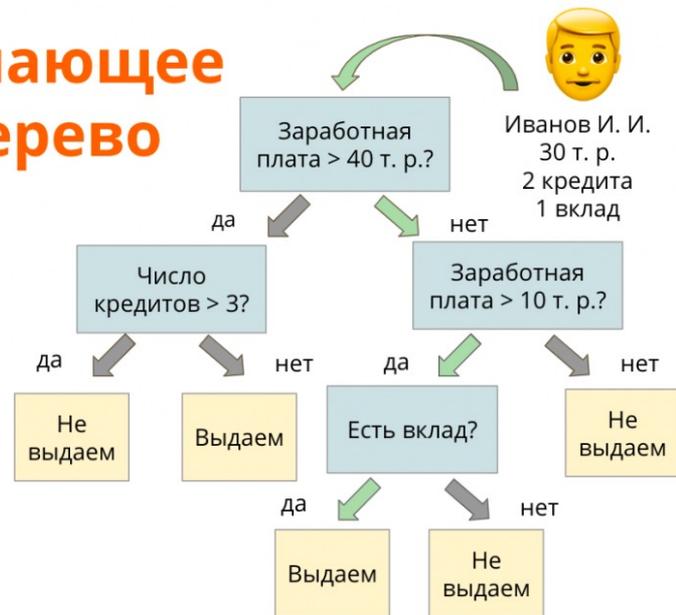


Лекция 3. ДЕРЕВЬЯ РЕШЕНИЙ

Решающее дерево — это алгоритм, который делает предсказания на основе серии вопросов об объекте.



Решающее дерево



К какой форме отнести публикацию?



Задачи

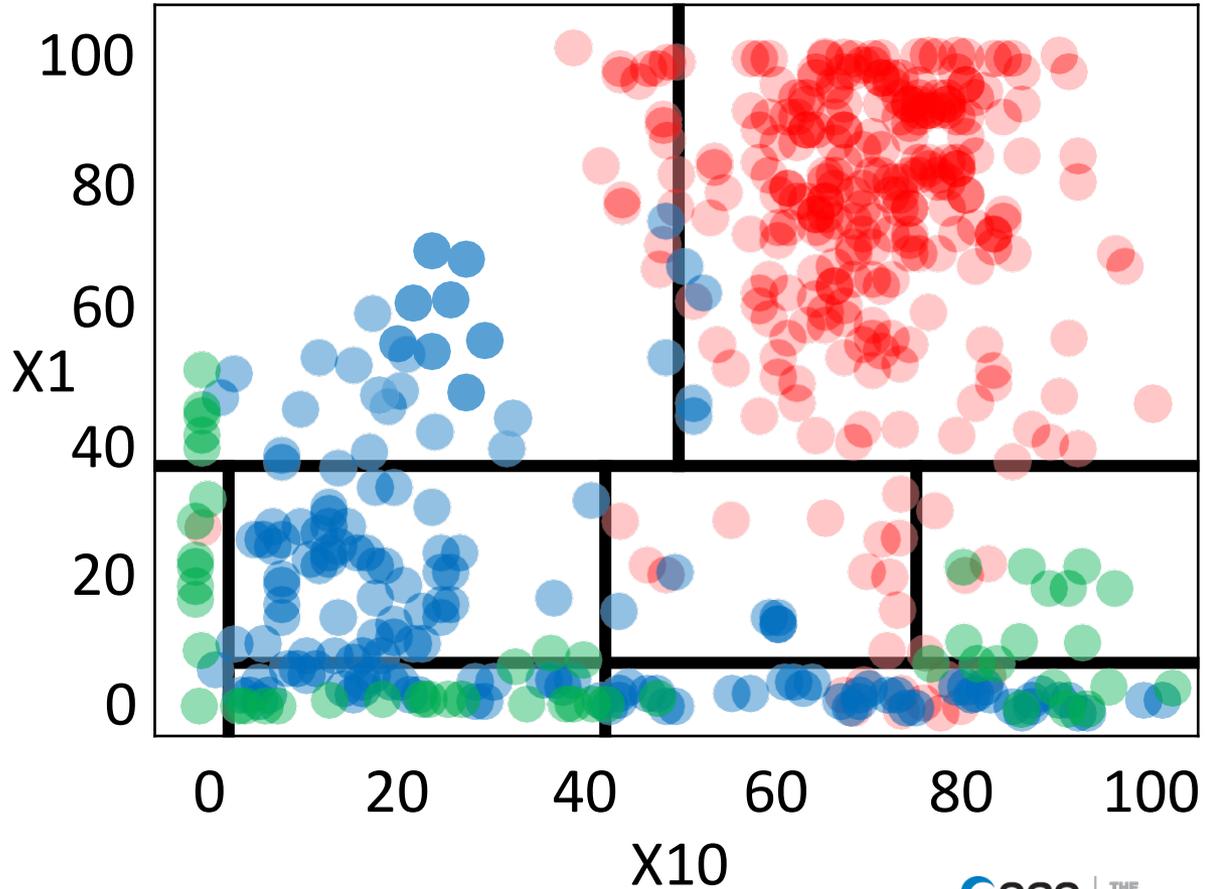
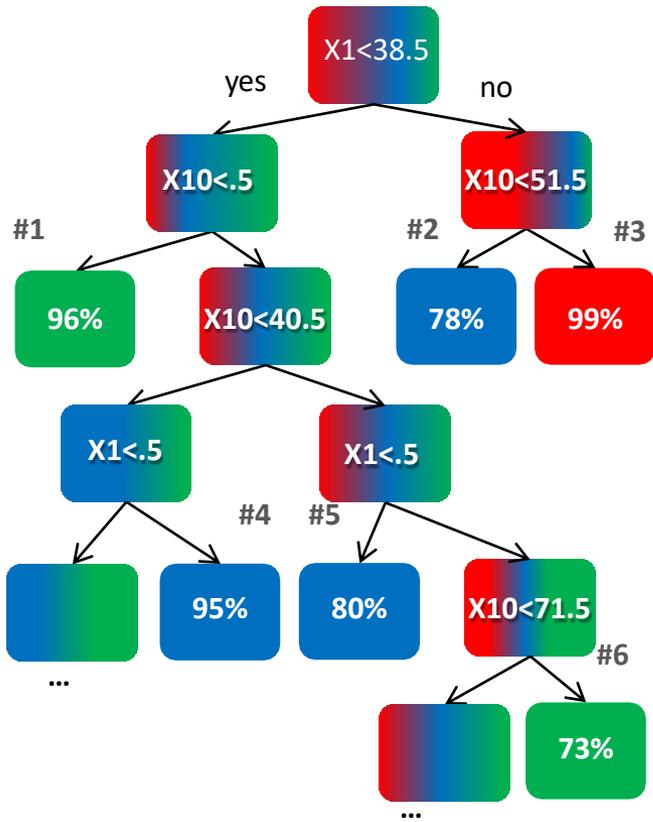
Основная сфера применения деревьев решений — поддержка процессов принятия управленческих решений, используемая в статистике, анализе данных и **машинном обучении**. Задачами, решаемыми с помощью данного аппарата, являются:

- **Классификация** — отнесение объектов к одному из заранее известных классов. Целевая переменная должна иметь дискретные значения.
- **Регрессия (численное предсказание)** — предсказание числового значения независимой переменной для заданного входного вектора.
- **Описание объектов** — набор правил в дереве решений позволяет компактно описывать объекты. Поэтому вместо сложных структур, описывающих объекты, можно хранить деревья решений.

Основные этапы построения

1. Выбор атрибута, по которому будет производиться разбиение в данном узле (атрибута разбиения).
2. Выбор критерия остановки обучения.
3. Выбор метода отсечения ветвей (упрощения).
4. Оценка точности построенного дерева.

НОМИНАЛЬНАЯ ЦЕЛЕВАЯ ПЕРЕМЕННАЯ



НОМИНАЛЬНАЯ ЦЕЛЕВАЯ ПЕРЕМЕННАЯ: РЕШЕНИЕ

#3

99%

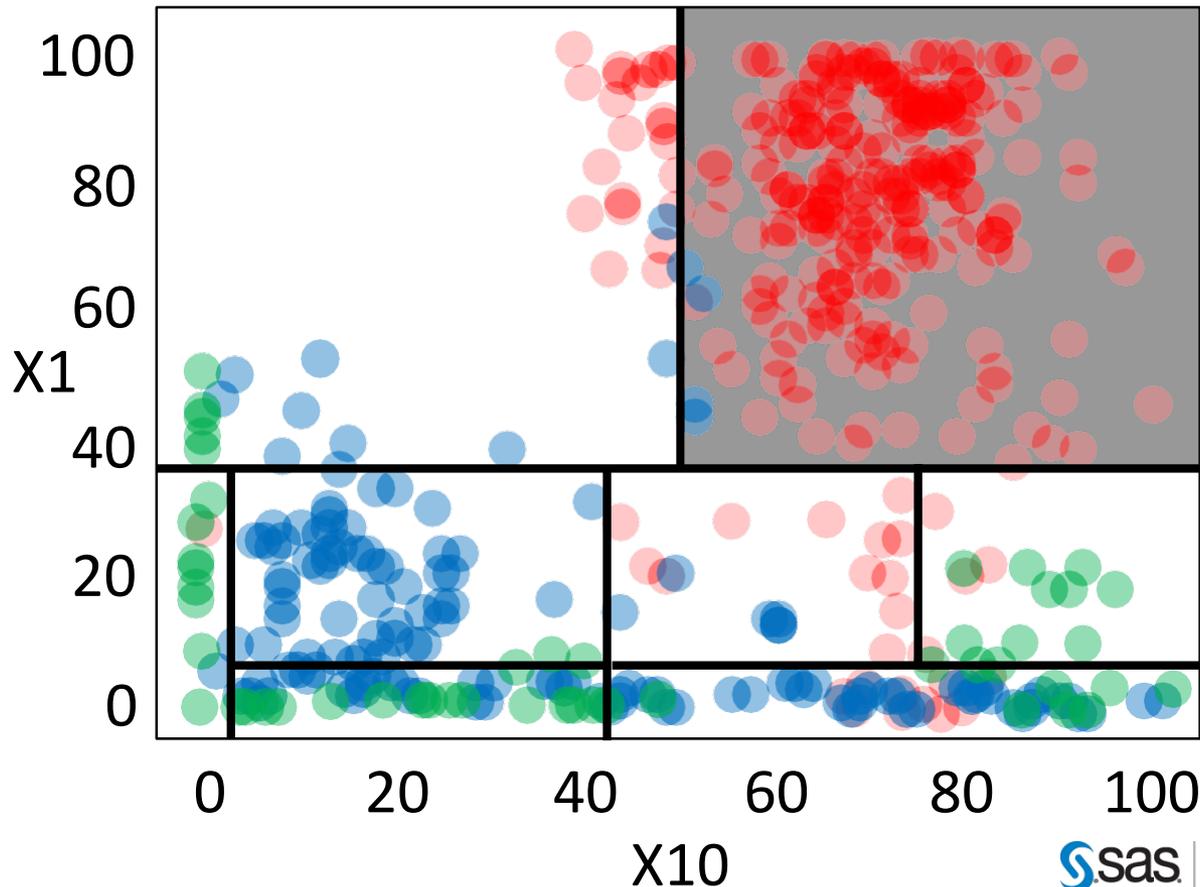
$P(\text{red}) = 0.99$

$P(\text{blue}) = 0.01$

$P(\text{green}) = 0.00$

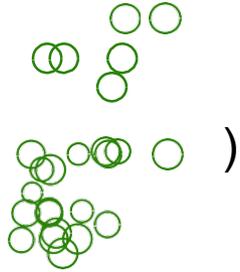


Prediction = 

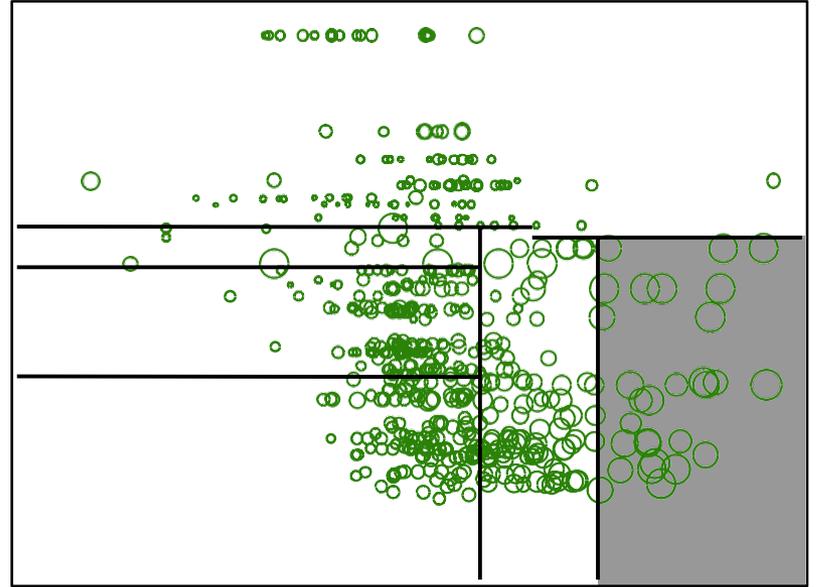


ИНТЕРВАЛЬНАЯ ЦЕЛЕВАЯ ПЕРЕМЕННАЯ: РЕШЕНИЕ

Prediction = Average()

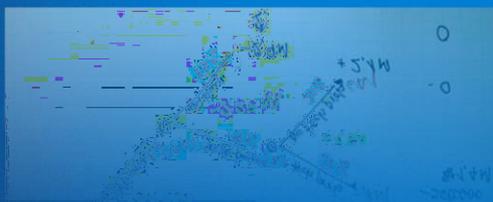


X1



X10

ПОИСК РАЗБИЕНИЙ

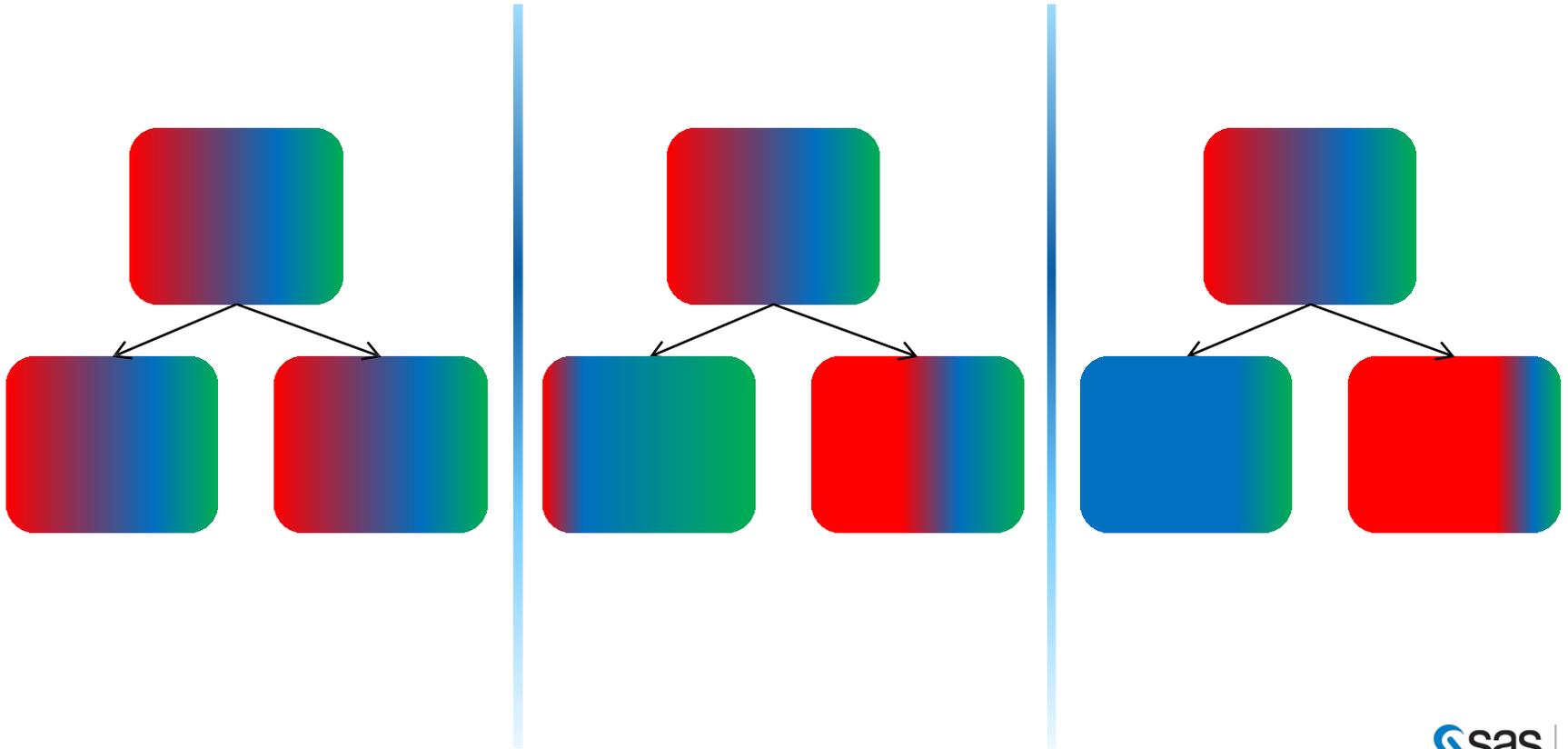


- Как выбрать «наилучшее» разбиение? – **КРИТЕРИЙ РАЗБИЕНИЯ**
- По какой переменной и какому значению производить разбиение? – **ПОИСК РАЗБИЕНИЯ**
- Когда нужно остановиться? – **ВЫБОР РАЗМЕРА ДЕРЕВА**

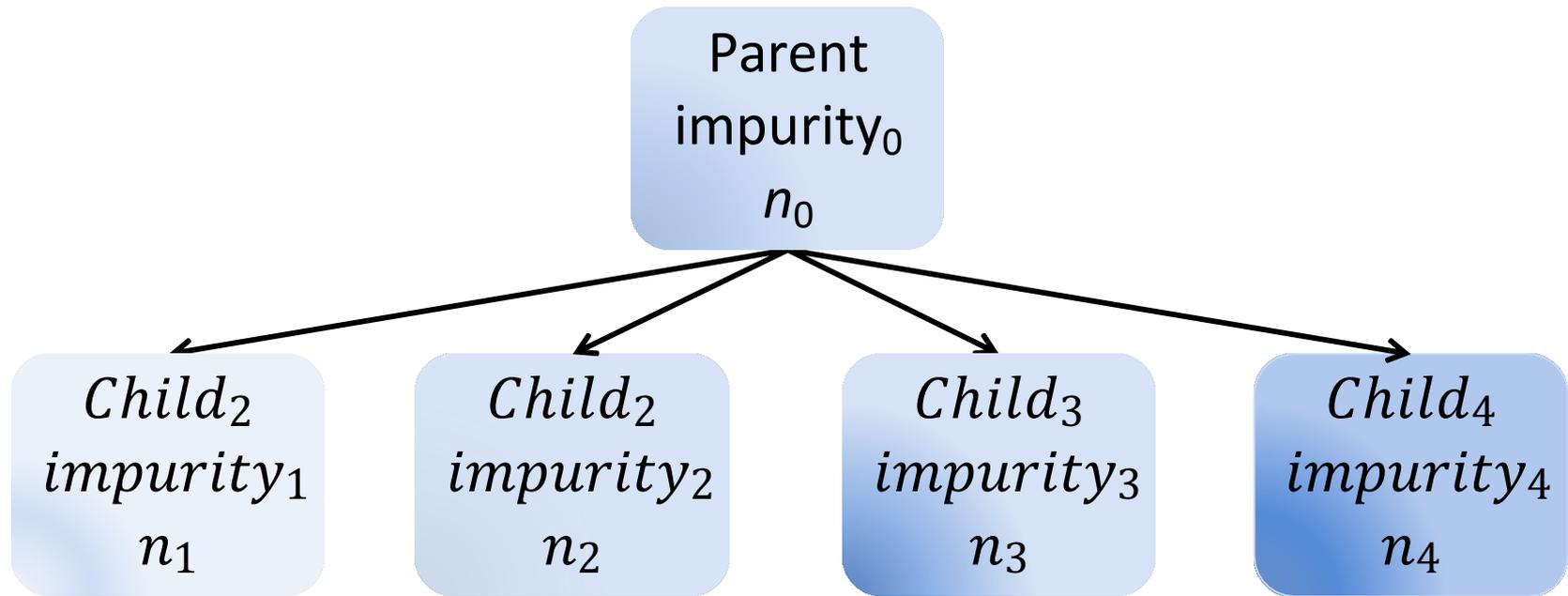


Критерии разбиения

КРИТЕРИИ РАЗБИЕНИЯ: ЧТО ТАКОЕ ХОРОШЕЕ РАЗБИЕНИЕ?



КРИТЕРИИ РАЗБИЕНИЯ: ЧТО ТАКОЕ ХОРОШЕЕ РАЗБИЕНИЕ?

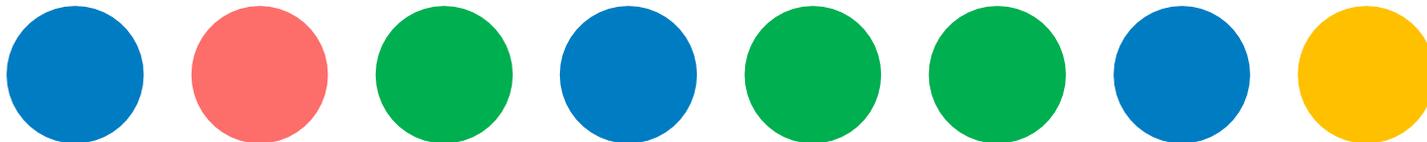


worth of split $\rightarrow \Delta i = i(0) - \left(\frac{n_1}{n_0} i(1) + \frac{n_2}{n_0} i(2) + \frac{n_3}{n_0} i(3) + \frac{n_4}{n_0} i(4) \right)$

КРИТЕРИИ РАЗБИЕНИЯ: GINI

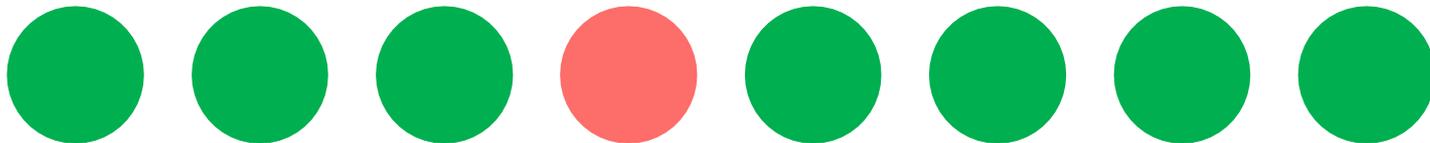
$$1 - \sum_{j=1}^r p_j^2 = 2 \sum_{j < k} p_j p_k$$

Большое
разнообразие,
низкая степень
чистоты



$$GINI = 1 - 2 \left(\frac{3}{8} \right)^2 - 2 \left(\frac{1}{8} \right)^2 = .69$$

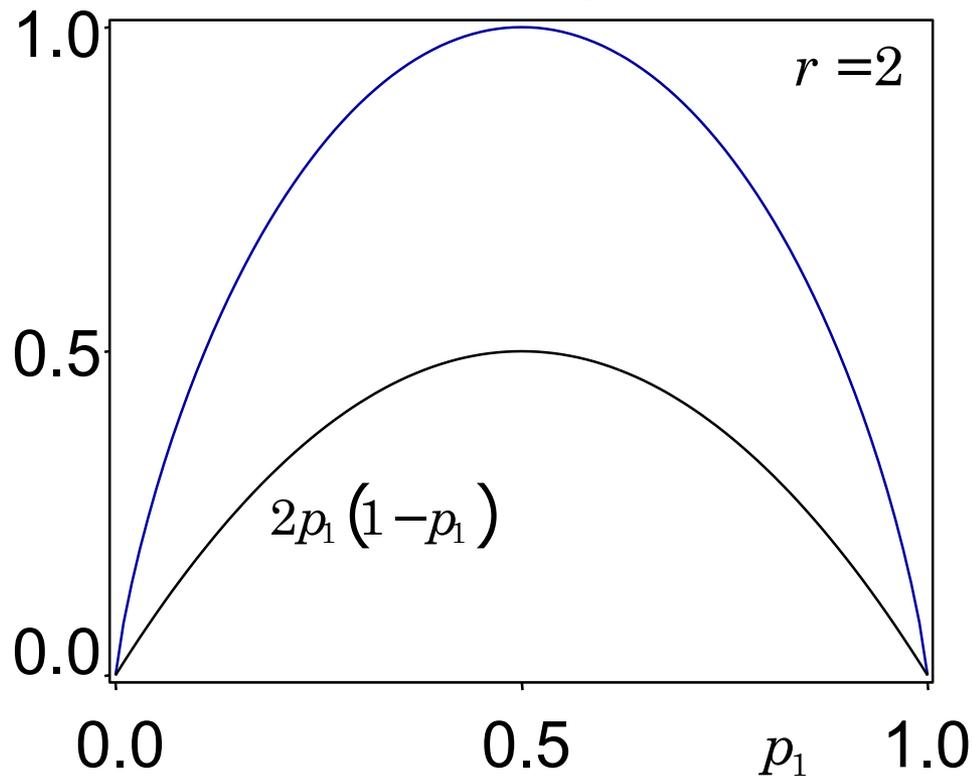
Небольшое
разнообразие,
высокая
степень чистоты



$$GINI = 1 - \left(\frac{7}{8} \right)^2 - \left(\frac{1}{8} \right)^2 = .22$$

КРИТЕРИИ РАЗБИЕНИЯ: ENTROPY

$$H(p_1, p_2, \dots, p_r) = -\sum_{i=1}^r p_i \log_2(p_i)$$



КРИТЕРИИ РАЗБИЕНИЯ: CHI-SQUARED

Observed
X1: <38.5 ≥38.5

	293	71	.342
	363	1	.342
	42	294	.316
	.656	.344	n=1064

Expected

239	125
239	125
225	116

$$\frac{(O - E)^2}{E}$$

12	23
64	123
149	273

$$X_v^2 = \sum \frac{(O - E)^2}{E} = 644 \rightarrow p\text{-value, worth} = -\log_{10}(p)$$

$$v = (3 - 1)(2 - 1) = 2$$

КРИТЕРИИ РАЗБИЕНИЯ: CHI-SQUARED

Поправки p-value

- Количество ветвей
- Количество различных значений переменной
- Количество переменных (кандидатов) для разбиения в данном узле
- Глубина текущего узла в дереве

Чем больше вариантов мы перебираем, тем больше шансы случайно выбрать неверное разбиение

Чем ниже (глубже) мы в дереве, тем меньше кол-во наблюдений в узле, p-value слишком оптимистичные

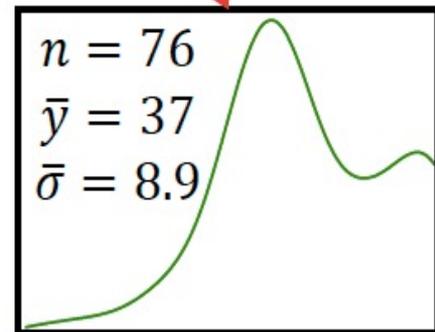
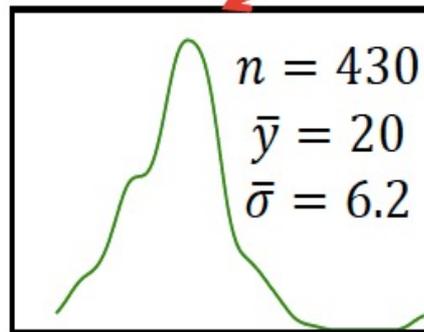
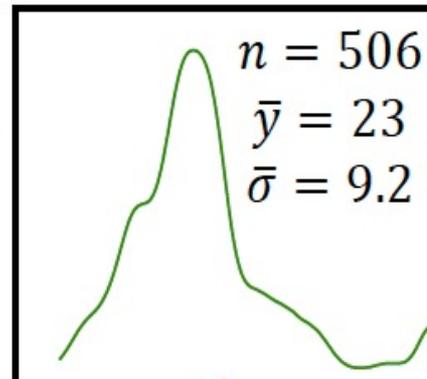
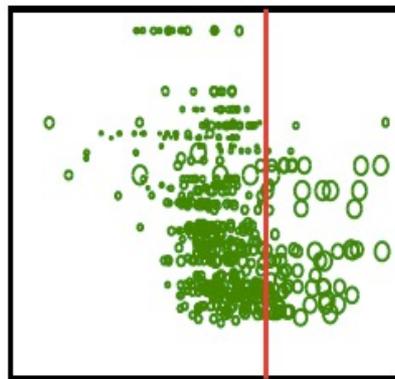
КРИТЕРИИ РАЗБИЕНИЯ: ИНТЕРВАЛЬНАЯ ЦЕЛЕВАЯ ПЕРЕМЕННАЯ

$$VARIANCE = \frac{\sigma_j^N}{1} = \frac{(Y_j - \mu)^2}{N}$$

$$F = \left(\frac{SS_{\text{between}}}{SS_{\text{within}}} \right) \left(\frac{n - B}{B - 1} \right)$$

p-value, worth = $-\log_{10}(p)$

поправки p-value



yes

$X_{10} < 6.94$

no



Поиск разбиения

ПОИСК РАЗБИЕНИЯ. ШАГ 1

2. Для переменной **X** удалить наблюдения:
 - с пропущенными значениями (кроме случая *Use In Search*)
 - если **X** – категориальная, то удалить наблюдения со значениями, встречающимися редко

Удаления всегда производятся независимо по каждой переменной и каждому узлу

ПОИСК РАЗБИЕНИЯ. ШАГ 2

1. Сортируем наблюдения:

- если **X** – интервальная или ординальная – просто сортировка по возрастанию, пропущенные значения отправляются в конец
- если **X** – номинальная, а целевая переменная **Y** – непрерывная, то сортируем уровни/значения $X: L_i$ по $avg(Y)$
- если **X** – номинальная, а целевая переменная **Y** – имеет два значения, то сортируем уровни/значения $X: L_i$ по доле одного из значений **Y**

В остальных случаях (**X** и **Y** – номинальная, **Y** имеет 3 и более значений) сортировка не производится

ПОИСК РАЗБИЕНИЯ. ШАГ 3

1. Рассматриваем все возможные

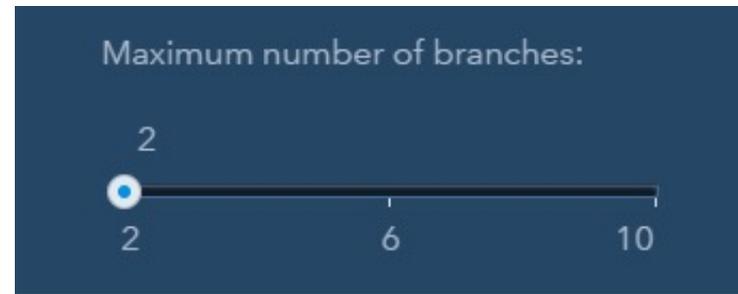
бинарные разбиения:

- если **X** – интервальная или ординальная – то разбиения между одинаковыми значениями X не допускаются
- если **X** – номинальная, а целевая переменная **Y** – непрерывная, то разбиение между разными значений X с одинаковым $avg(Y)$ не допускаются
- Разбиение, при котором в листе останется мало наблюдений не допускаются

Теперь найдены все возможные
БИНАРНЫЕ РАЗБИЕНИЯ

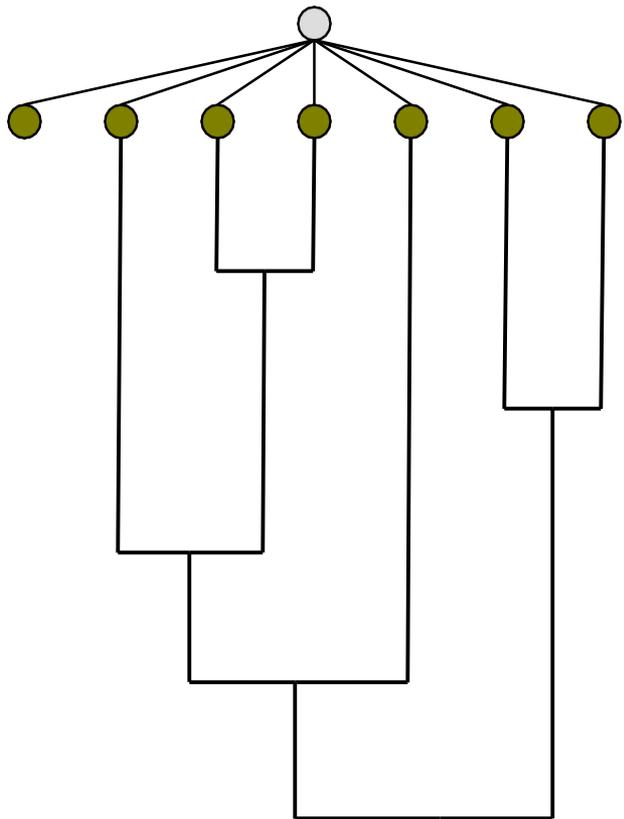
ПОИСК РАЗБИЕНИЯ. ШАГ 4

1. Рассматриваем все возможные **M**-арные разбиения ($2 \leq N \leq M$):
 - Перебираем все – выбираем одно лучшее
 - Или проводим консолидацию (*merge-and-shuffle algorithm*)



Perform clustering-based split search

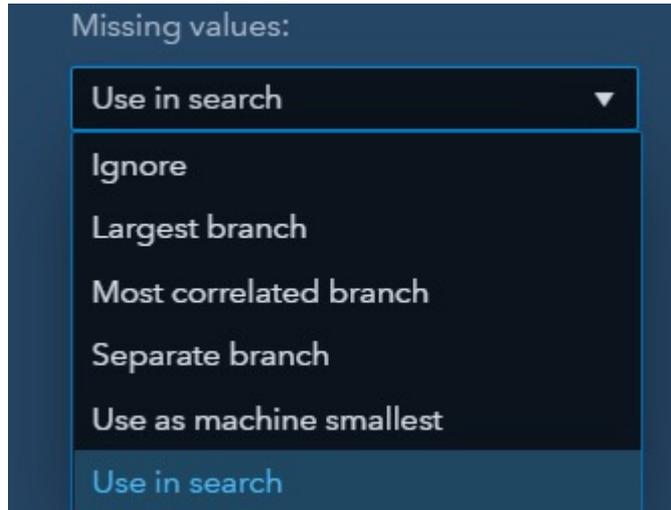
ПОИСК РАЗБИЕНИЯ. ШАГ 5



Merge-and-shuffle algorithm

1. Для каждого значения X создаем свою **узел**
 - Интервальные переменные разбиваем на N групп
2. Поочередно рассматриваем пары **узлов**:
 - Объединяем пару
 - Обрато разбиваем на две
 - Измеряем полезность такого разбиения
 - Объединяем пару с минимальной полезностью

РАБОТА С ПРОПУЩЕННЫМИ ЗНАЧЕНИЯМИ



Что можно сделать с
пропущенными значениями

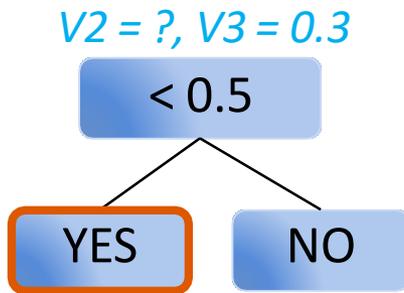
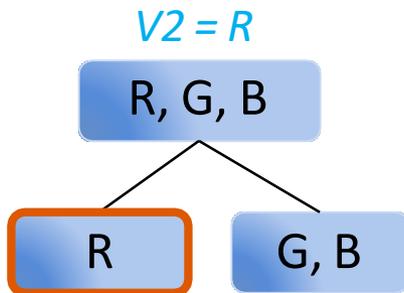
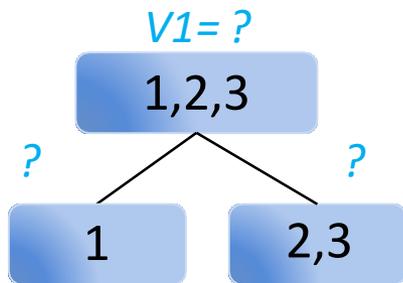
РАБОТА С ПРОПУЩЕННЫМИ ЗНАЧЕНИЯМИ



Суррогатные правила:

- в какую ветку отправится наблюдение, если значение *основной переменной* - пропущенное
- из всех возможных выбирается правило, максимально *соответствующее* основному

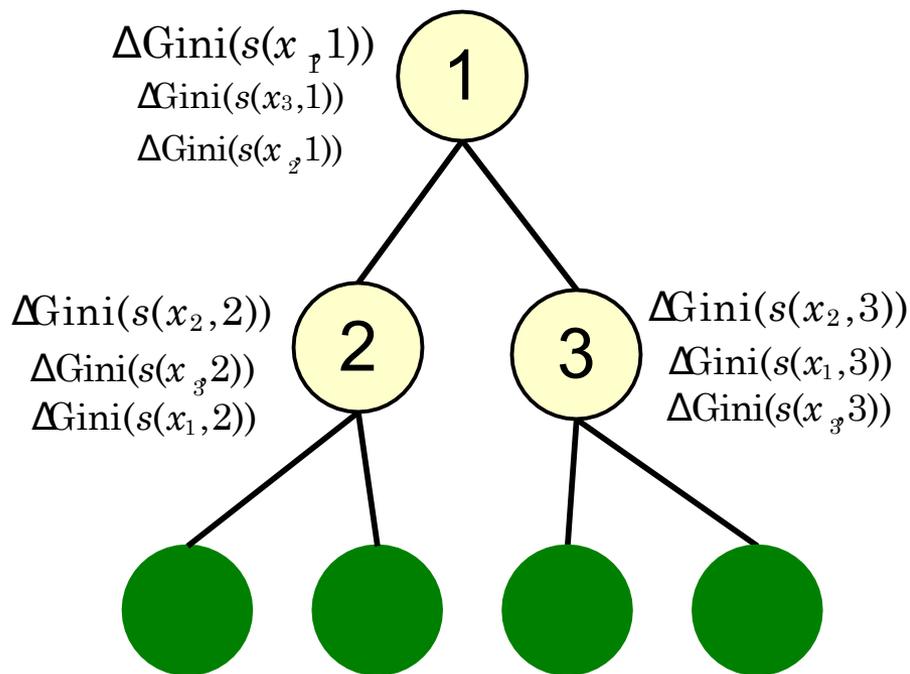
разбиение по **основной переменной V1**



...

Значимость переменных

ЗНАЧИМОСТЬ ПЕРЕМЕННЫХ: НА ОСНОВЕ РАЗБИЕНИЙ



Значимость (важность) переменной x_i :

$$Importance(x_j) = \sqrt{\sum_{t=1}^T \frac{n_t}{n} \times \alpha(s_x, t_j) \times \Delta i(s(x_j, t))}$$

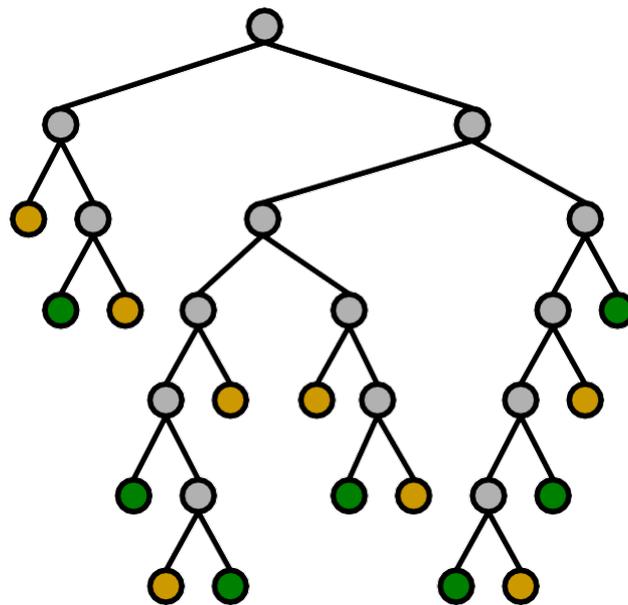
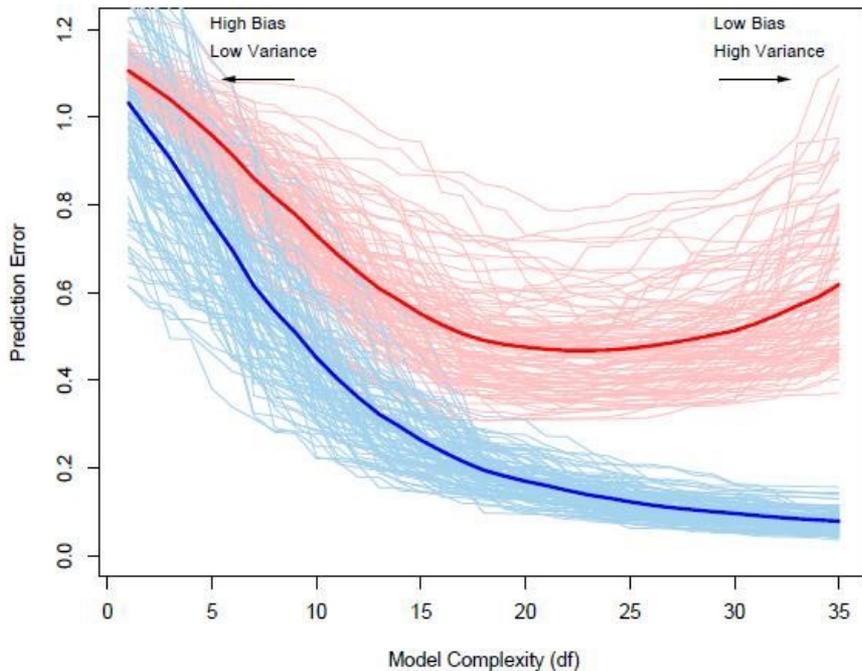
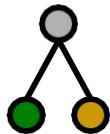
↓ вес ↓ соответствие ↓ улучшение

Нормируется: $Importance \in [0, 1]$



Остановка роста дерева & обрубание

КОГДА НУЖНО ПРЕКРАТИТЬ РОСТ ДЕРЕВА

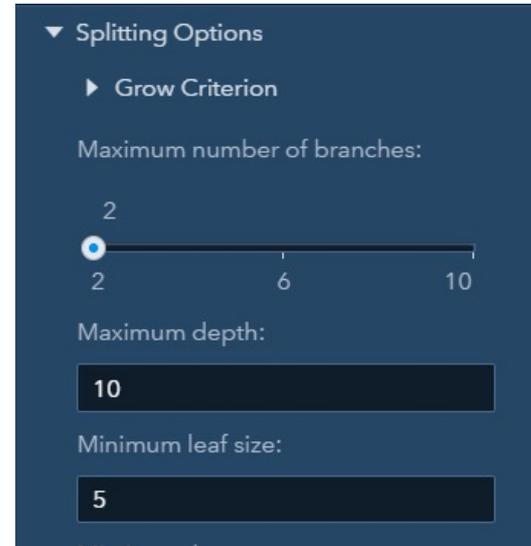


↓
Сложность = $f(\# \text{ листов, } \# \text{ ветвей, глубина})$

КОГДА НУЖНО ПРЕКРАТИТЬ РОСТ ДЕРЕВА: ОБРУБАНИЕ / PRUNING

1. Сверху вниз (Top-Down) – критерии остановки роста:

- Размер листа
- Глубина дерева



▼ Splitting Options

▶ Grow Criterion

Maximum number of branches:

2

2 6 10

Maximum depth:

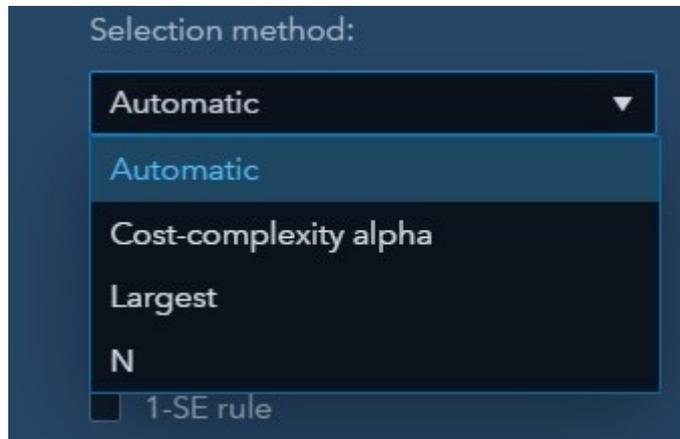
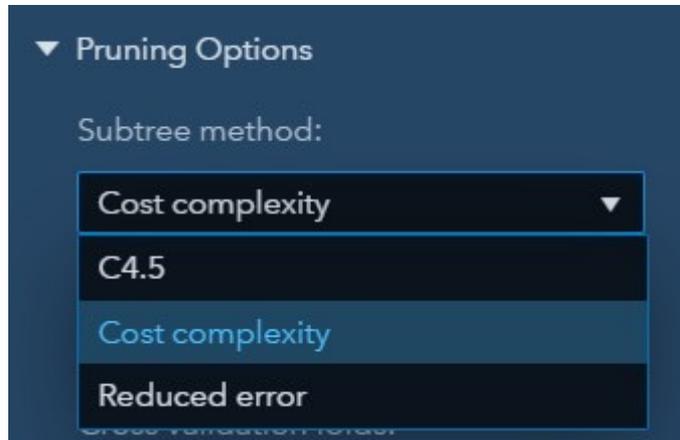
10

Minimum leaf size:

5

КОГДА НУЖНО ПРЕКРАТИТЬ РОСТ ДЕРЕВА: ОБРУБАНИЕ / PRUNING

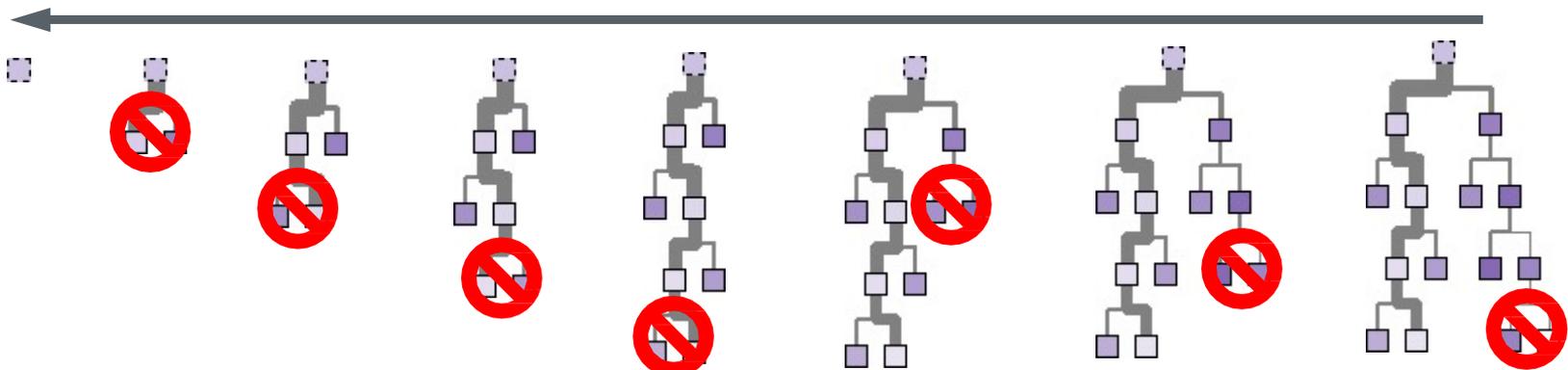
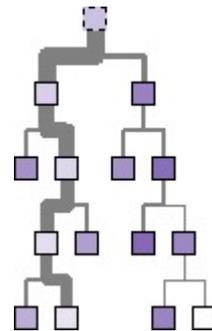
2. Снизу вверх (Bottom-Up) – критерий выбора «лучшего» дерева
- Точность модели
 - На валидационном наборе
 - Используя CV

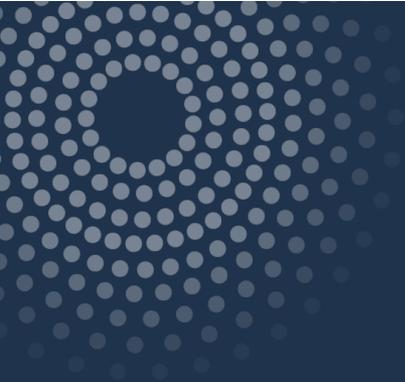


КОГДА НУЖНО ПРЕКРАТИТЬ РОСТ ДЕРЕВА: ОБРУБАНИЕ / PRUNING

2. Снизу вверх (Bottom-Up) – критерий выбора «лучшего» дерева

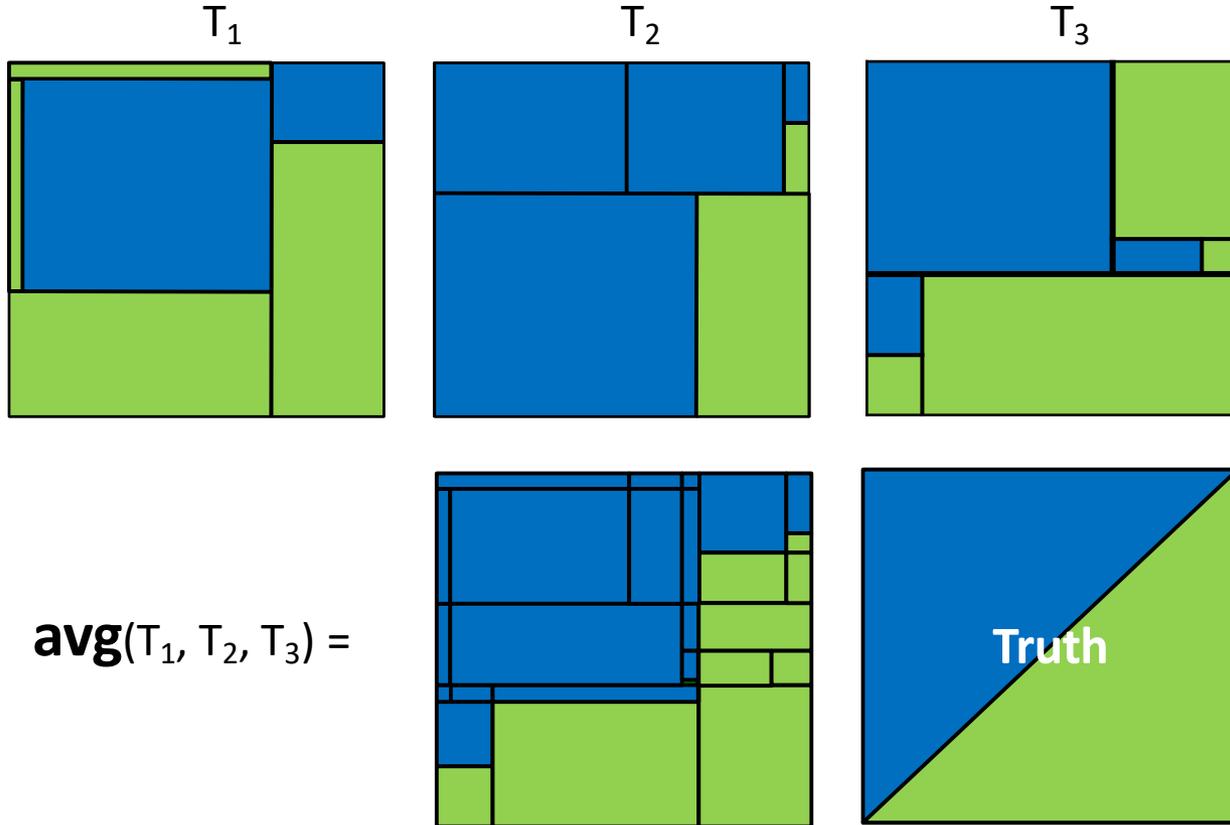
- Вырастить максимальное дерево
- Обрубая по одному листу (*min. valid-error* OR *min. CV-error*), построить последовательность деревьев:





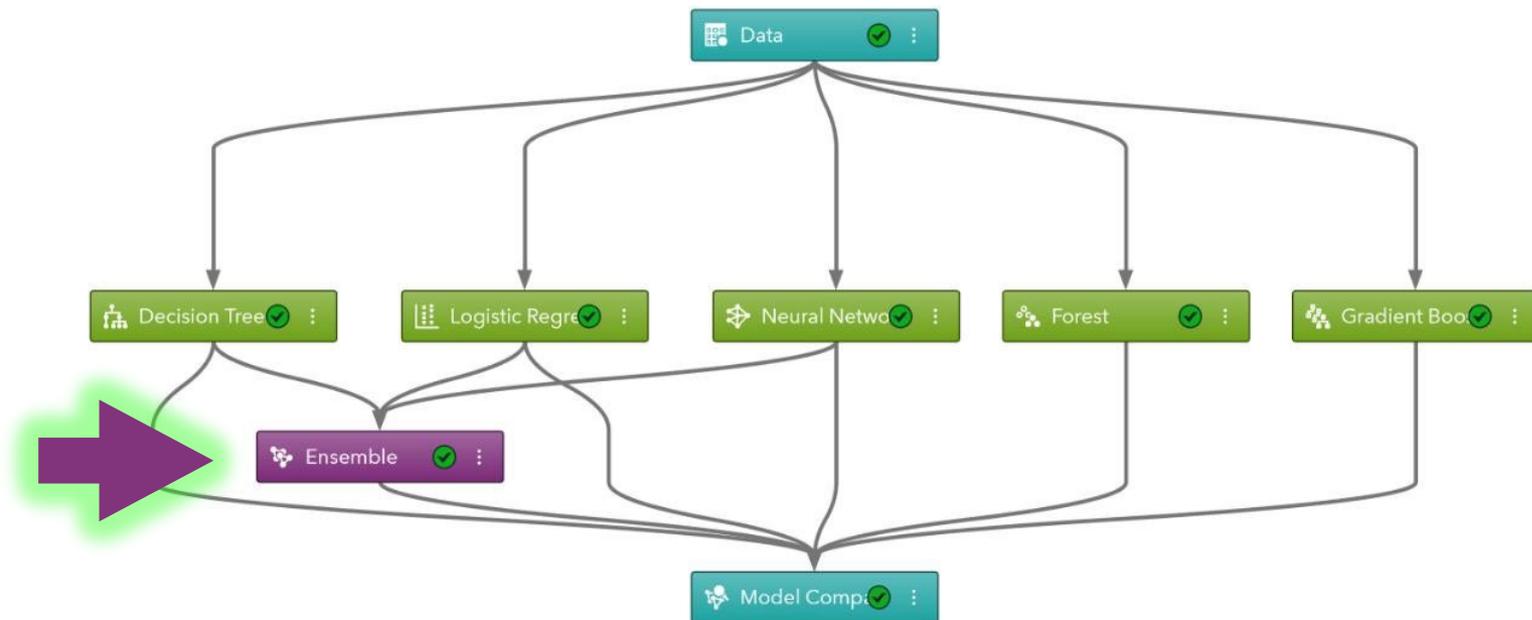
Ансамбли моделей. Случайный лес.
Деревья как вспомогательный инструмент.

АНСАМБЛЬ – КОМБИНАЦИЯ МОДЕЛЕЙ

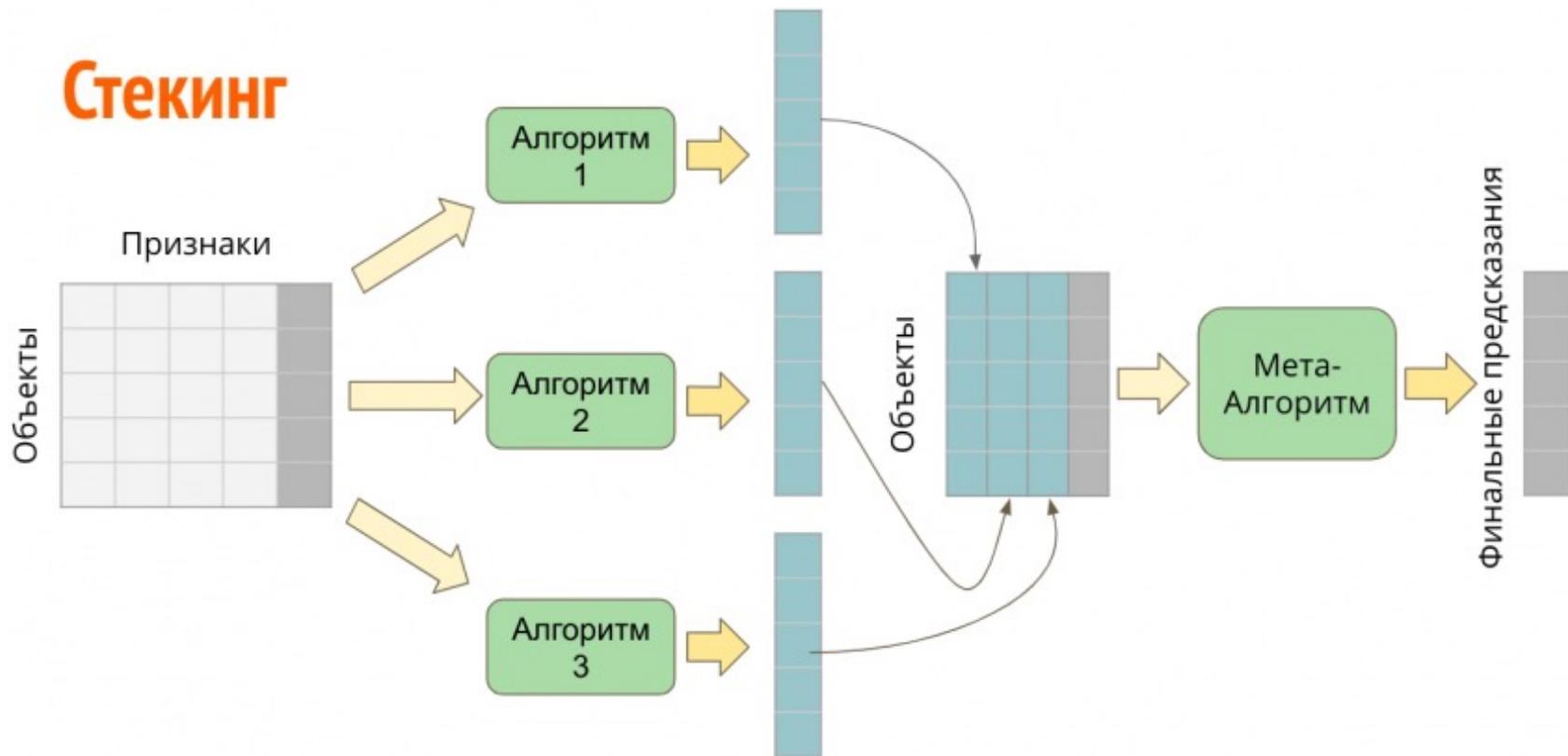


АНСАМБЛЬ – КОМБИНАЦИЯ МОДЕЛЕЙ

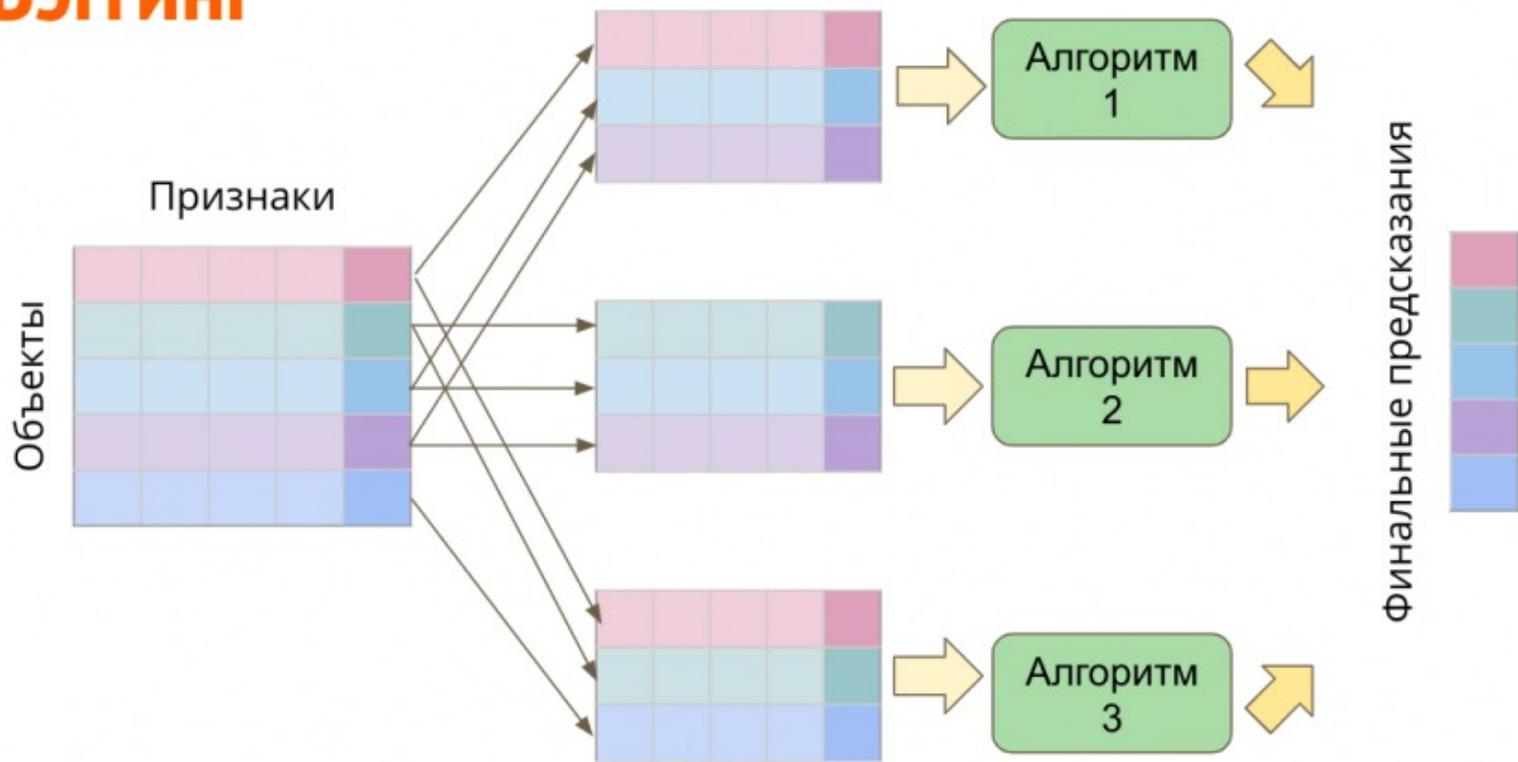
ПРОСТОЕ «ГОЛОСОВАНИЕ»: среднее вероятности на выходе моделей
(арифметическое и геометрическое среднее), максимум вероятности



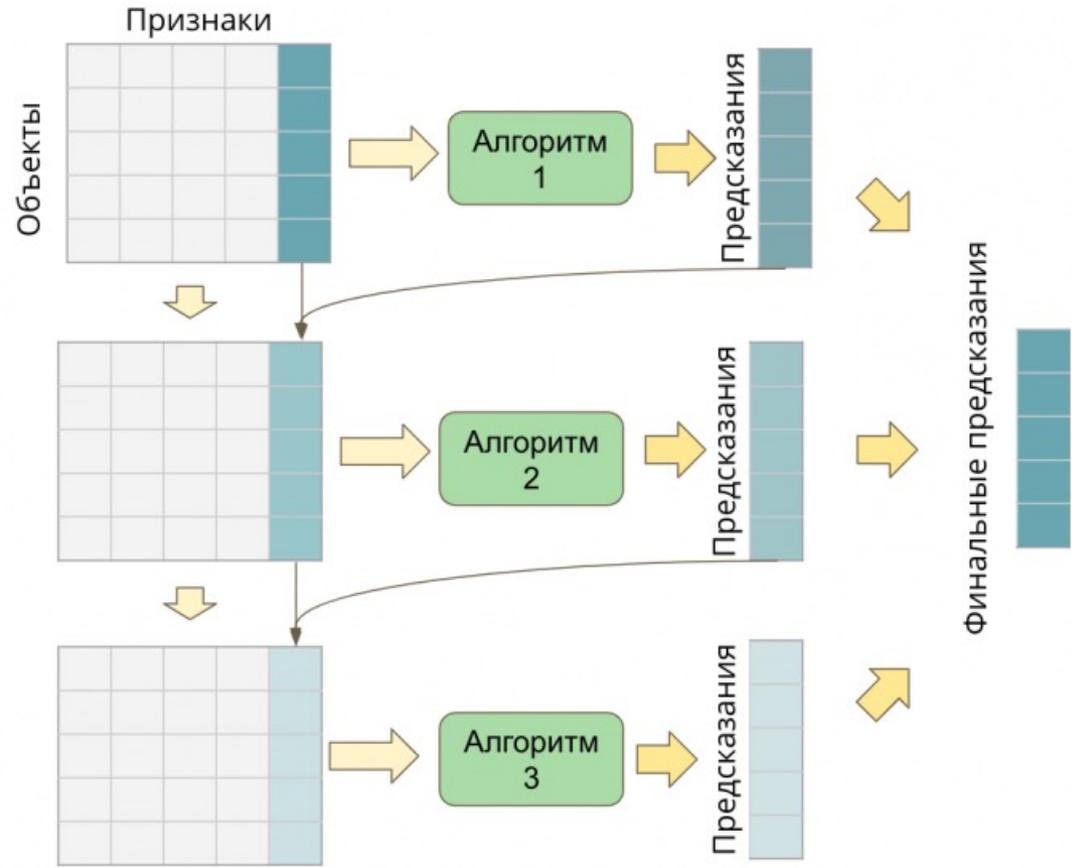
Стекинг



Бэггинг



Бустинг



УЛЕЗ ENSEMBLE: НАСТРОЙКИ

Ensemble  

Description:

Creates a new model by taking a function of posterior probabilities (for class targets) or the predicted values

▼ Interval Target

Predicted values:

Average ▼

▼ Average

Maximum

Posterior probabilities:

Average

class targets) or the predicted values ▼

▼ Interval Target

Predicted values:

Average ▼

▼ Class Target

Posterior probabilities:

Average ▼

▶ Average

▶ Geometric Mean

Maximum



Random forest

RANDOM FOREST – СЛУЧАЙНЫЙ ЛЕС

- ЛЕС

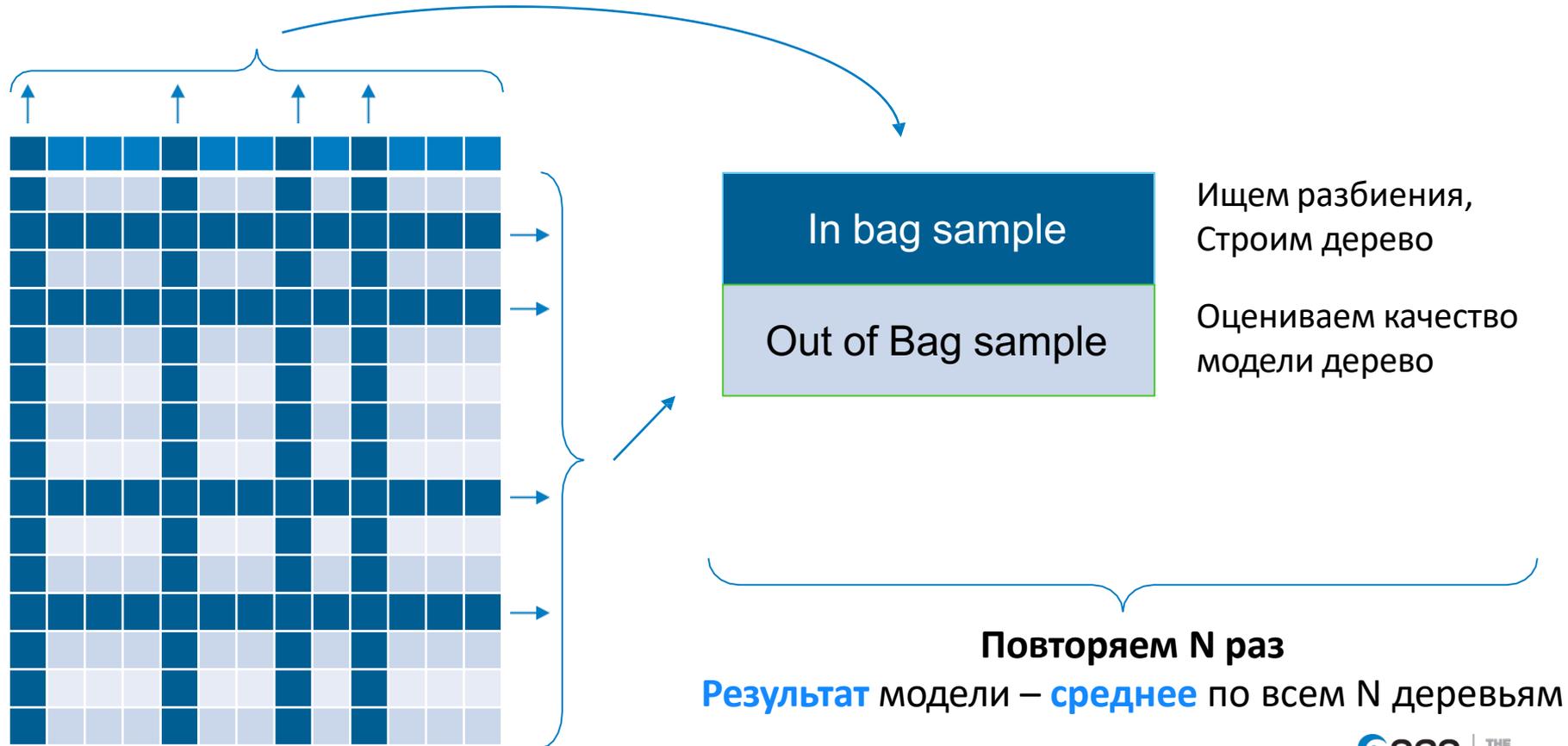
- Строим много деревьев

- СЛУЧАЙНЫЙ

- Используем (случайную) подвыборку наблюдений тренировочного набора
- Используем (случайную) подвыборку переменных

Хороший лес - много разнообразных деревьев

RANDOM FOREST



RANDOM FOREST

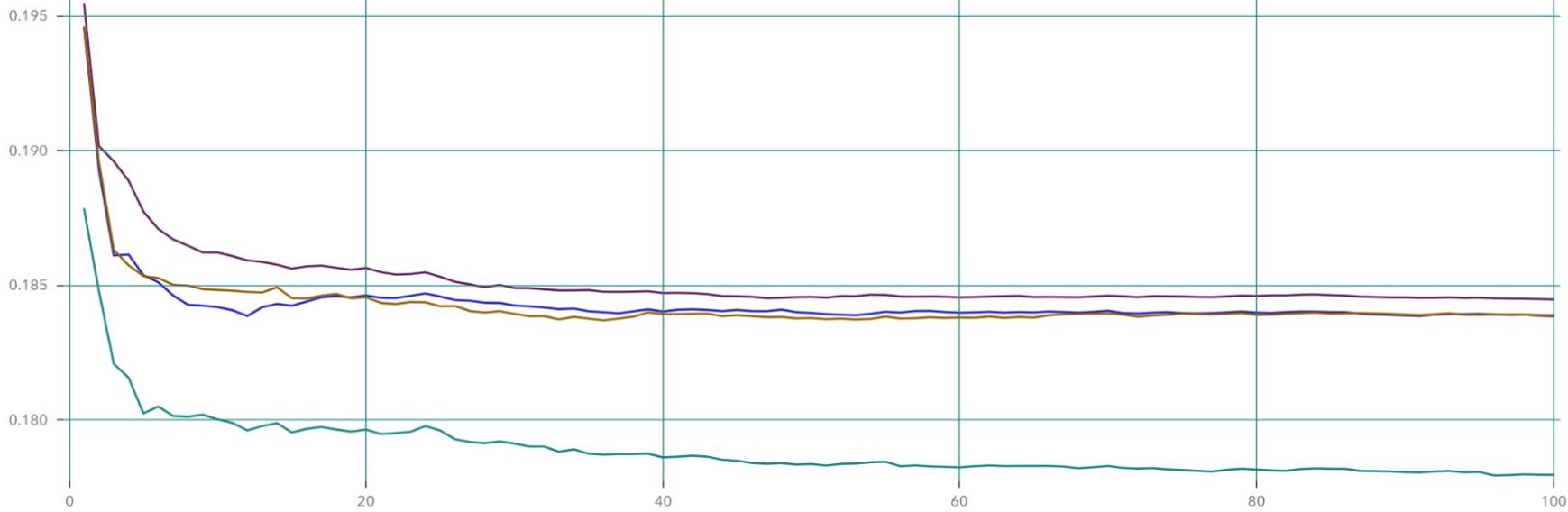
- Сколько нужно деревьев?
 - Больше деревьев – меньше дисперсия – надежнее предсказание
- Сколько нужно переменных?
 - Default - $\sqrt{N_{vars}}$
 - Много сильных предикторов – делаем меньше
 - Много слабых предикторов – делаем больше
 - дополнительно (опция) из выбранных случайно отбираем с лучшим p-value хи-квадрат (спец. вариант)

RANDOM FOREST

Error Plot

Average Squared Error

Average Squared Error



Legend: TRAIN (teal), VALIDATE (blue), OUT OF BAG (purple), TEST (orange)

RANDOM FOREST

Variable Importance

Variable Name	Train Importance	Relative Importance
FREQUENCY_STATUS_97NK	12.7705	1
CARD_PROM_12	10.3678	0.8119
RECENT_RESPONSE_COUNT	9.9015	0.7753
RECENT_CARD_RESPONSE_COUNT	9.5688	0.7493
RECENT_CARD_RESPONSE_PROP	6.8346	0.5352
WEALTH_RATING	6.5134	0.5100
MONTHS_SINCE_LAST_GIFT	5.3773	0.4211
INCOME_GROUP	5.1718	0.4050
RECENT_RESPONSE_PROP	4.8848	0.3825
MONTHS_SINCE_LAST_PROM_RESP	4.5640	0.3574
NUMBER_PROM_12	4.3387	0.3397
MEDIAN_HOUSEHOLD_INCOME	3.9939	0.3127
PCT_ATTRIBUTE4	3.9486	0.3092
LIFETIME_CARD_PROM	3.6924	0.2891
LIFETIME_GIFT_COUNT	3.6143	0.2830
PCT_ATTRIBUTE2	3.5952	0.2815
URBANITY	3.5138	0.2752



Применение в качестве вспомогательного инструмента

Применение в
качестве
вспомогательного
инструмента

- **Выявление аномалий**
- **Уменьшение размерности**
 - Выбор важных / значимых переменных для других методов моделирования
 - Группировка значений категориальных переменных
- **Преобразование входных данных**
 - Нелинейные пересечения переменных
 - Замена пропущенных значений
 - Дискретизация непрерывных переменных

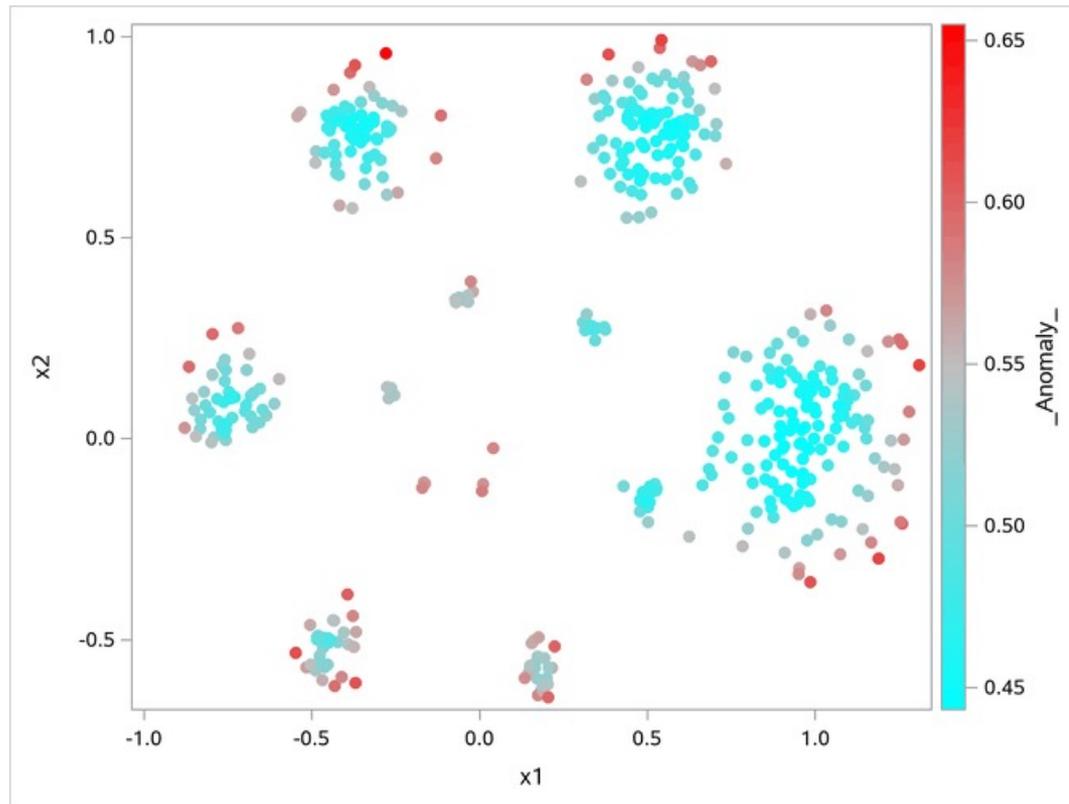
Случайно выбираем одну переменную:

- **Интервальная переменная:**
разбиение по случайному значению
- **Номинальная:** каждое значение
отправляется в случайную ветку/лист

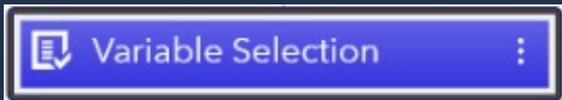
При таком построении аномалии будут скорее всего ближе к корню, чем нормальные наблюдения. Считаем ***anomaly score*** для каждого наблюдения x :

$$s(x) = 2\{-h(x)\}$$

$h(x)$ - средняя (по всем деревьям)
длина пути от корня до листа с
наблюдением x



ВЫБОР ВАЖНЫХ / ЗНАЧИМЫХ ПЕРЕМЕННЫХ ДЛЯ ДРУГИХ МЕТОДОВ



Относительная важность

Variable Selection  

Description:

Performs unsupervised and several supervised methods of variable selection to reduce the number of

▶ Pre-screen Input Variables

Combination criterion:

Selected by at least 1

▶ Unsupervised Selection

▶ Fast Supervised Selection

▶ Linear Regression Selection

▶ Decision Tree Selection

▶ Forest Selection

▼ Gradient Boosting Selection

▼ Tree-splitting Options



ДИСКРЕТИЗАЦИЯ НЕПРЕРЫВНЫХ ПЕРЕМЕННЫХ

Transformations  

Description:

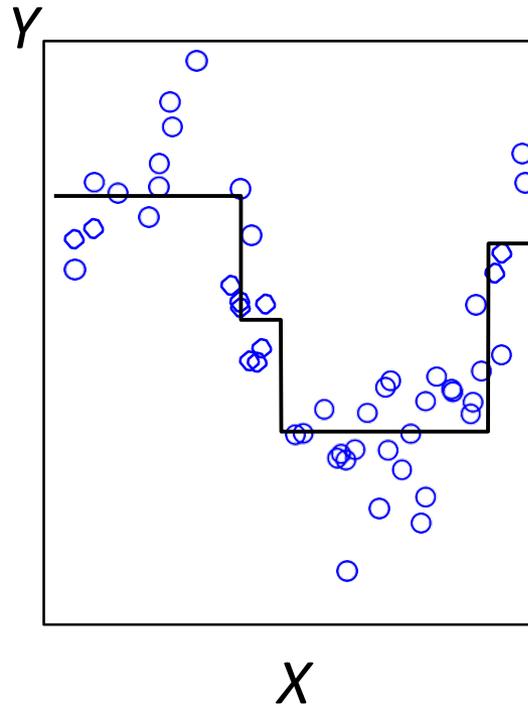
Applies numerical or binning transformations to input variables.

Interval Inputs

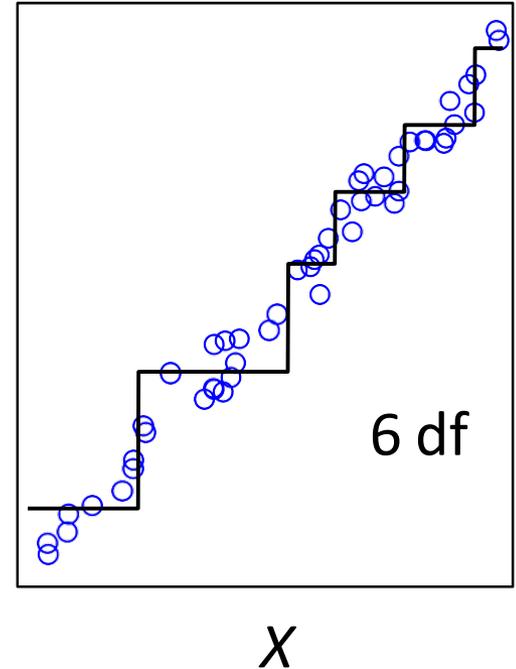
Default interval inputs method:

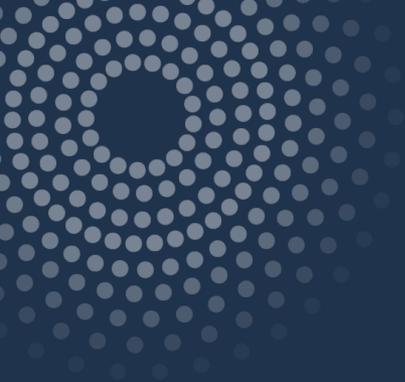
- None
- Exponential
- Inverse
- Inverse square
- Inverse square root
- Log
- Log10
- None**
- Quantile binning
- Range standardization
- Square
- Square root
- Standardization
- Tree-based binning

Mis Se



Dimension Inflation





Спасибо!

sas.com