

# Юстина Иванова

Программист, data scientist

Статистика в python. Кейс-стади №1.  
Датасеты: faulty steel plates,  
Titanic, Boston houses.

Спикер



**Юстина Иванова,**  
Data scientist по  
Компьютерному зрению  
в компании ОЦРВ,  
Выпускница МГТУ им. Баумана  
Магистр по Artificial Intelligence  
В University of Southampton

# Классификационный анализ

Среди самых популярных задач Т в машинном обучении:

**классификация** – отнесение объекта к одной из категорий на основании его признаков

**регрессия** – прогнозирование количественного признака объекта на основании прочих его признаков

**кластеризация** – разбиение множества объектов на группы на основании признаков этих объектов так, чтобы внутри групп объекты были похожи между собой, а вне одной группы – менее похожи

**детекция аномалий** – поиск объектов, "сильно непохожих" на все остальные в выборке либо на какую-то группу объектов и много других, более специфичных.

Хороший обзор дан в главе "Machine Learning basics"

книги "Deep Learning" (Ian Goodfellow, Yoshua Bengio, Aaron Courville, 2016)

# Логистическая регрессия.

Задача логистической регрессии – определить вероятность принадлежности к классу.

Построена на основе линейной функции.

$$h(x) = \theta^T x$$

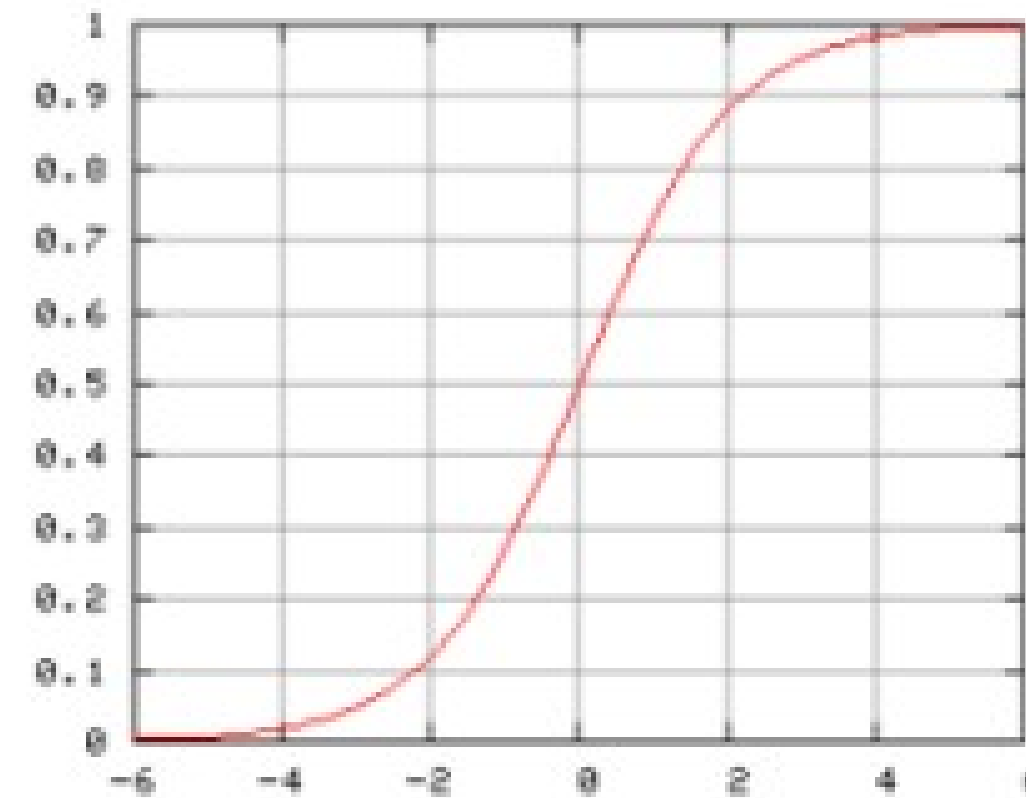
К линейной функции применяется функция активации:

$$h(x) = \sigma(\theta^T x)$$

Функция активации:

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

# Сигмоида.



$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Производная сигмоиды:

$$\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$$

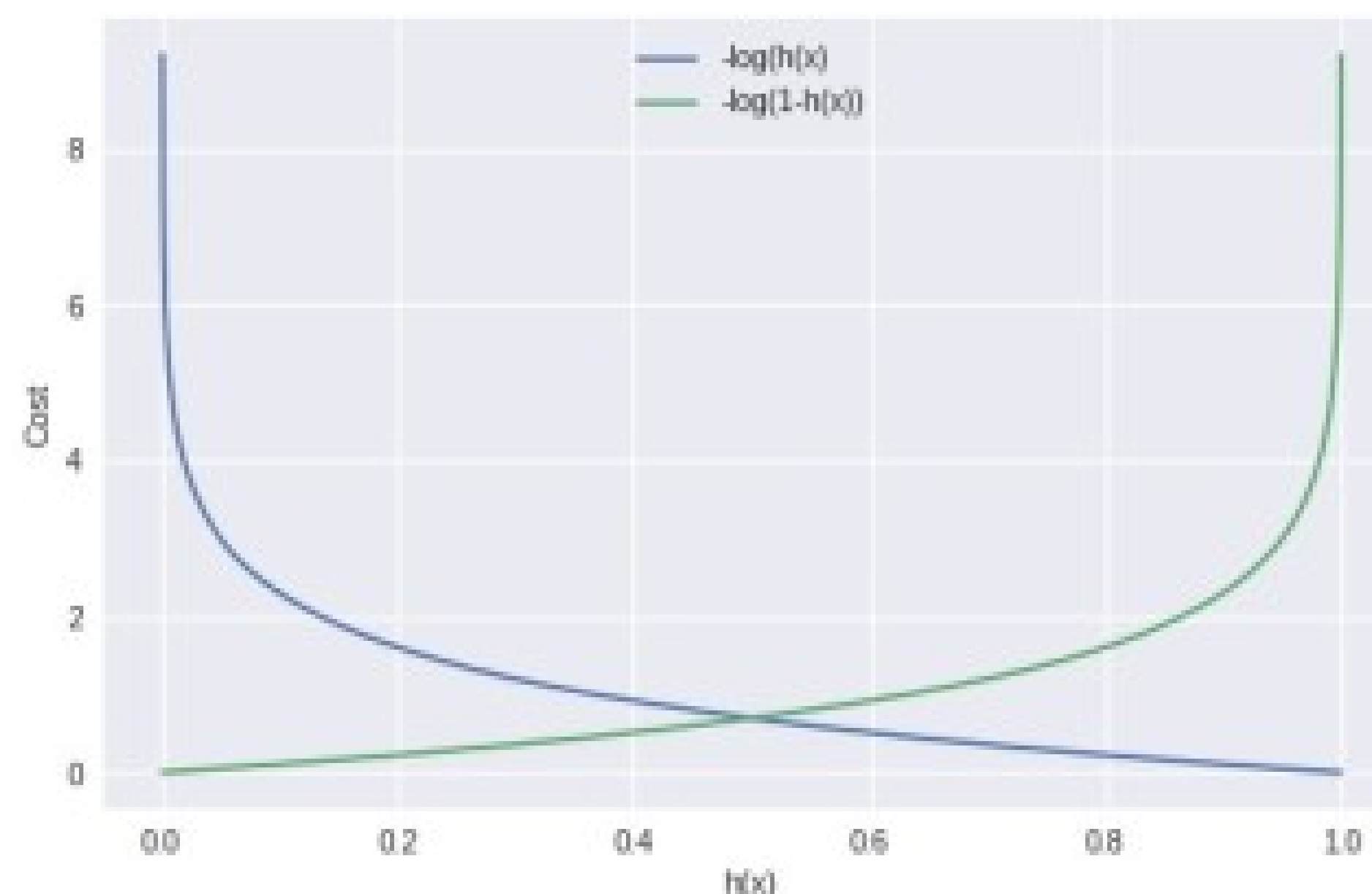
<https://ru.wikipedia.org/wiki/Сигмоида>

# Функция ошибки в логистической регрессии.

Модель ищет параметры, которые минимизируют функцию ошибки:

$$cost = \begin{cases} -\log(h(x)), & \text{if } y = 1 \\ -\log(1 - h(x)), & \text{if } y = 0 \end{cases}$$

Чем выше вероятность определения класса 1 при верном классе 0, тем выше стоимость ошибки.



## Функция ошибки в логистической регрессии.

Общий вид функции ошибки для модели:

$$\text{cost}(h(x), y) = -y \cdot \log(h(x)) - (1 - y)\log(1 - h(x))$$

Ошибка для всех данных датасета:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(h(x^i)) + (1 - y^i)\log(1 - h(x^i))]$$

Где  $m$  – количество элементов.

# Градиентный спуск.

Будем искать минимум функции относительно параметров:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h(x^i) - y^i) x_j^i$$

Для поиска минимума используется градиентный спуск.

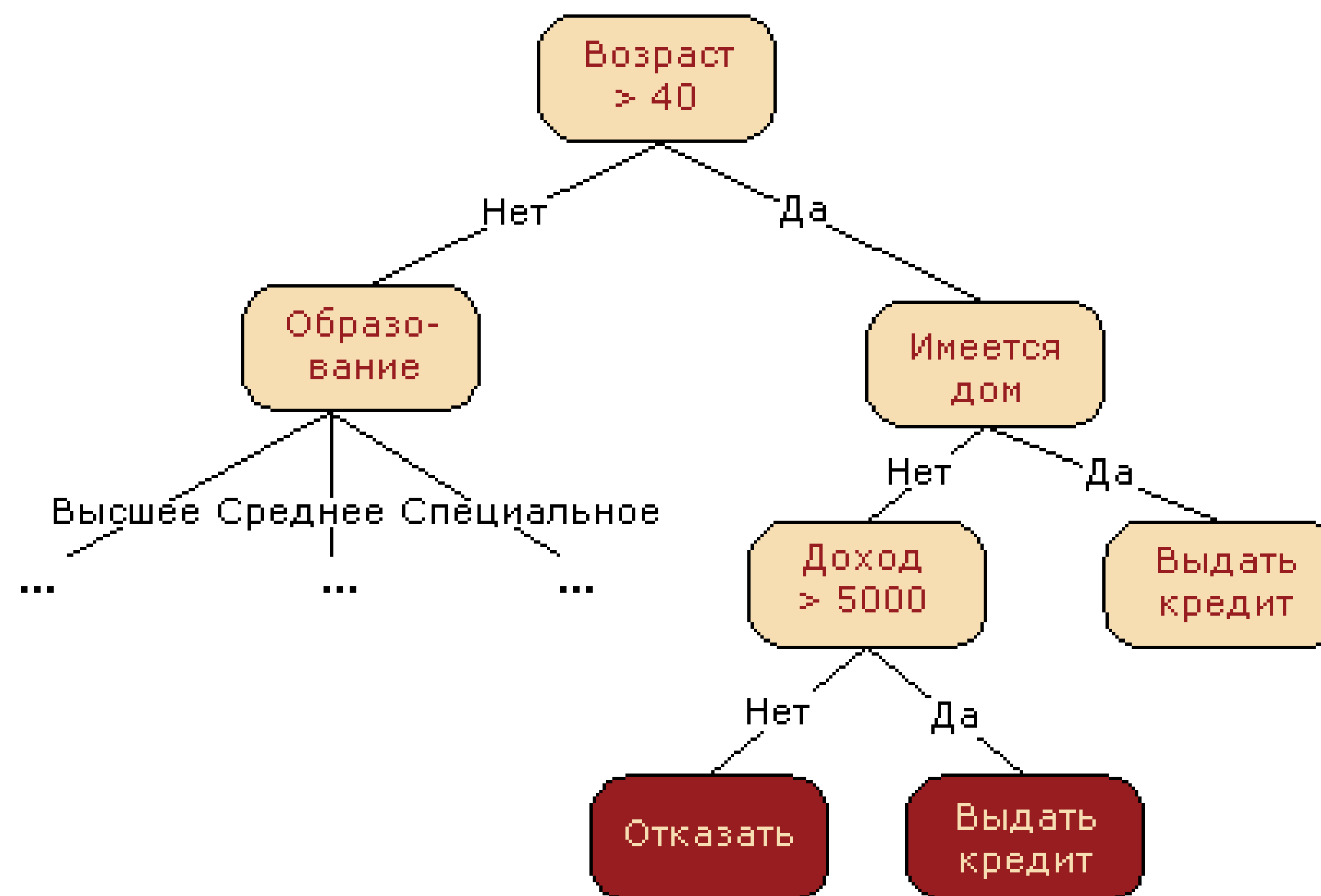
Это метод нахождения локального минимума (максимума) функции с помощью движения вдоль градиента.



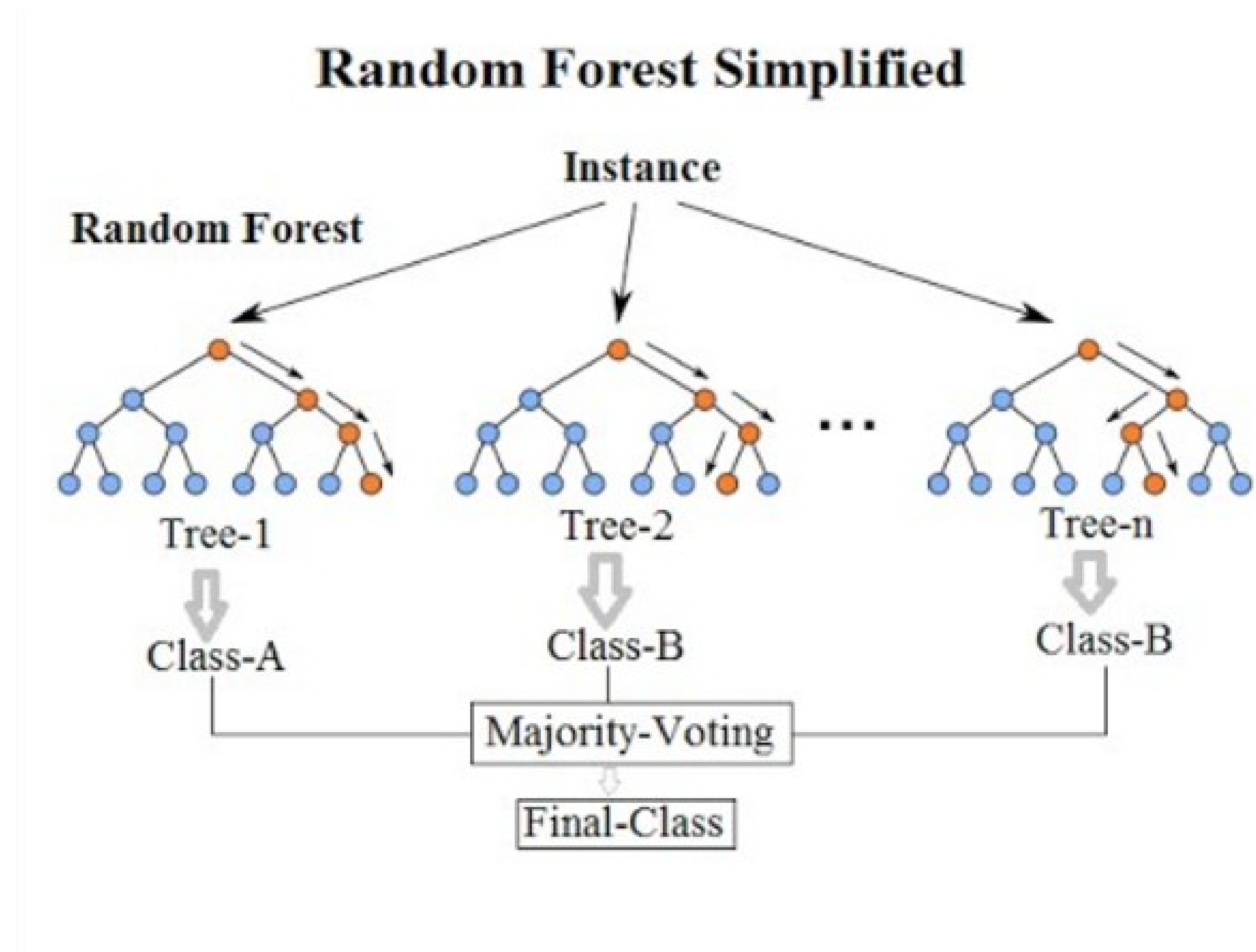
# Деревья решений.



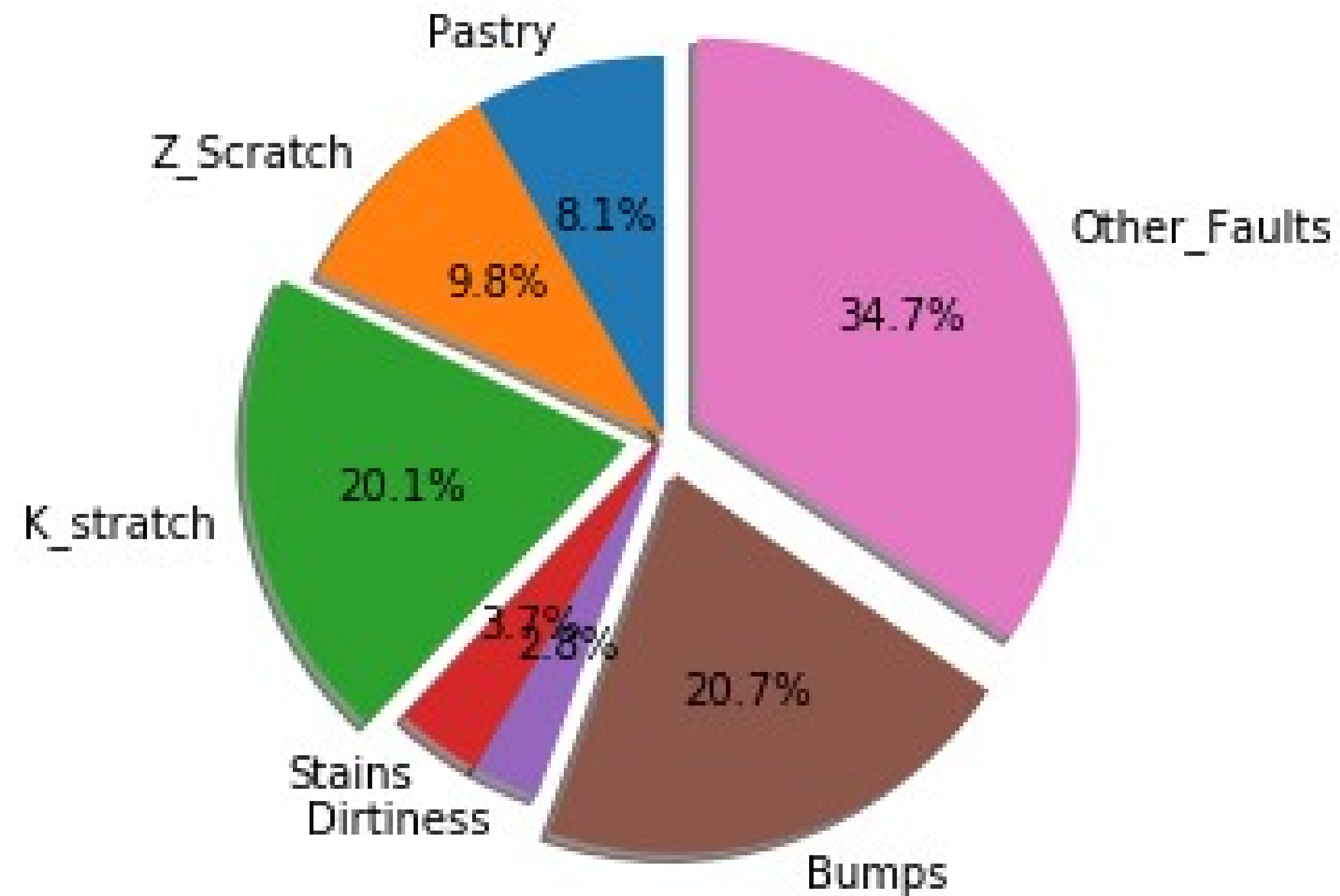
# Деревья решений.



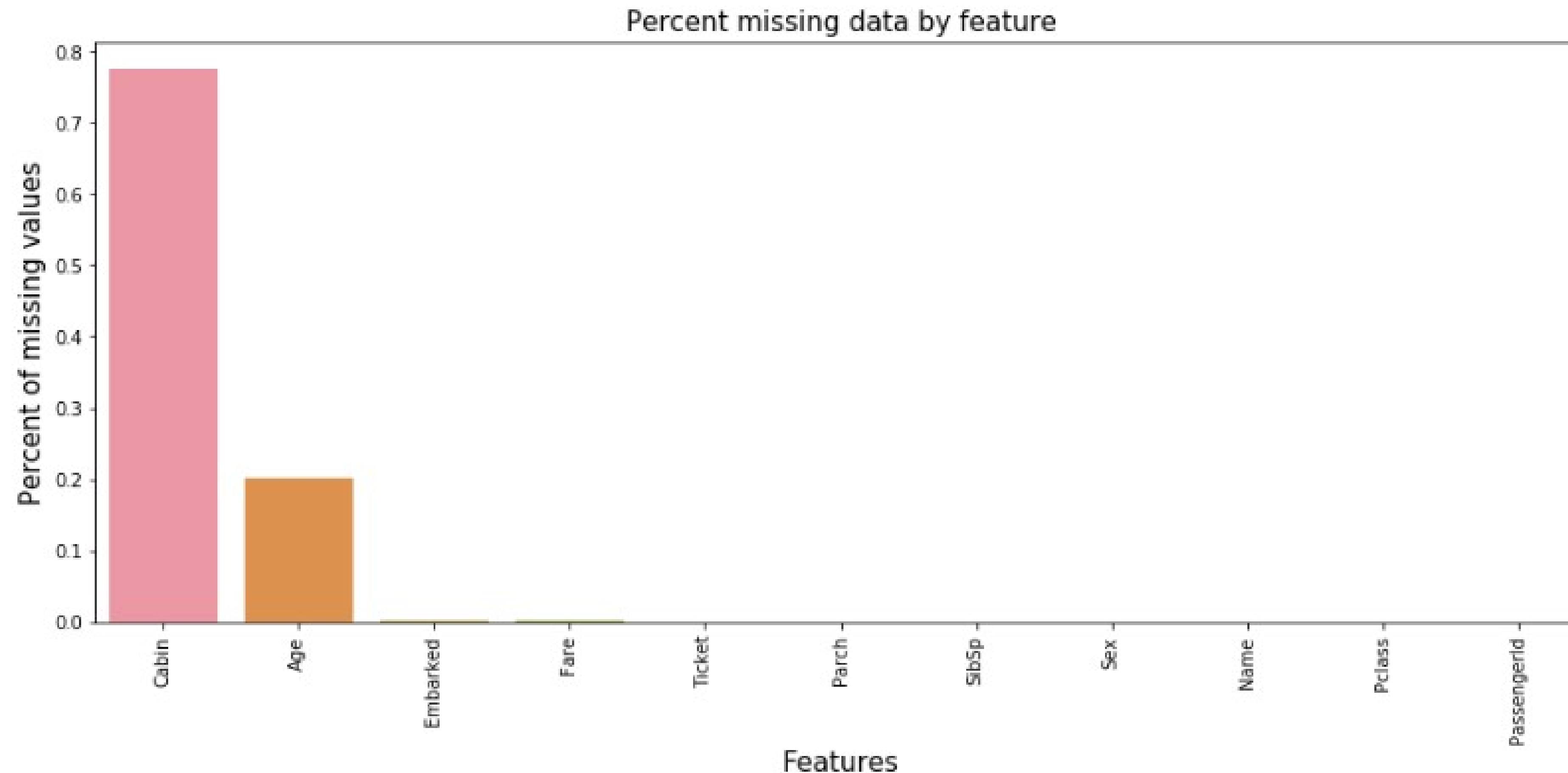
# Случайный лес.



# Проблема несбалансированности классов.



# Проблема нехватки информации для элементов.



# Методы заполнения недостающих данных.

Удаление элементов

Заполнение данными предыдущего или последующего элемента

Заполнение некой константой (выходящей за пределы интервала значений)

Заполнение средним значением или модой

Создать атрибут: отсутствующее значение

# Вопросы?

Контакты спикера:  
[yustiks@gmail.com](mailto:yustiks@gmail.com)