — CART (Classification and Regression Trees) и ID3 (Iterative Dichotomiser 3).

### 3.1 Алгоритм CART

CART — один из наиболее популярных алгоритмов для построения решающих деревьев. Он применим как для задач классификации, так и для регрессии. В основе CART лежит идея двоичного разделения, где каждый узел делится на две ветви. Основные шаги алгоритма CART включают:

1. Определение признака и порогового значения, которые максимизируют качество разделения.
2. Разделение данных на два подмножества.
3. Рекурсивное повторение процесса для каждого подмножества до тех пор, пока не будет достигнут критерий остановки (например, минимальный размер подмножества).

Для классификации CART использует критерий Джини, а для регрессии — среднеквадратическую ошибку.

### 3.2 Алгоритм ID3

ID3 — это алгоритм, разработанный для задач классификации, и использует прирост информации в качестве критерия для выбора признаков. Основные шаги:

1. Вычисление прироста информации для каждого признака.
2. Выбор признака с наибольшим приростом информации для текущего узла.
3. Рекурсивное повторение процесса для каждого подмножества данных.

## 4. Критерии разбиения

Для эффективного разбиения данных на каждом этапе построения дерева применяются различные метрики, такие как прирост информации и критерий Джини.

### 4.1 Прирост информации

Прирост информации (Information Gain) измеряет уменьшение неопределенности после разбиения узла по определенному признаку. В основе прироста информации лежит энтропия, которая характеризует степень хаоса или неопределенности в данных. Энтропия определяется следующим образом:

$$H(S) = -\sum_{i=1}^n p_i \log_2(p_i) H(S) = -\sum_{i=1}^n p_i \log_2(p_i) H(S) = -\sum_{i=1}^n p_i \log_2(p_i)$$

где  $p_i$  — вероятность класса  $i$  в текущем узле. Прирост информации вычисляется как разность энтропии до и после разбиения:

$$IG(S, A) = H(S) - \sum_{v \in Values(A)} |S_v| / |S| H(S_v) IG(S, A) = H(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} H(S_v) IG(S, A) = H(S) - \sum_{v \in Values(A)} |S_v| H(S_v)$$

где  $S_v$  — подмножество  $S$ , соответствующее значению  $v$  признака  $A$ .

## 4.2 Критерий Джини

Критерий Джини (Gini Index) используется для измерения «чистоты» узлов и часто применяется в алгоритме CART. Индекс Джини для набора данных определяется как:

$$Gini(S) = 1 - \sum_{i=1}^n p_i^2 Gini(S) = 1 - \sum_{i=1}^n p_i^2 Gini(S) = 1 - \sum_{i=1}^n p_i^2$$

где  $p_i$  — вероятность того, что случайно выбранный элемент принадлежит классу  $i$ . Значение критерия Джини варьируется от 0 (полностью чистый узел) до 0.5 (максимально смешанный узел). В процессе построения дерева алгоритм выбирает разбиения, минимизирующие критерий Джини.

## 5. Преимущества и недостатки решающих деревьев

Решающие деревья имеют ряд преимуществ, которые делают их популярным выбором для задач классификации и регрессии:

- Интерпретируемость:** Логика принятия решений в дереве легко понимается и объясняется.
- Гибкость:** Решающие деревья не требуют предварительного нормирования данных или кодирования категориальных признаков.
- Простота использования:** Алгоритм прост в реализации и требует минимальной настройки гиперпараметров.

Однако решающие деревья также обладают недостатками:

- Переобучение:** Деревья имеют тенденцию к переобучению на тренировочных данных, особенно если не применяются критерии остановки.
- Чувствительность к выбросам:** Выбросы могут существенно повлиять на структуру дерева.

- **Неустойчивость:** Небольшие изменения в данных могут привести к значительным изменениям в структуре дерева.

## 6. Методы предотвращения переобучения

Для повышения устойчивости и предотвращения переобучения в решающих деревьях применяются методы подрезки (pruning) и ограничения глубины дерева.

### 6.1 Подрезка

Подрезка (Pruning) — это метод удаления неинформативных узлов в дереве для уменьшения его сложности и предотвращения переобучения. Существует два основных типа подрезки:

- **До подрезки (Pre-Pruning):** Останавливает рост дерева, если достигнуты определенные критерии, такие как максимальная глубина или минимальный размер узла.
- **После подрезки (Post-Pruning):** Сначала строится полное дерево, а затем удаляются узлы, которые незначительно влияют на точность модели.

### 6.2 Ограничение глубины дерева

Ограничение максимальной глубины дерева — это еще один способ контроля переобучения. Чем глубже дерево, тем больше вероятность переобучения, поскольку модель может слишком точно подстраиваться под тренировочные данные. Глубина дерева является важным гиперпараметром, который следует оптимизировать с помощью перекрестной проверки.

## 7. Ансамблирование решающих деревьев

Для повышения точности и стабильности модели часто используются методы ансамблирования, такие как бэггинг и бустинг.

### 7.1 Случайные леса

Случайный лес (Random Forest) — это ансамблевый метод, который создает несколько решающих деревьев и усредняет их предсказания. Каждый из деревьев обучается на случайном подмножестве признаков и данных, что повышает разнообразие моделей и снижает вероятность переобучения. Случайные леса отличаются высокой точностью и устойчивостью к выбросам, а также снижают вероятность переобучения.

### 7.2 Градиентный бустинг

Градиентный бустинг (Gradient Boosting) создает последовательность деревьев, где каждое последующее дерево корректирует ошибки предыдущего. Модель работает итеративно, уменьшая ошибку на каждом шаге. Градиентный бустинг показал высокую точность в задачах классификации и регрессии, но требует больше вычислительных ресурсов, чем случайные леса, и более подвержен переобучению.

## 8. Применение решающих деревьев

Решающие деревья находят применение в самых различных областях:

- **Медицина:** Диагностика заболеваний на основе симптомов и клинических данных.
- **Финансы:** Оценка кредитоспособности и предсказание дефолтов.
- **Маркетинг:** Сегментация клиентов на основе их поведения и предпочтений.
- **Биоинформатика:** Классификация белков и генов по их функциональным признакам.

## 9. Заключение

Решающие деревья представляют собой эффективный и понятный метод машинного обучения, способный решать задачи как классификации, так и регрессии. Они обладают рядом преимуществ, таких как интерпретируемость и гибкость, однако их недостатки, такие как склонность к переобучению, требуют применения методов регуляризации и ансамблирования.