

The image features a dark, textured background. Three paper airplanes are scattered across the frame: one bright yellow one is positioned in the upper right, and two grey ones are in the lower left and bottom right. A dashed white line, resembling chalk, winds through the scene, starting from the bottom left, looping around the grey airplane, extending towards the yellow one, and then looping back towards the bottom right. The title text is centered in the lower half of the image.

Градиентный бустинг: имплементации

МАКСИМОВСКАЯ
АНАСТАСИЯ

Extreme Gradient Boosting

XGBoost

В градиентном бустинге на каждой итерации вычисляется вектор сдвигов s :

$$s = \left(- \frac{\partial L}{\partial z} \Big|_{z=a_{N-1}(x_i)} \right)_{i=1}^{\ell} = -\nabla_z \sum_{i=1}^{\ell} L(y_i, z_i) \Big|_{z_i=a_{N-1}(x_i)}$$

После этого обучается новый базовый алгоритм:

$$b_N(x) = \arg \min_{b \in \mathcal{A}} \sum_{i=1}^{\ell} (b(x_i) - s_i)^2$$

$$a_N = \sum b_i(x)$$

y-target

$$1) b_1(x) \quad \frac{1}{p} \sum (b(x_i) - y_i)^2 \rightarrow \min_{b_i}$$

$$2) \quad S_i = \cancel{(y_i - b_i(x))} \quad S_i = -\frac{\partial L}{\partial z} b(x_i) \quad b_2(x) \quad \frac{1}{q} \sum (b(x_i) - S_i)^2 \rightarrow \min_{b_i}$$

$$a(x) = b_1(x) + b_2(x)$$

$$3) z_i = y_i - b_1(x_i) - b_2(x_i)$$

XGBoost

В XGBoost на каждом шаге решается задача:

$$\sum_{i=1}^l \left(-s_i b(x_i) + \frac{1}{2} h_i b^2(x_i) \right) + \gamma J + \frac{\lambda}{2} \sum_{j=1}^J b_j^2 \rightarrow \min_b$$
$$h_i = \left. \frac{\partial^2 L}{\partial z^2} \right|_{a_{N-1}(x_i)}$$

XGBoost

Мы хотим найти алгоритм $b(x)$, решающий следующую задачу:

$$\sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b(x_i)) \rightarrow \min_b$$

XGBoost

Разложим функцию L в каждом слагаемом в ряд Тейлора до второго члена с центром в ответе композиции $a_{N-1}(x_i)$:

$$\begin{aligned} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b(x_i)) &\approx \\ &\approx \sum_{i=1}^{\ell} \left(L(y_i, a_{N-1}(x_i)) - s_i b(x_i) + \frac{1}{2} h_i b^2(x_i) \right), \end{aligned}$$

XGBoost

Где h_i – производные по сдвигам:

$$h_i = \frac{\partial^2}{\partial z^2} L(y_i, z) \Big|_{a_{N-1}(x_i)}$$

XGBoost

Выкидываем первое слагаемое, т.к. оно не зависит от нового базового алгоритма:

$$\sum_{i=1}^{\ell} \left(-s_i b(x_i) + \frac{1}{2} h_i b^2(x_i) \right) \rightarrow \min_b$$

XGBoost

$$\begin{aligned} & \sum_{i=1}^{\ell} (b(x_i) - s_i)^2 \\ &= \sum_{i=1}^{\ell} (b^2(x_i) - 2s_i b(x_i) + s_i^2) = \{\text{последнее слагаемое не зависит от } b\} \\ &= \sum_{i=1}^{\ell} (b^2(x_i) - 2s_i b(x_i)) \\ &= 2 \sum_{i=1}^{\ell} \left(-s_i b(x_i) + \frac{1}{2} b^2(x_i) \right) \rightarrow \min_b \end{aligned}$$

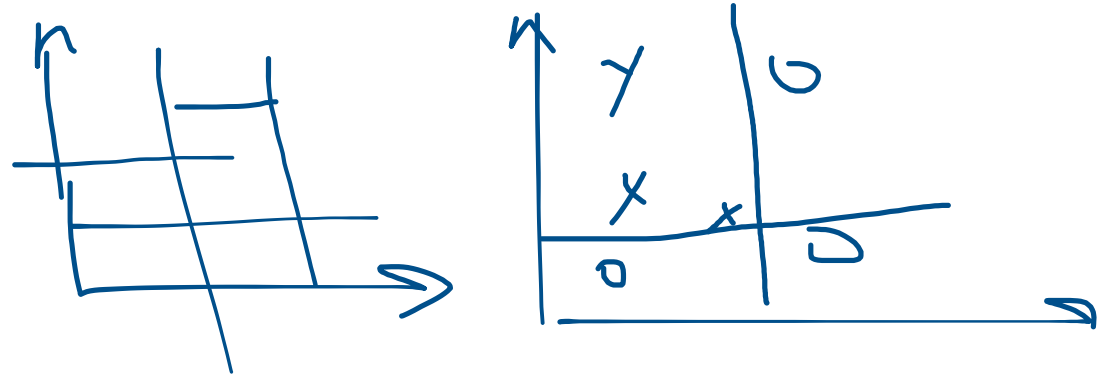
XGBoost

$$b(x) = \sum w_i [x \in R_i]$$

В XGBoost на каждом шаге решается задача:

$$\sum_{i=1}^l \left(-s_i b(x_i) + \frac{1}{2} h_i b^2(x_i) \right) + \gamma J + \frac{\lambda}{2} \sum_{j=1}^J b_j^2 \rightarrow \min_b$$
$$h_i = \left. \frac{\partial^2 L}{\partial z^2} \right|_{a_{N-1}(x_i)}$$

XGBoost

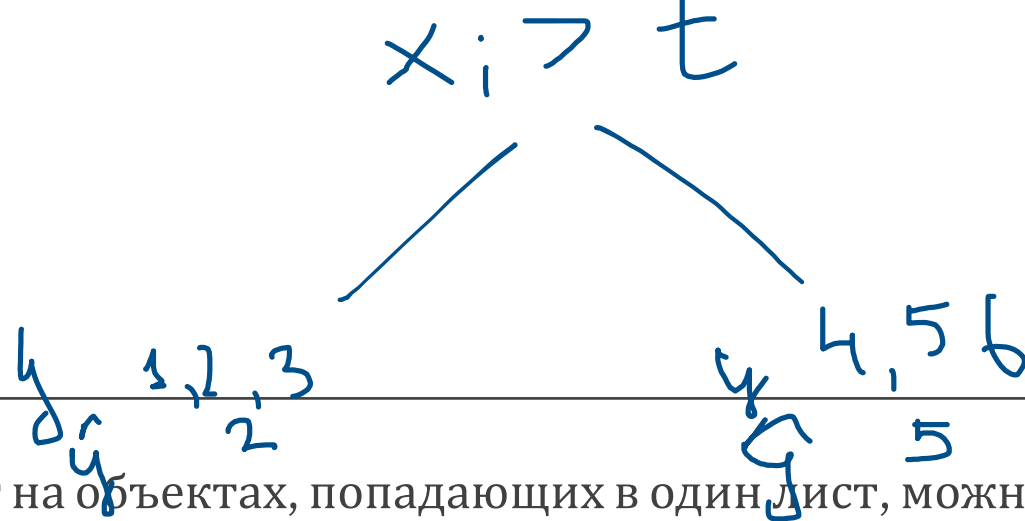


Сложность дерева $b(x)$ зависит от:

- Число листьев J . Чем больше листьев имеет дерево, тем сложнее его разделяющая поверхность, тем больше у него параметров и тем выше риск переобучения
- Норма коэффициентов в листьях. Чем сильнее коэффициенты отличаются от нуля, тем сильнее данный базовый алгоритм будет влиять на итоговый ответ композиции.

$$\|b\|_2^2 = \sum_{j=1}^J b_j^2$$

XGBoost



Т.к. дерево дает одинаковый ответ на объектах, попадающих в один лист, можно упростить функционал:

$$\sum_{j=1}^J \left\{ \underbrace{\left(- \sum_{i \in R_j} s_i \right)}_{=-S_j} b_j + \frac{1}{2} \left(\lambda + \underbrace{\sum_{i \in R_j} h_i}_{=H_j} \right) b_j^2 + \gamma \right\} \rightarrow \min_b$$

$$b = \frac{S_j}{H_j + \lambda}$$

XGBoost

Отдельное слагаемое представляет собой параболу относительно b_j , благодаря чему можно аналитически найти оптимальные коэффициенты в листьях:

$$b_j = \frac{S_j}{H_j + \lambda}$$

XGBoost

Ошибка дерева с оптимальными коэффициентами:

$$H(b) = -\frac{1}{2} \sum_{j=1}^J \frac{S_j^2}{H_j + \lambda} + \gamma J$$

XGBoost

Этот функционал подходит в качестве критерия информативности:

$$Q = H(R) - H(R_\ell) - H(R_r) \rightarrow \max.$$

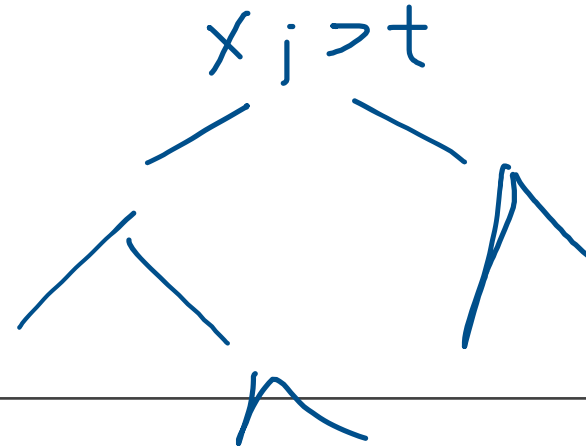
$$H(R) = -\frac{1}{2} \left(\sum_{(h_i, s_i) \in R} s_j \right)^2 / \left(\sum_{(h_i, s_i) \in R} h_j + \lambda \right) + \gamma.$$

XGBoost

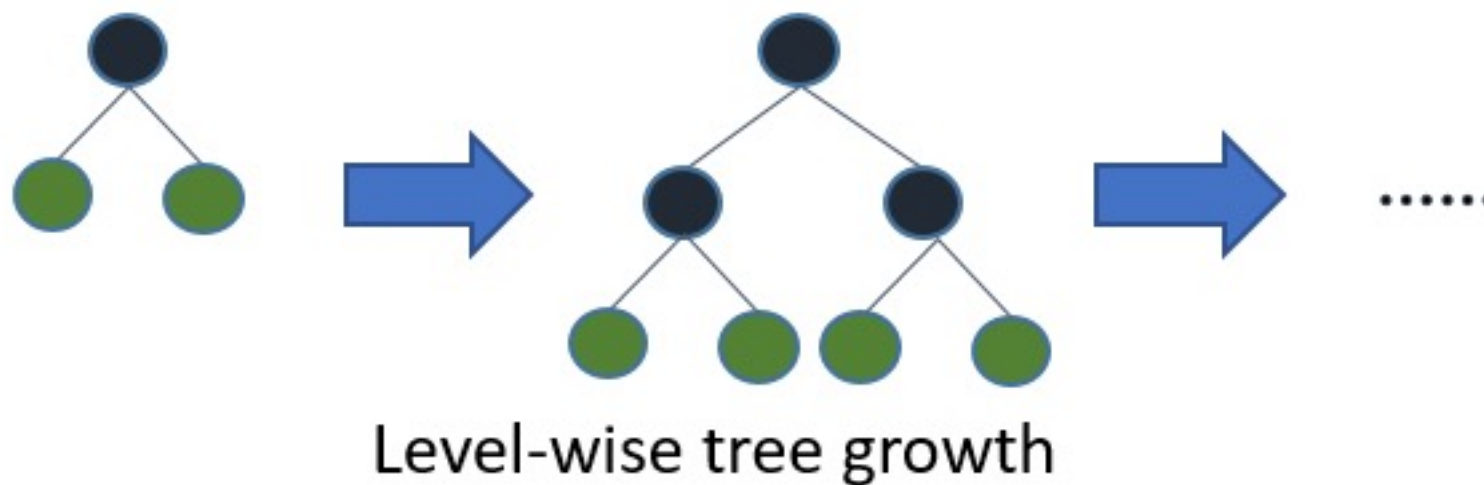
- Базовый алгоритм приближает направление, посчитанное с учетом вторых производных функции потерь;
- Функционал регуляризуется — добавляются штрафы за количество листьев и за норму коэффициентов;
- При построении дерева используется критерий информативности, зависящий от оптимального вектора сдвига;
- Критерий останова при обучении дерева также зависит от оптимального сдвига.

lightGBM

lightGBM

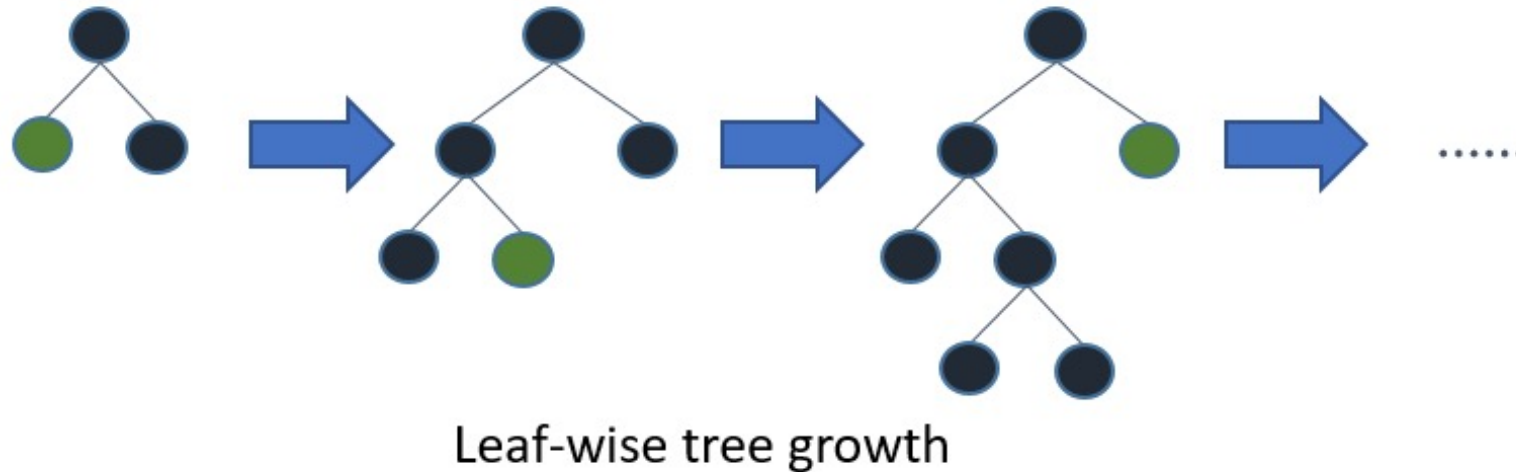


В других реализациях градиентного бустинга деревья строятся по уровням:



lightGBM

LightGBM строит деревья, добавляя на каждом шаге один лист. Это позволяет добиться более высокой точности решения задачи оптимизации



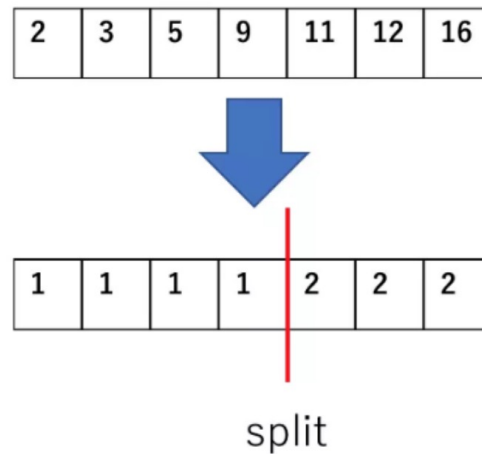
lightGBM

Кодирование категориальных признаков.

- LightGBM разбивает значения категориального признака на два подмножества в каждой вершине дерева, находя при этом наилучшее разбиение
- Если категориальный признак имеет k различных значений, то возможных разбиений $2^{k-1} - 1$. В LightGBM реализован способ поиска оптимального разбиения за $O(k \log k)$ операций.

lightGBM

Ускорение построения деревьев за счёт бинаризации признаков:



An example of how binning can reduce the number of splits to explore. The features must be sorted in advance for this method to be effective.

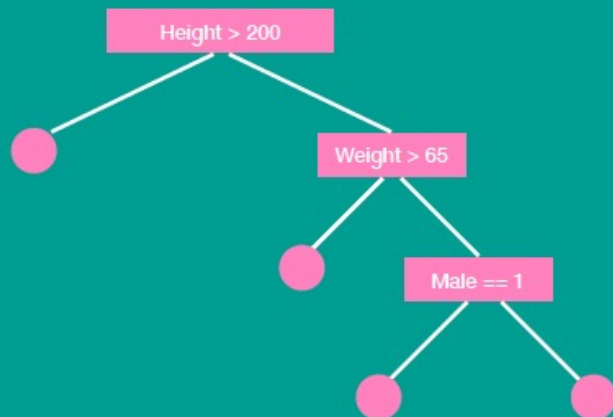
CatBoost

CatBoost

- Алгоритм Яндекса, является оптимизацией Xgboost и в отличие от Xgboost умеет обрабатывать категориальные признаки.

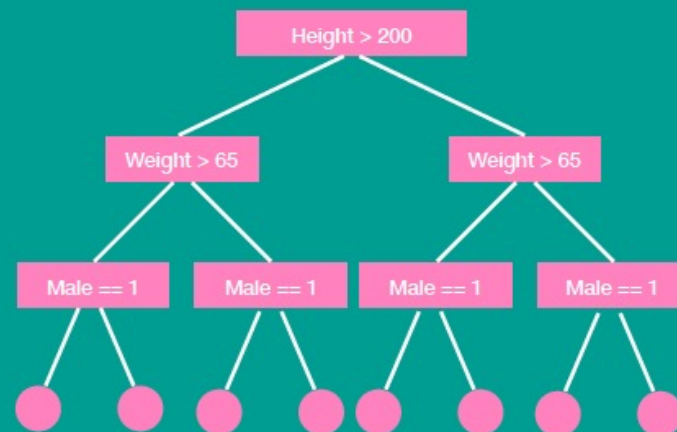
Симметричные деревья

Asymmetric trees



XGBoost
LightGBM

Oblivious trees



CatBoost

Выбор лучшего сплита

$$\text{score}(\text{split}) = \frac{\sum_{doc} \text{leafValue}(doc) * \text{gradient}(doc) * w(doc)}{\sqrt{\sum_{doc} w(doc) * \text{leafValue}(doc)^2}}$$

$$\text{leafValue}(doc) = \frac{\text{sumWeightedDer}}{\text{sumWeights}}$$

Бутстрап

Бернулли:

$$w(doc) = 0 \text{ or } 1 \ (P(1) = sample_rate)$$

Байесовский:

$$w(doc) *= (-\log(rand(0,1)))^{bagging_temperature}$$

Только на этапе выбора структуры дерева

Числовые факторы

Рост > 170 см?



Да



Нет

Категориальные факторы



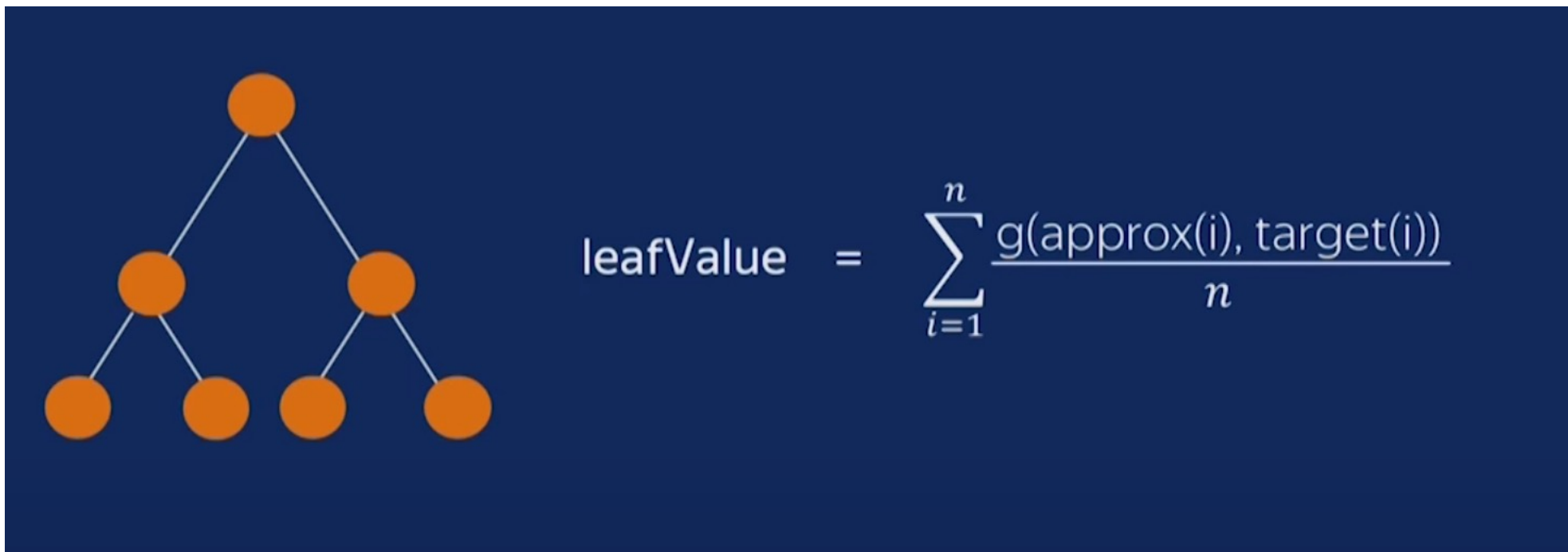
Статистики по категориальным факторам

- › One-hot кодирование
- › Статистики без использования таргета
- › Статистики по случайным перестановкам
- › Комбинации факторов

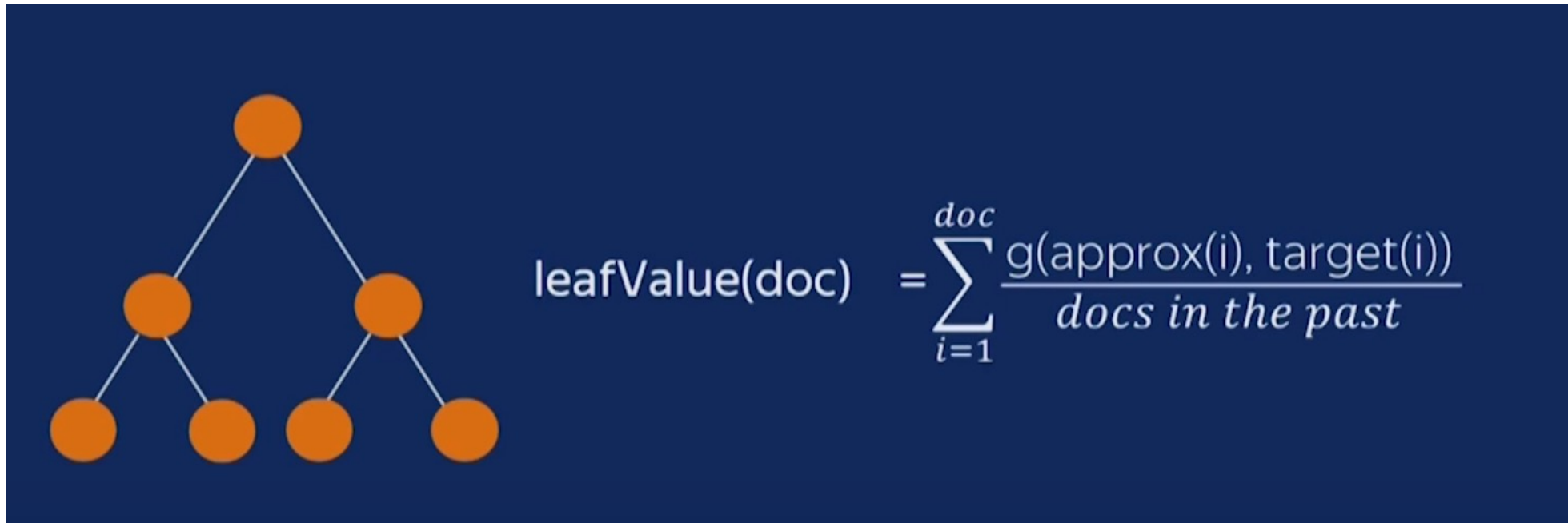
прошлое		SDE		1
		SDE		1
		SDE		0
		PR		
	i	SDE		1
		PR		

$$i \rightarrow \frac{1+1+0}{3}$$

Обычный бустинг



Динамический бустинг



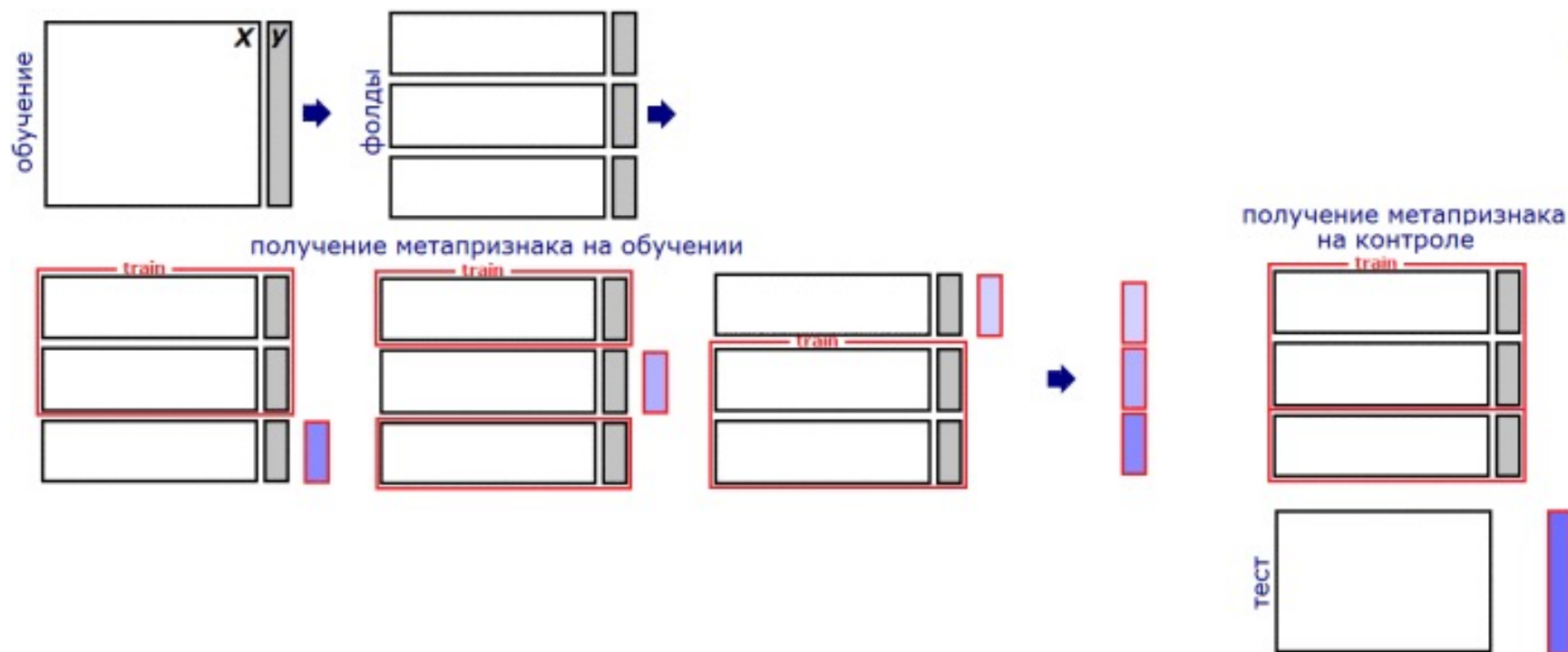
Как еще можно
строить композиции?

Стекинг

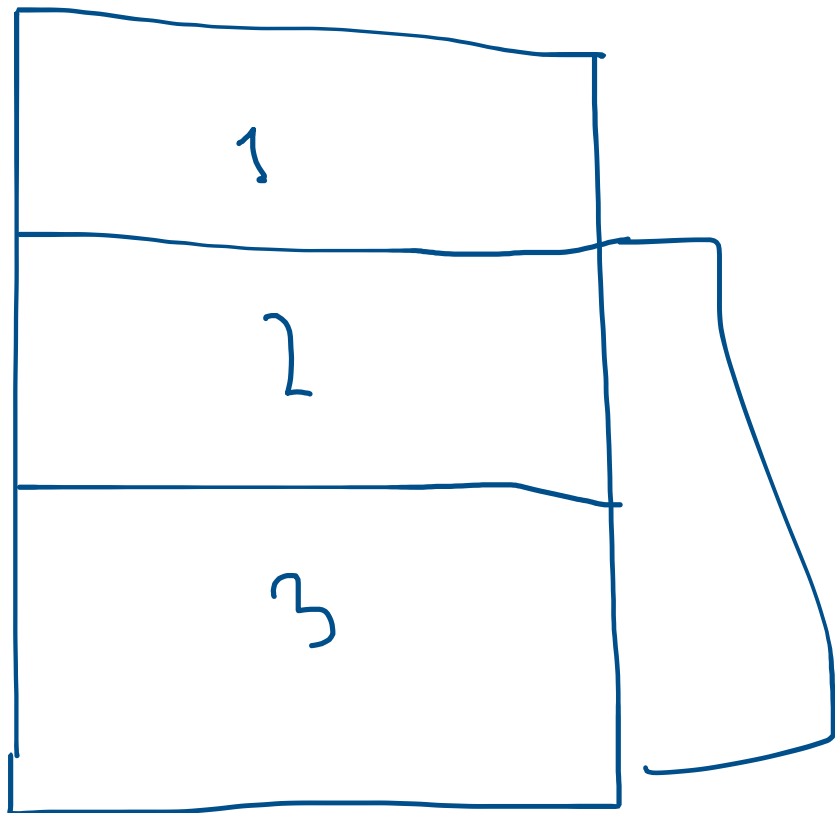
$$N \quad b_1(x) \dots b_N(x) \quad a(x)$$

$$\sum L(y_i, a(b_1(x), b_2(x), \dots, b_N(x))) \rightarrow \min_q$$

- Прогнозы алгоритмов объявляются новыми признаками, и поверх них обучается ещё один алгоритм (который иногда называют мета-алгоритмом)



X



$$b_i^{-1}(x) \rightarrow x_2 x_3$$

$$\sum_{k=1}^K \sum L(y_i) \left(b_1^{-k}(x), b_2^{-k}(x), \dots \right)$$

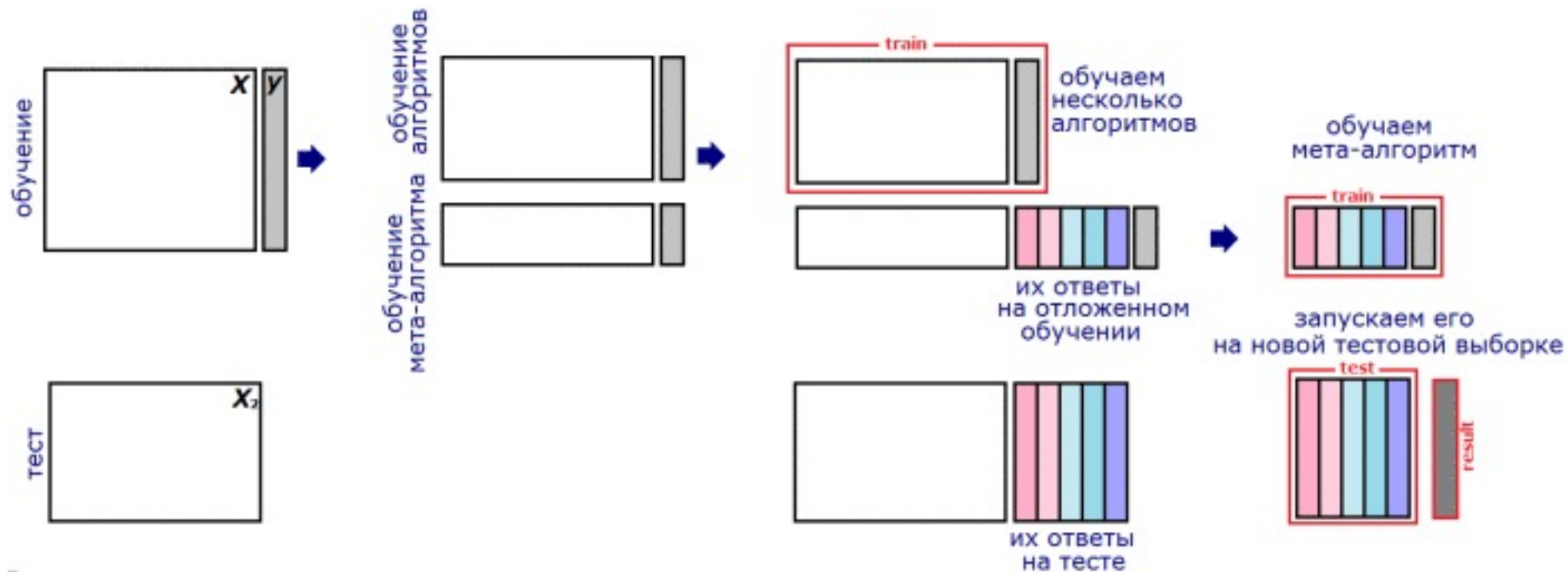
Блендинг

- Частным случаем стекинга является блендинг, в котором мета-алгоритм является линейным:

$$a(x) = \sum_{n=1}^N w_n b_n(x)$$

Это самый простой способ объединить несколько алгоритмов в композицию. Иногда даже блендинг без обучения весов (то есть вариант с $w_1 = \dots = w_n = 1/N$) позволяет улучшить качество по сравнению с отдельными базовыми алгоритмами.

Блендинг



Дополнительные материалы

- [Статья про стекинг и блендинг \(рус.\)](#)
- [Про приближение функции потерь в XGBoost рядом Тейлора \(англ.\)](#)
- [Реализация XGBoost с нуля \(рус.\)](#)
- [Сравнение CatBoost, lightGBM, XGBoost \(англ.\)](#)
- [Простой пример использования lightGBM \(рус.\)](#)
- [Примеры использования lightGBM из документации \(англ.\)](#)
- [Советы по настройке параметров lightGBM \(англ.\)](#)

Дополнительные материалы

- [Очень классная лекция про CatBoost \(рус.\)](#)
- [Тutorial по CatBoost \(англ.\)](#)
- [Тutorial по CatBoost \(рус.\)](#)
- [Тutorial по CatBoost с заданиями \(англ.\)](#) За решение заданий до конца апреля можно получить 5 бонусных баллов.
- [Лекция про AdaBoost \(англ.\)](#)