

---

# Занятие № 9

## Ансамблирование моделей



---

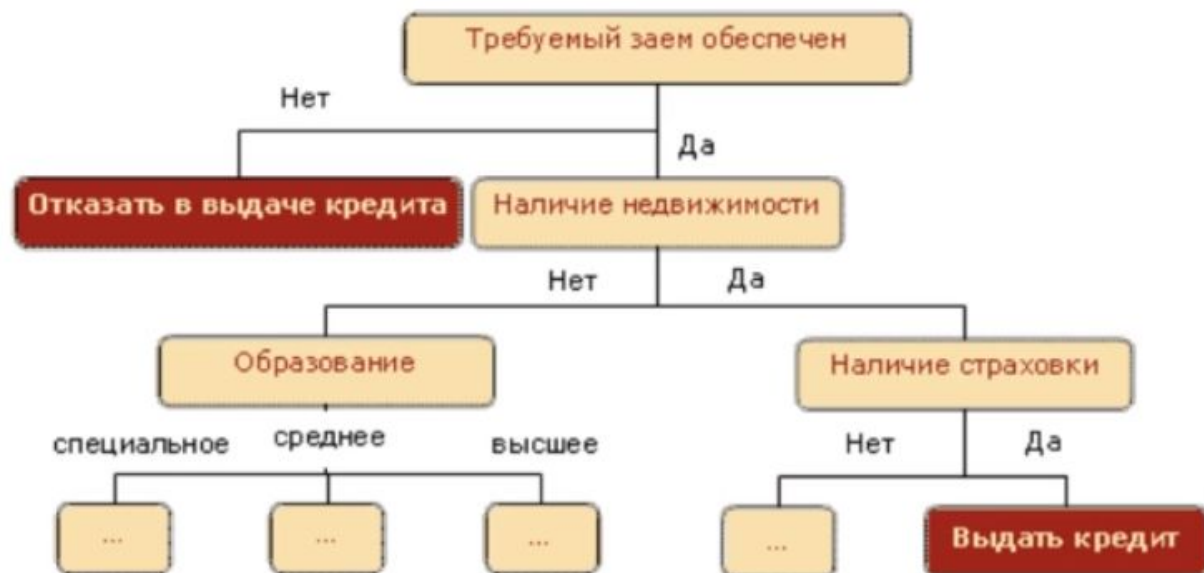
# Содержание

---

- 1 Введение
- 2 Что не так с деревьями?
- 3 Ансамбли
- 4 Бутстреп. Бэггинг. Случайный лес. Блэндинг. Стекинг. Бустинг
- 5 Практика.



# Введение



# Что не так с деревьями?

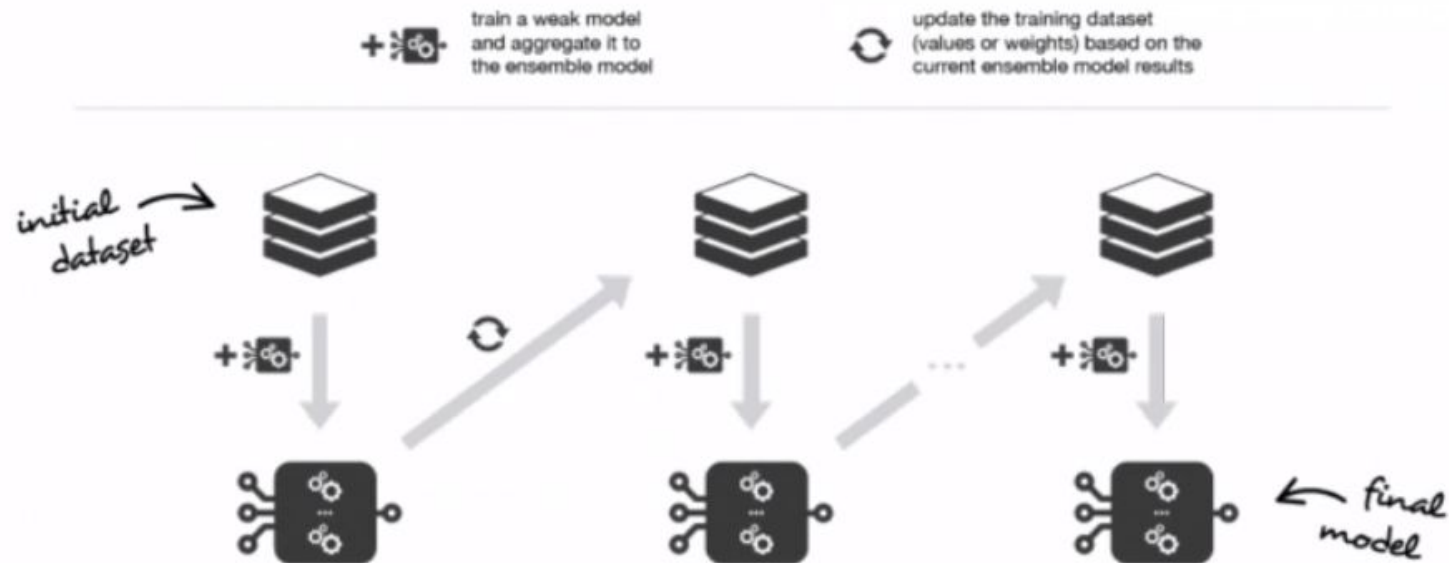
1. Острая проблема переобучения
2. Неустойчивость
3. Не учитывает нелинейные зависимости или даже простые линейные, которые идут не по осям координат (f.e., представьте дерево для классификатора вида  $y > x$ )
4. Чувствителен к несбалансированным классам
5. Хорошо интерполирует, плохо экстраполирует



# Звать ансамбль!



# Звать ансамбль!



**«Мудрость толпы»**



# Фрэнсис Гальтон «Мудрость толпы»

Фрэнсис Гальтон в 1906 году посетил рынок, где проводилась некая лотерея для крестьян.

Их собралось около 800 человек, и они пытались угадать вес быка, который стоял перед ними. Бык весил 1198 фунтов. Ни один крестьянин не угадал точный вес быка, но если посчитать среднее от их предсказаний, то получим 1197 фунтов.

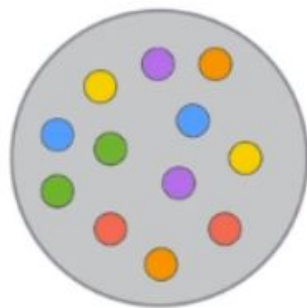
Эту идею уменьшения ошибки применили и в машинном обучении.





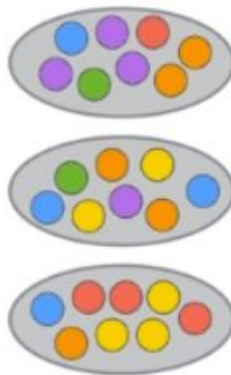
# Бутстреп

Исходная выборка



Статистика по  
выборке

Бутстрэп выборки



Статистики по  
бутстрэп выборкам

Статистика 1

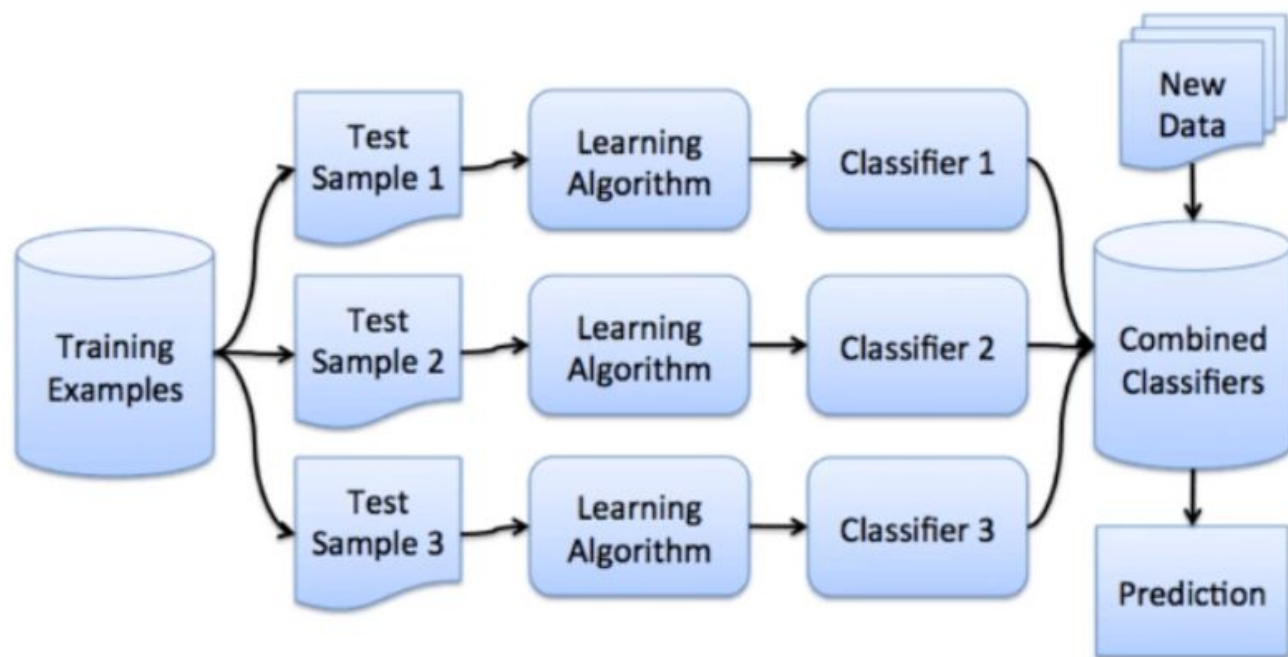
Статистика 2

Статистика 3

Бутстрэп  
распределение



# БЭГИНГ



# Бэгинг

$$\varepsilon_i(x) = b_i(x) - y(x), i = 1, \dots, n$$

$$E_x(b_i(x) - y(x))^2 = E_x \varepsilon_i^2(x).$$

$$E_1 = \frac{1}{n} E_x \sum_{i=1}^n \varepsilon_i^2(x)$$



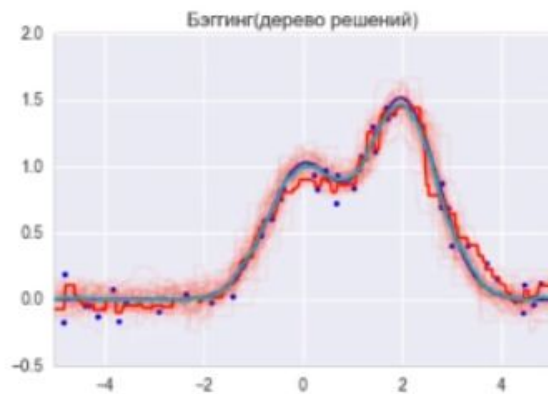
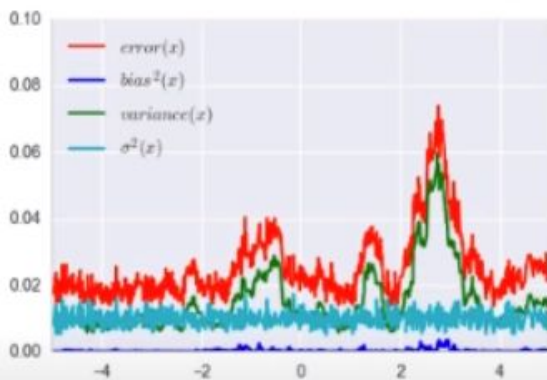
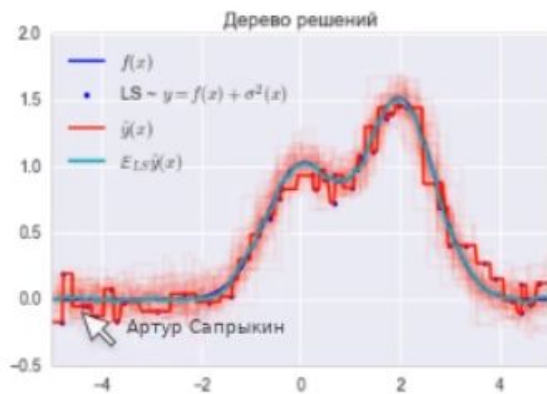
# Бэгинг

$$a(x) = \frac{1}{n} \sum_{i=1}^n b_i(x)$$

$$\begin{aligned} E_n &= E_x \left( \frac{1}{n} \sum_{i=1}^n b_i(x) - y(x) \right)^2 \\ &= E_x \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i \right)^2 \\ &= \frac{1}{n^2} E_x \left( \sum_{i=1}^n \varepsilon_i^2(x) + \sum_{i \neq j} \varepsilon_i(x) \varepsilon_j(x) \right) \\ &= \frac{1}{n} E_1 \end{aligned}$$



# Бэггинг



# Бэгинг

```
from sklearn.ensemble import BaggingRegressor
```

**base\_estimator** object or None – модель регрессии из sklearn (по умолчанию деревья решений)

**n\_estimators** – количество моделей

**max\_samples** int or float, optional (default=1.0) Количество сэмплов для обучения

**max\_features** int or float, optional (default=1.0) Количество признаков для обучения



# Случайный лес

Алгоритм построения случайного леса, состоящего из **N** деревьев, выглядит следующим образом:

Для каждого  **$n=1, \dots, N$** :

Сгенерировать выборку  **$X_n$**  с помощью бутстрэпа;

Построить решающее дерево  **$bn$**  по выборке  **$X_n$** :

— по заданному критерию мы выбираем лучший признак, делаем разбиение в дереве по нему и так до исчерпания выборки

— дерево строится, пока в каждом листе не более  **$n_{min}$**  объектов или пока не достигнем определенной высоты дерева

— при каждом разбиении сначала выбирается  **$m$**  случайных признаков из  **$n$**  исходных,

и оптимальное разделение выборки ищется только среди них.



# Случайный лес

```
from sklearn.ensemble import DecisionTreeRegressor
```

**criterion** – метод оценки ошибки (MSE по умолчанию)

**splitter** str – метод разделения "better" или "random"

**n\_estimators** – количество моделей

**max\_depth** – максимальная глубина деревьев

**max\_samples** int or float, optional (default=1.0) Количество сэмплов для обучения

**max\_features** int or float, optional (default=1.0) Количество признаков для обучения





# Случайный лес



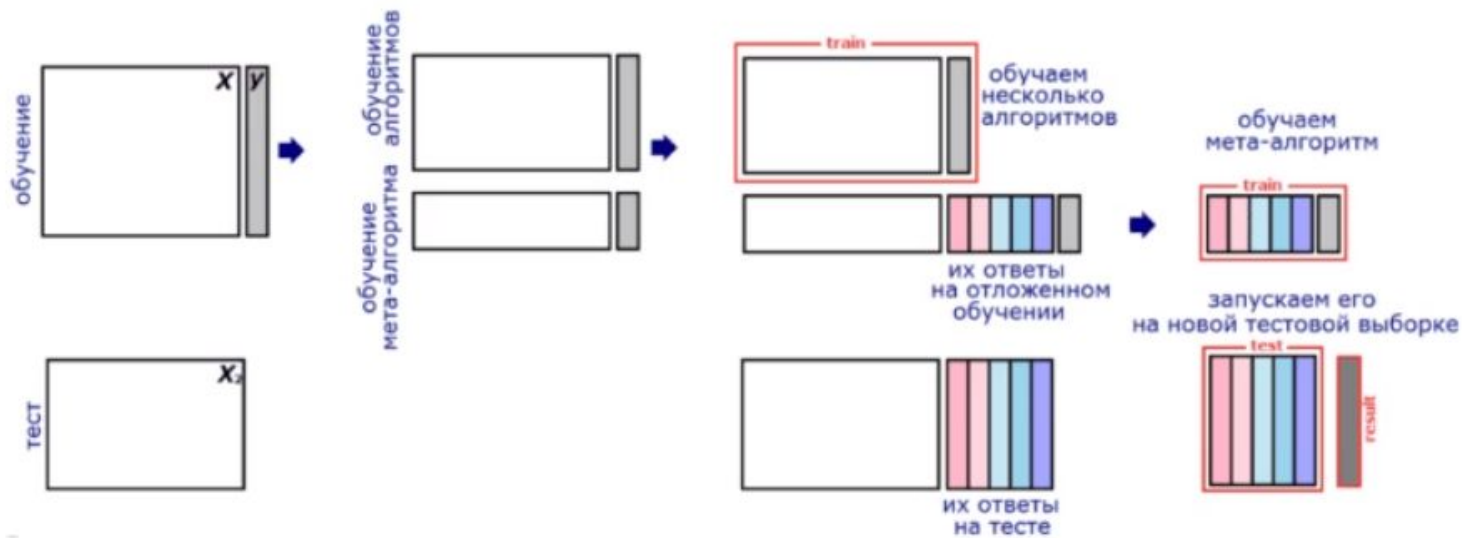
# Случайный лес



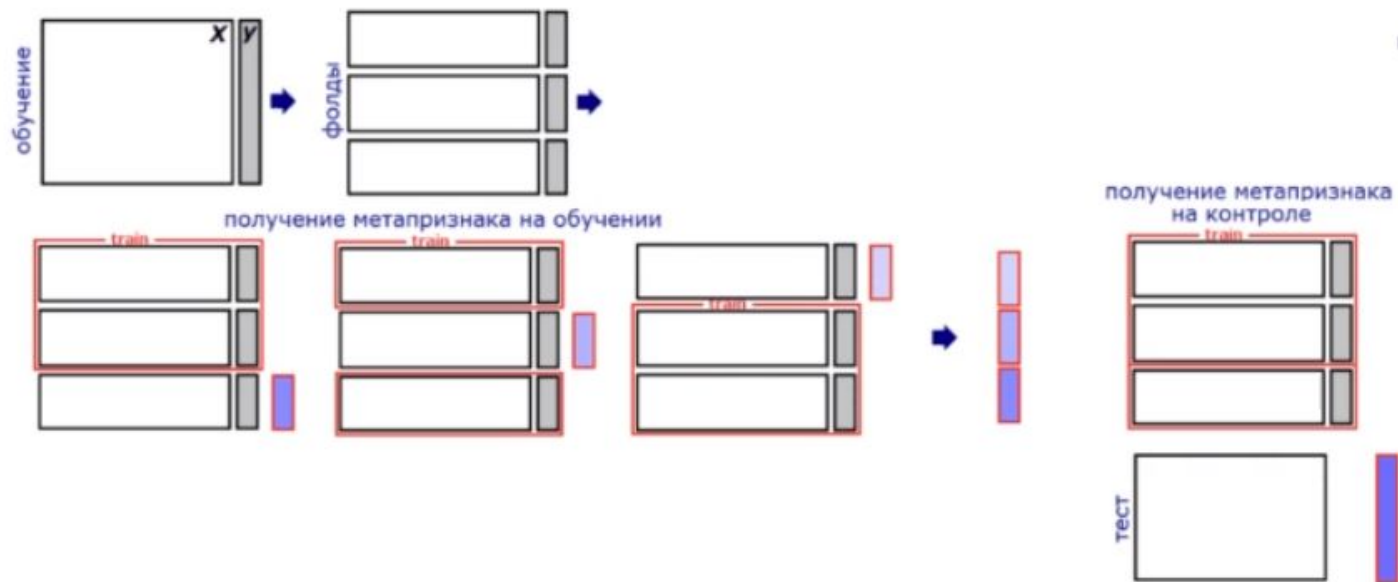
# Случайный лес



# Блендинг



# Стекинг



# Бустинг



train a weak model  
and aggregate it to  
the ensemble model



update the training dataset  
(values or weights) based on the  
current ensemble model results

