

Кейс стади 1



Задача классификации. Поиск причины дефекта в металле. Тесты на согласие: поиск распределения. Проблема мультиколлинеарности. Проблема несбалансированности классов. Пример решения задачи классификации с помощью RandomForest. Метрики классификации: precision, recall, F1. Принцип минимальных компонент, PCA. Кросс-валидация. ROC-кривая.

Юстина Иванова

Специалист по Анализу Данных



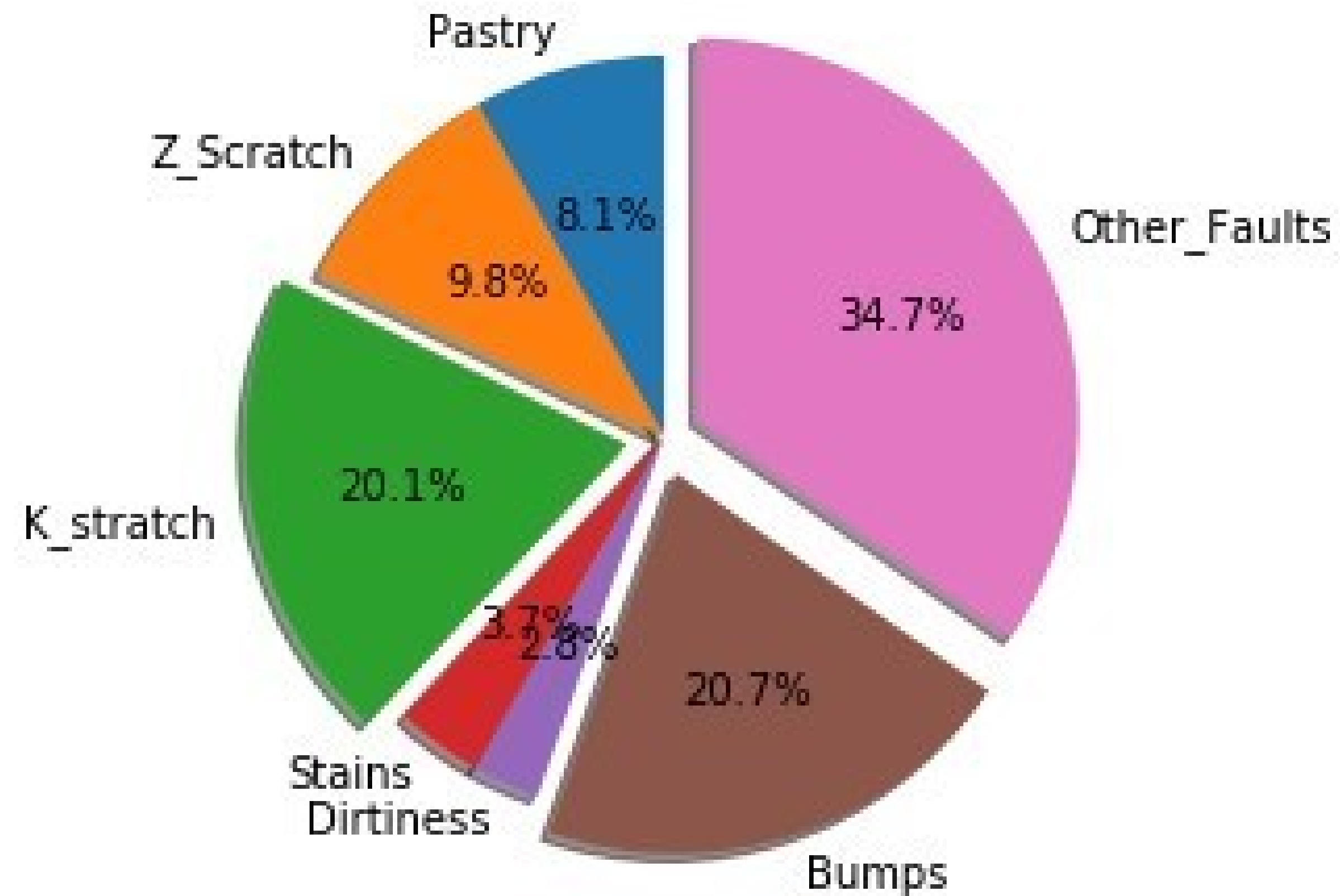
Инженер-программист МГТУ им. Баумана

Master of Science in Artificial Intelligence
University of Southampton

Специалист по анализу данных
в компании ОЦРВ.

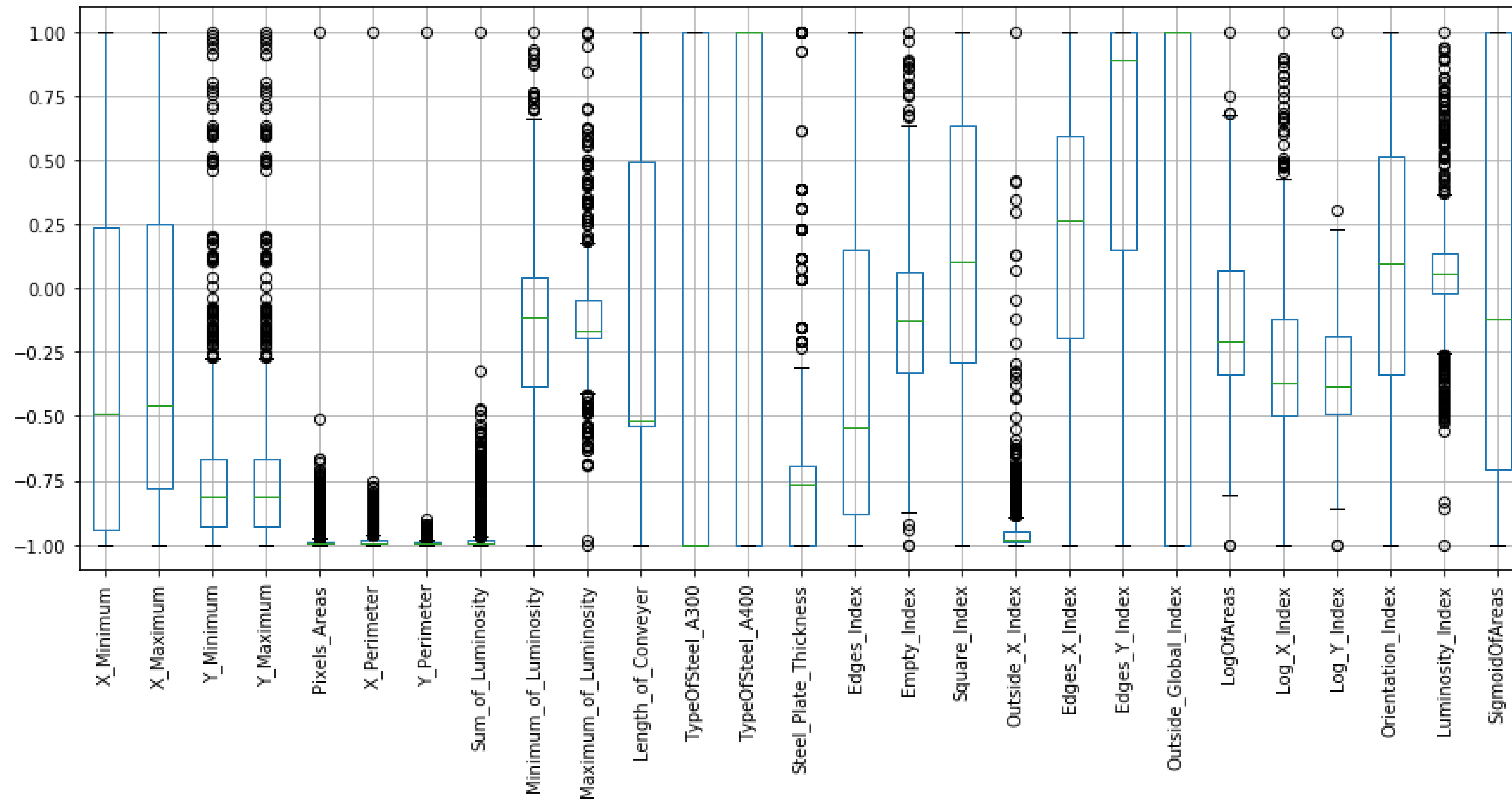
Юстина Иванова
студент-аспирант
University of Bolzano

Датасет Faulty Steel Plates.



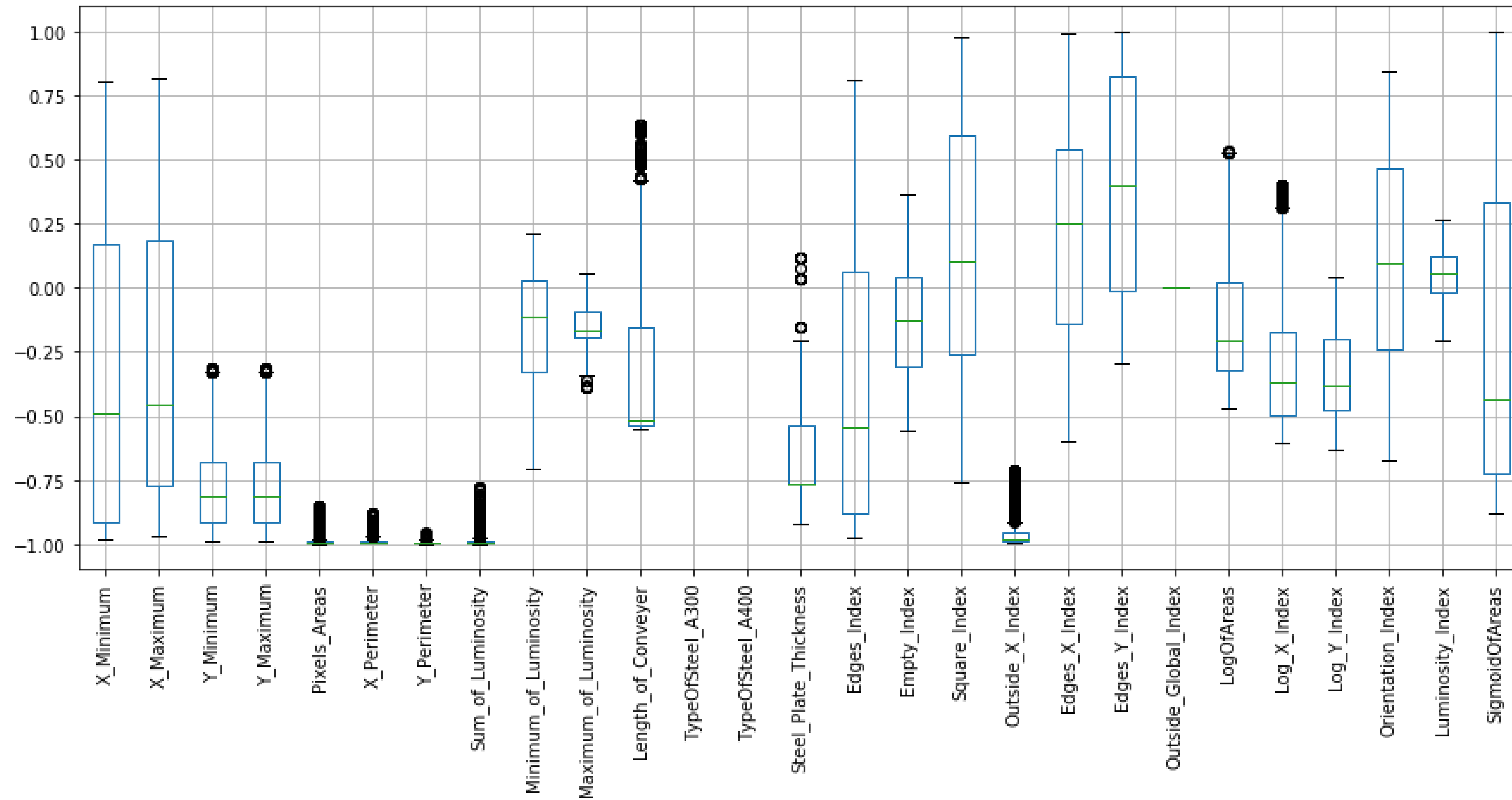


Boxplot → выбросы.

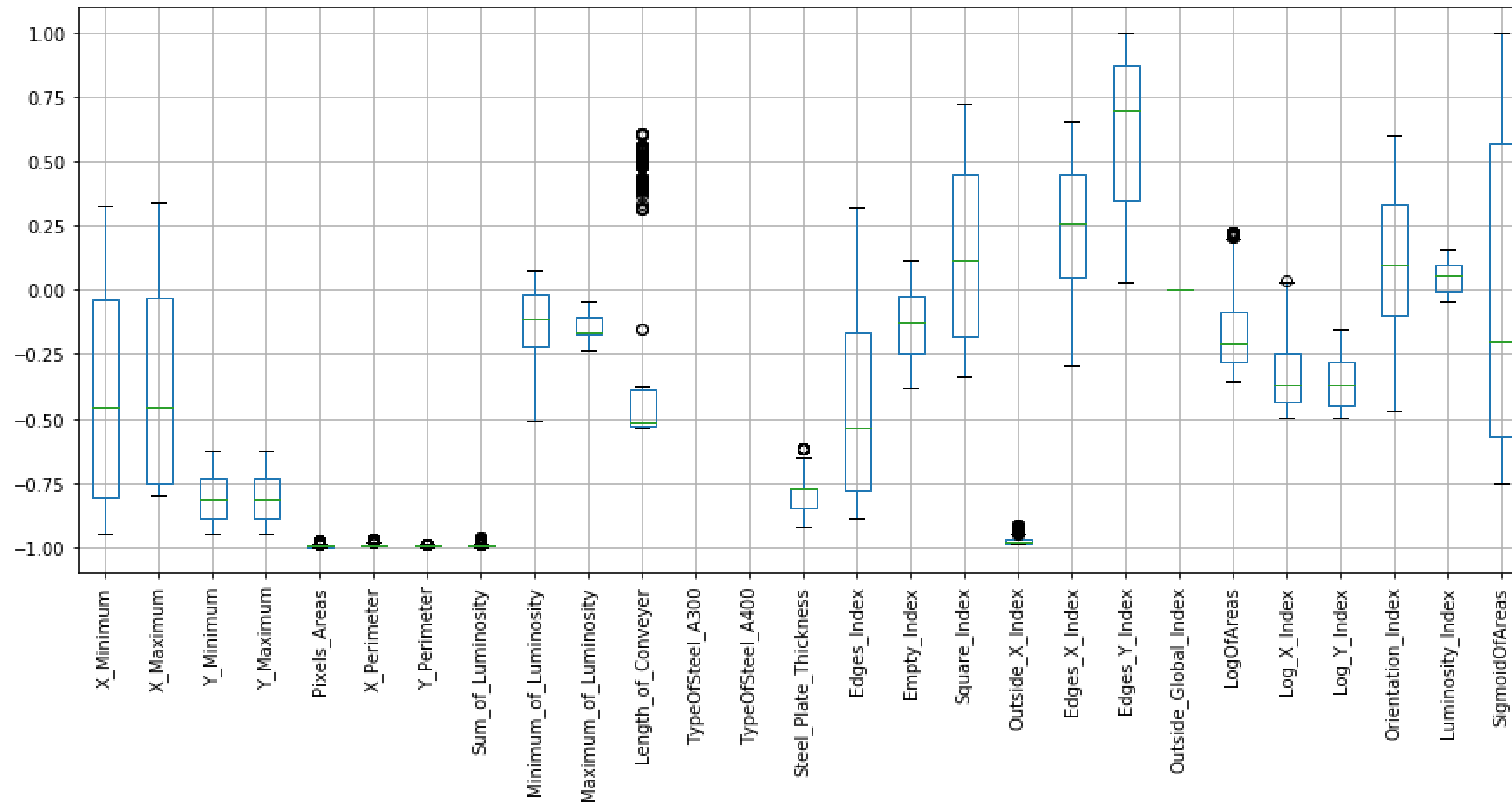


<https://pandas.pydata.org/pandas-docs/version/0.23/generated/pandas.DataFrame.boxplot.html>

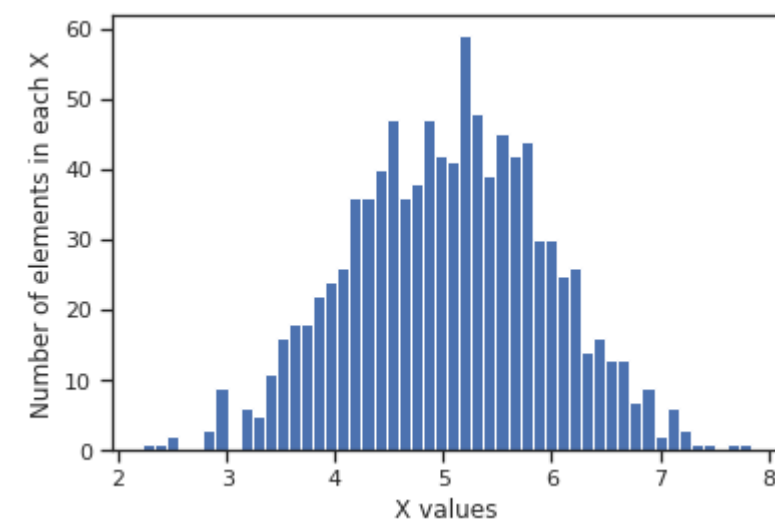
Удаление элементов вне интерквартильного интервала..



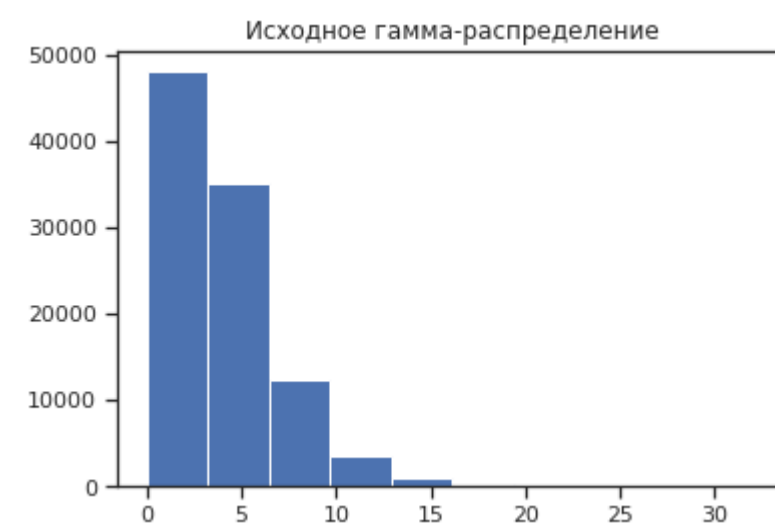
Удаление элементов вне квантилей 20% и 80%.



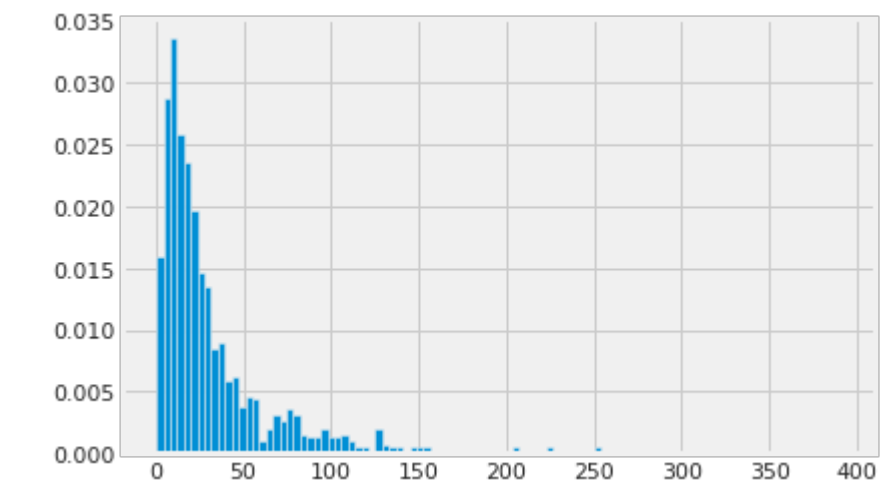
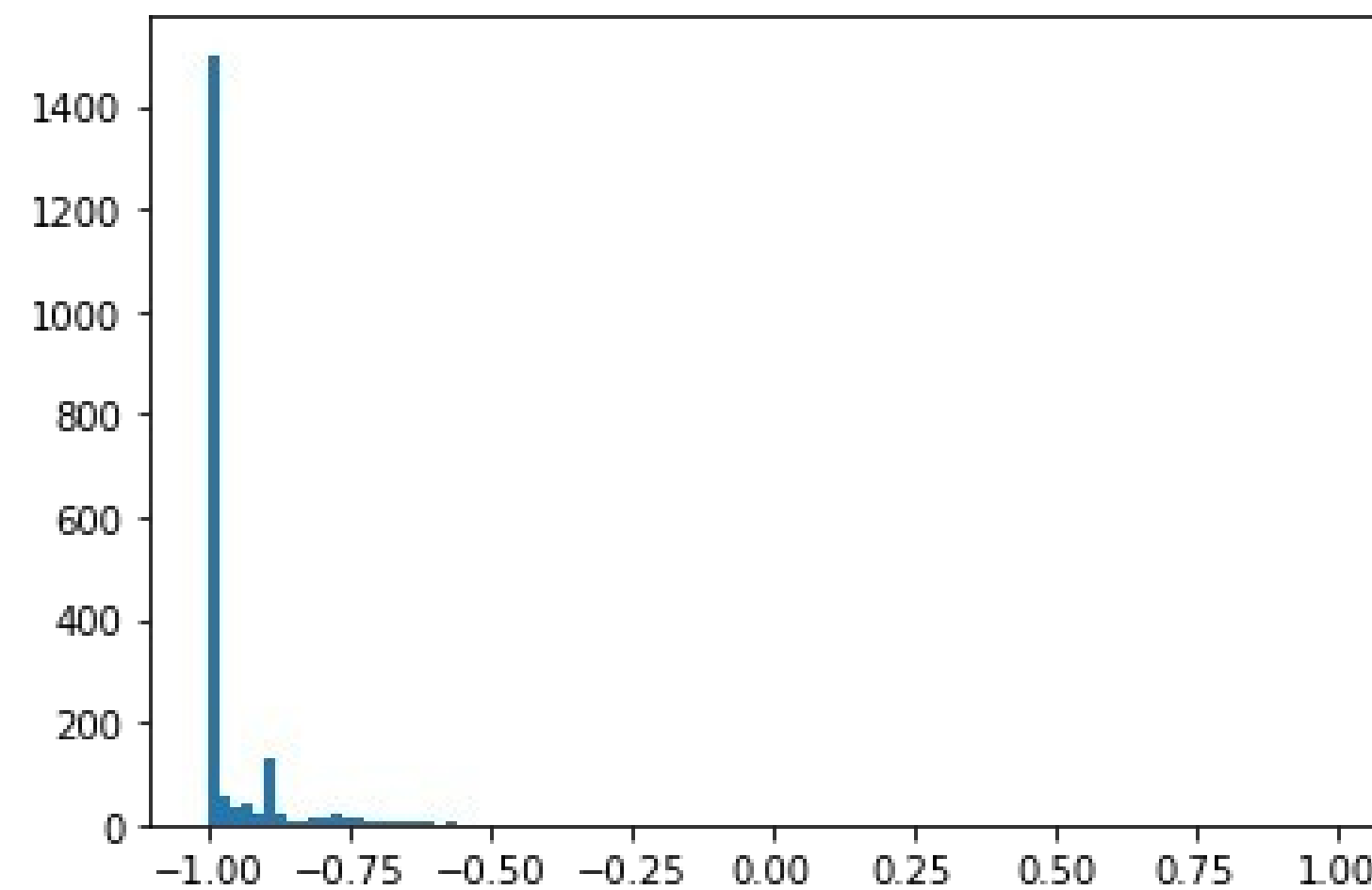
Тесты на согласие: какое это распределение?



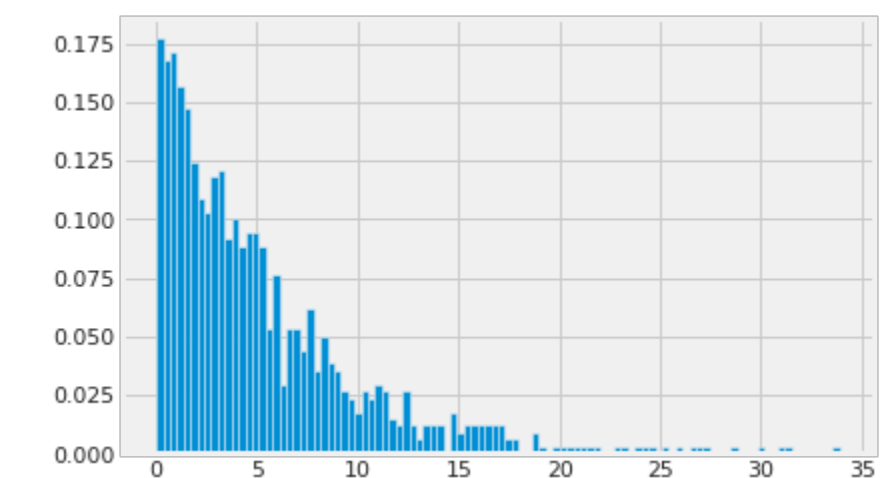
нормальное



гамма

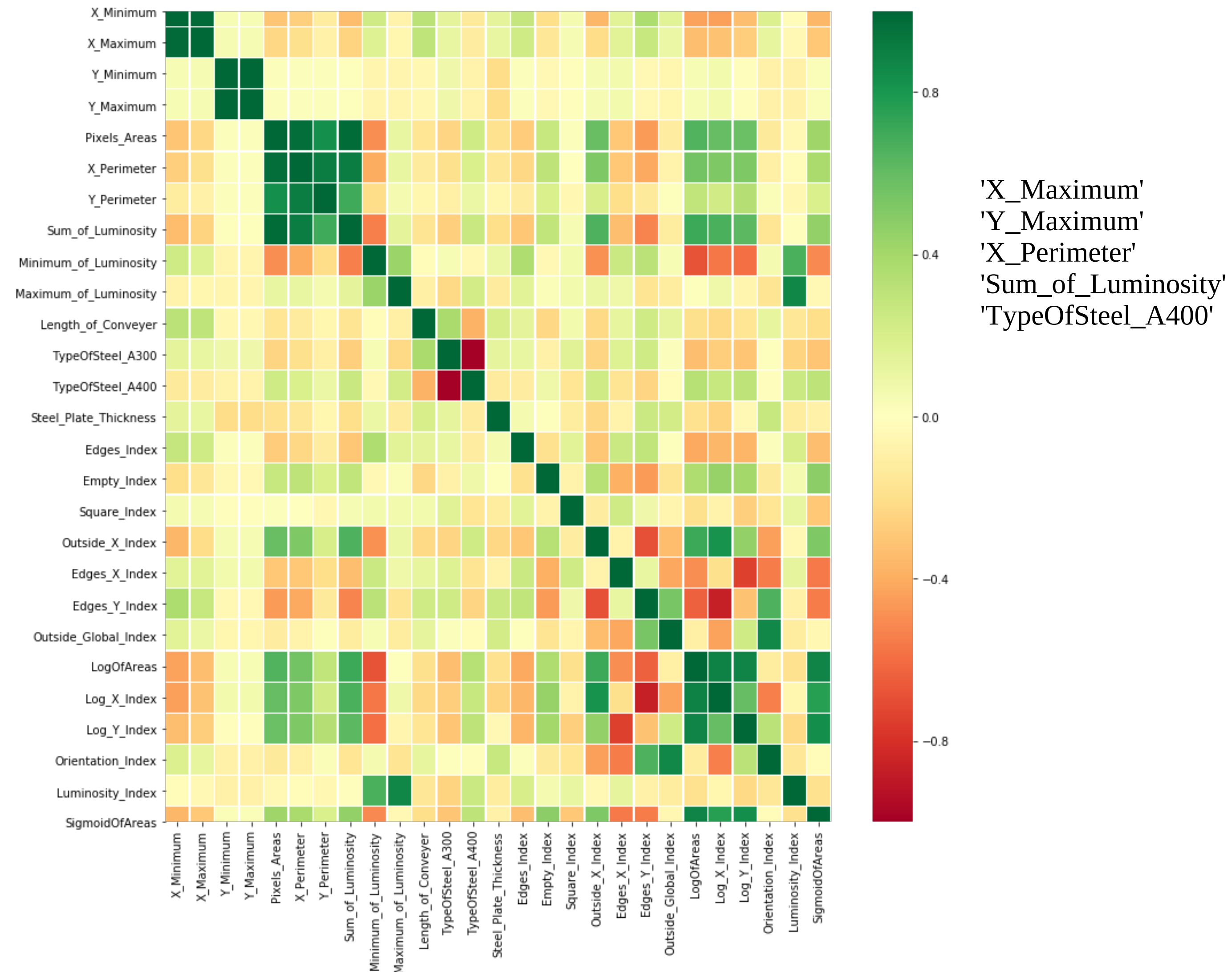


экспоненциальное



экспоненциальное

Удаление мультиколлинеарности

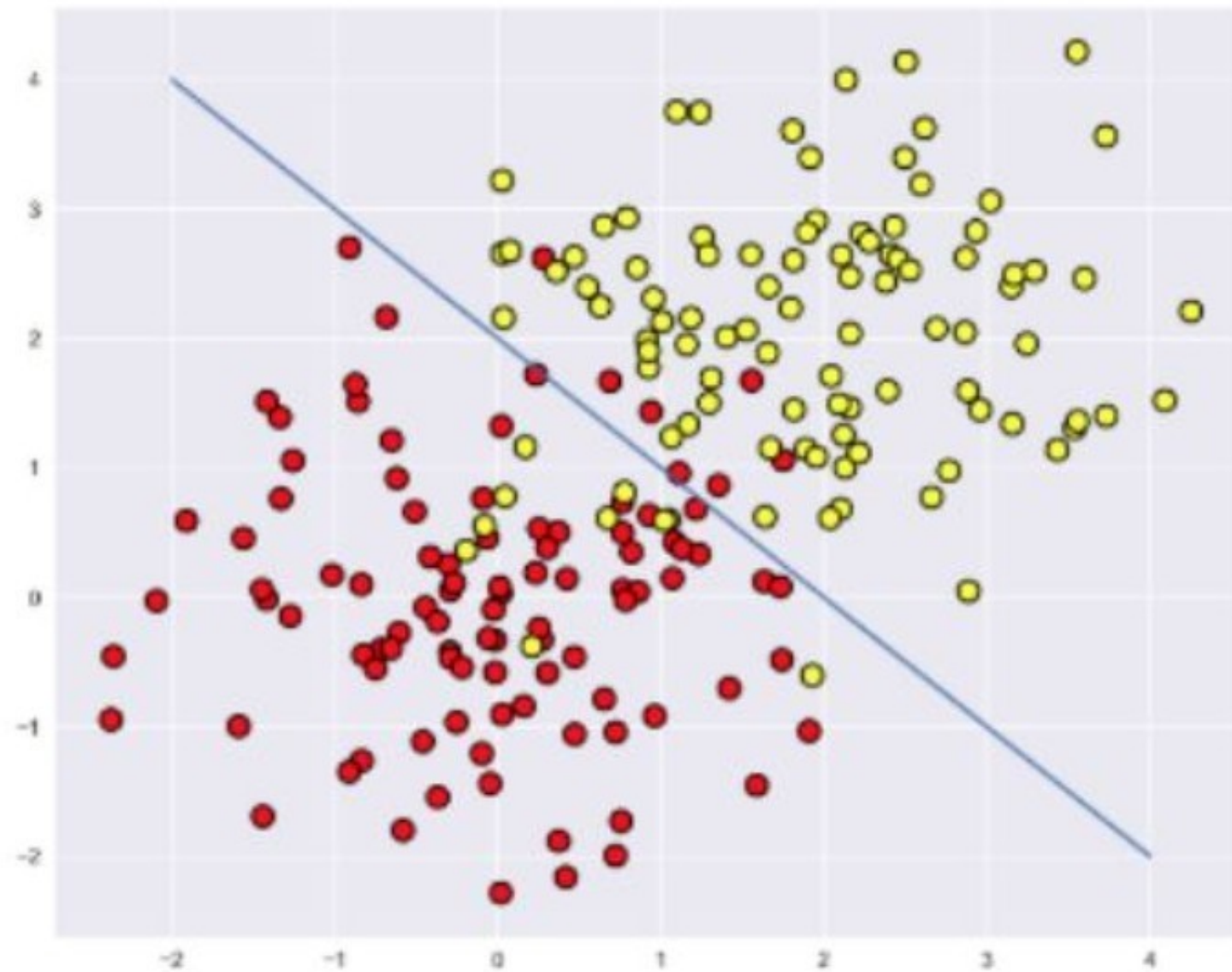


Классическое Обучение



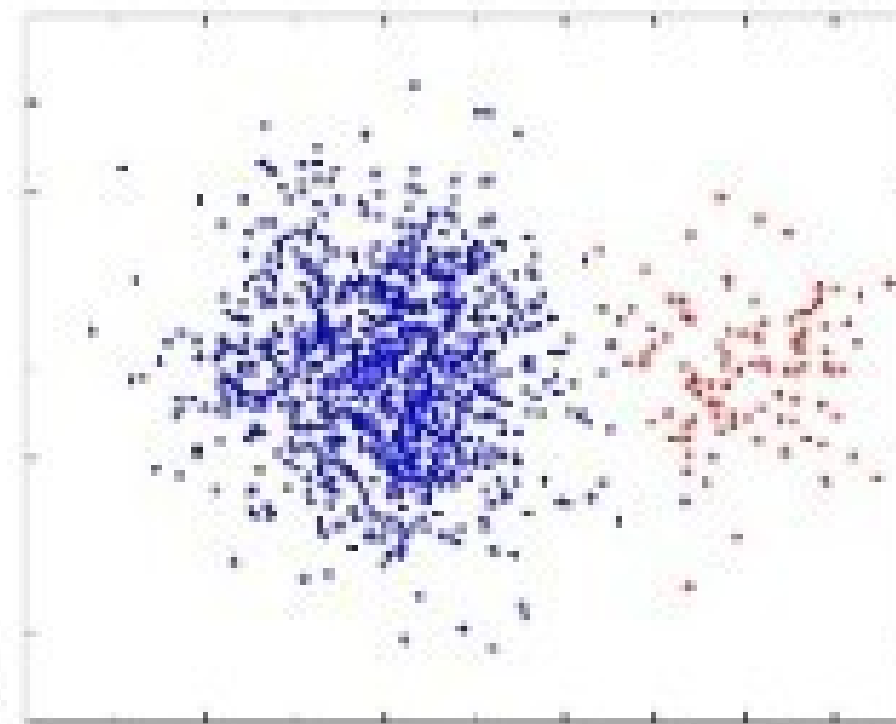
Классификация

Множество допустимых ответов конечно. Их называют метками классов (class label). Класс — это множество всех объектов с данным значением метки.

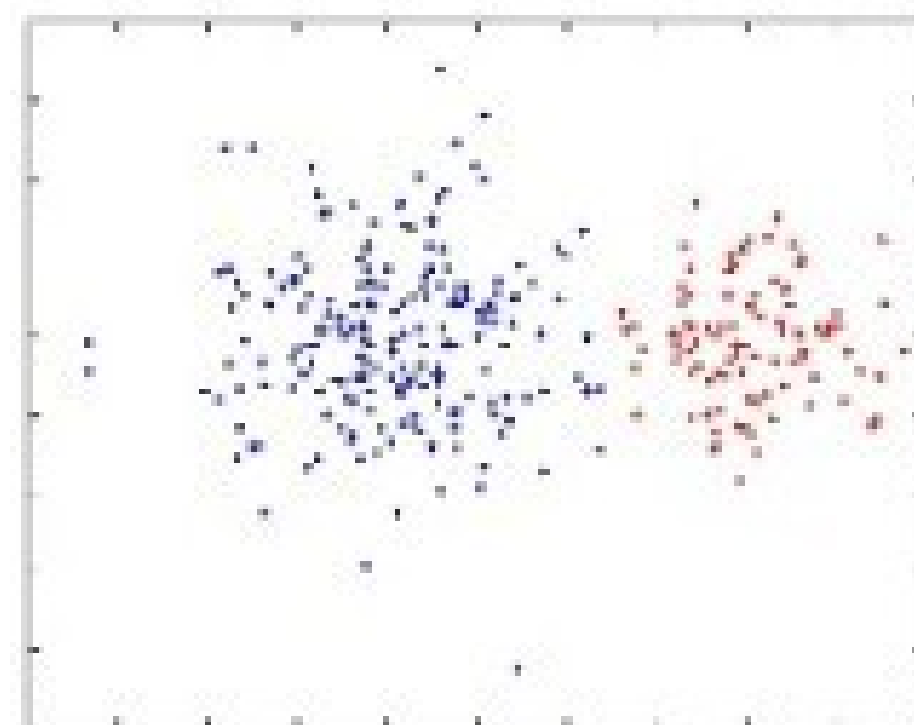


Проблема несбалансированности классов.

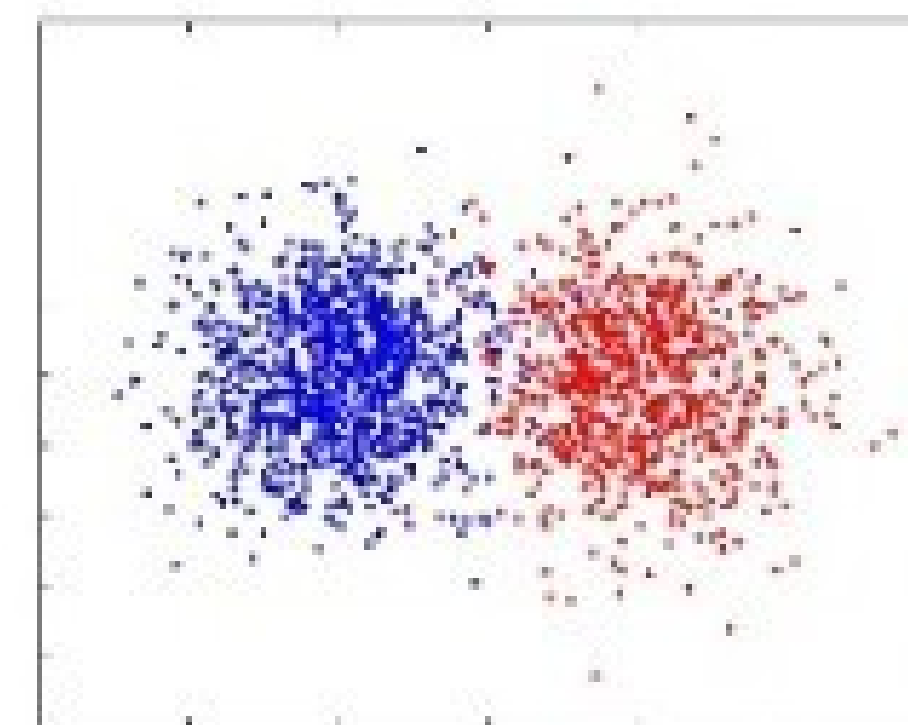
Sampling: Rebalancing the dataset



Under-sampling



Over-sampling



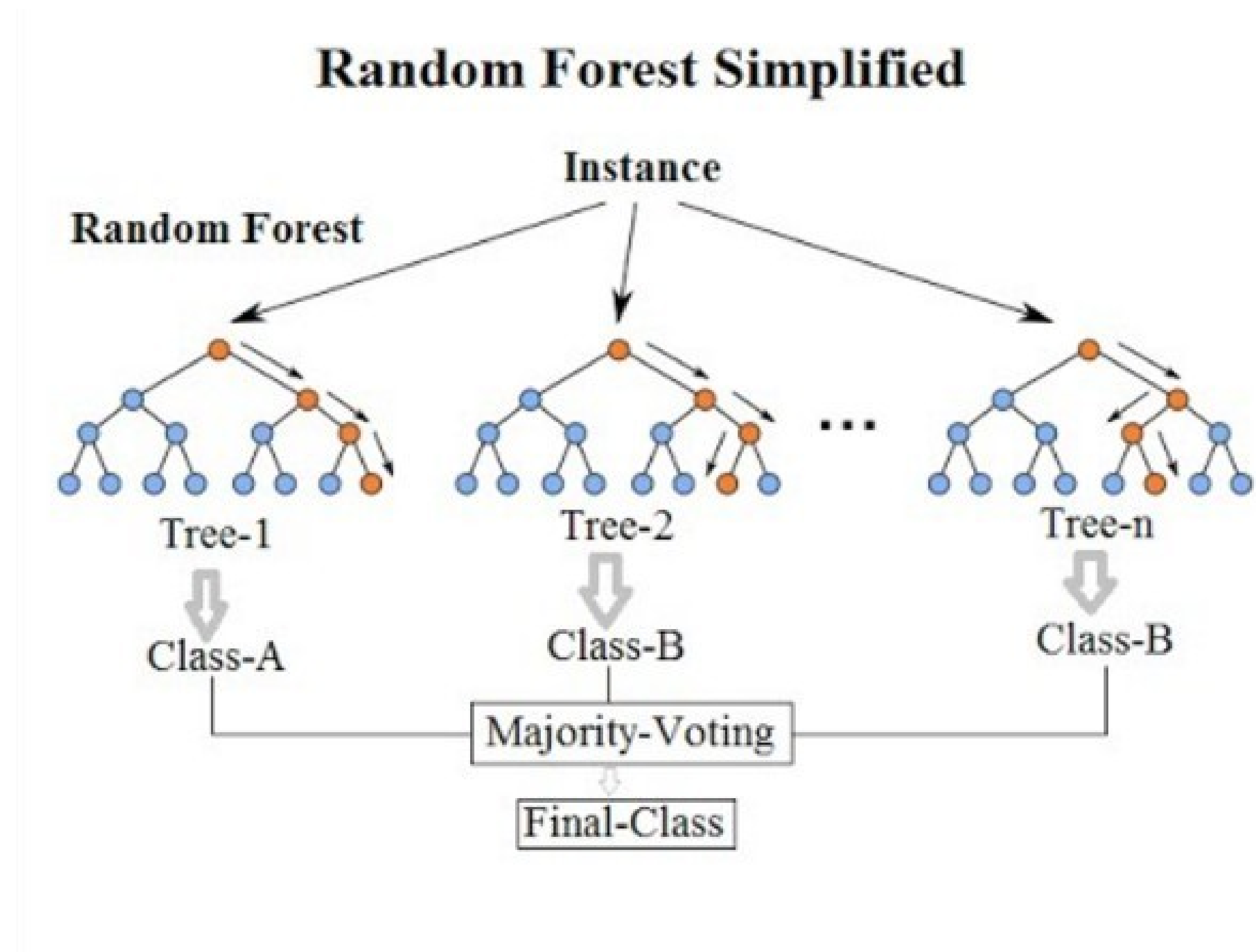
Дерево решений.

Давать ли кредит?

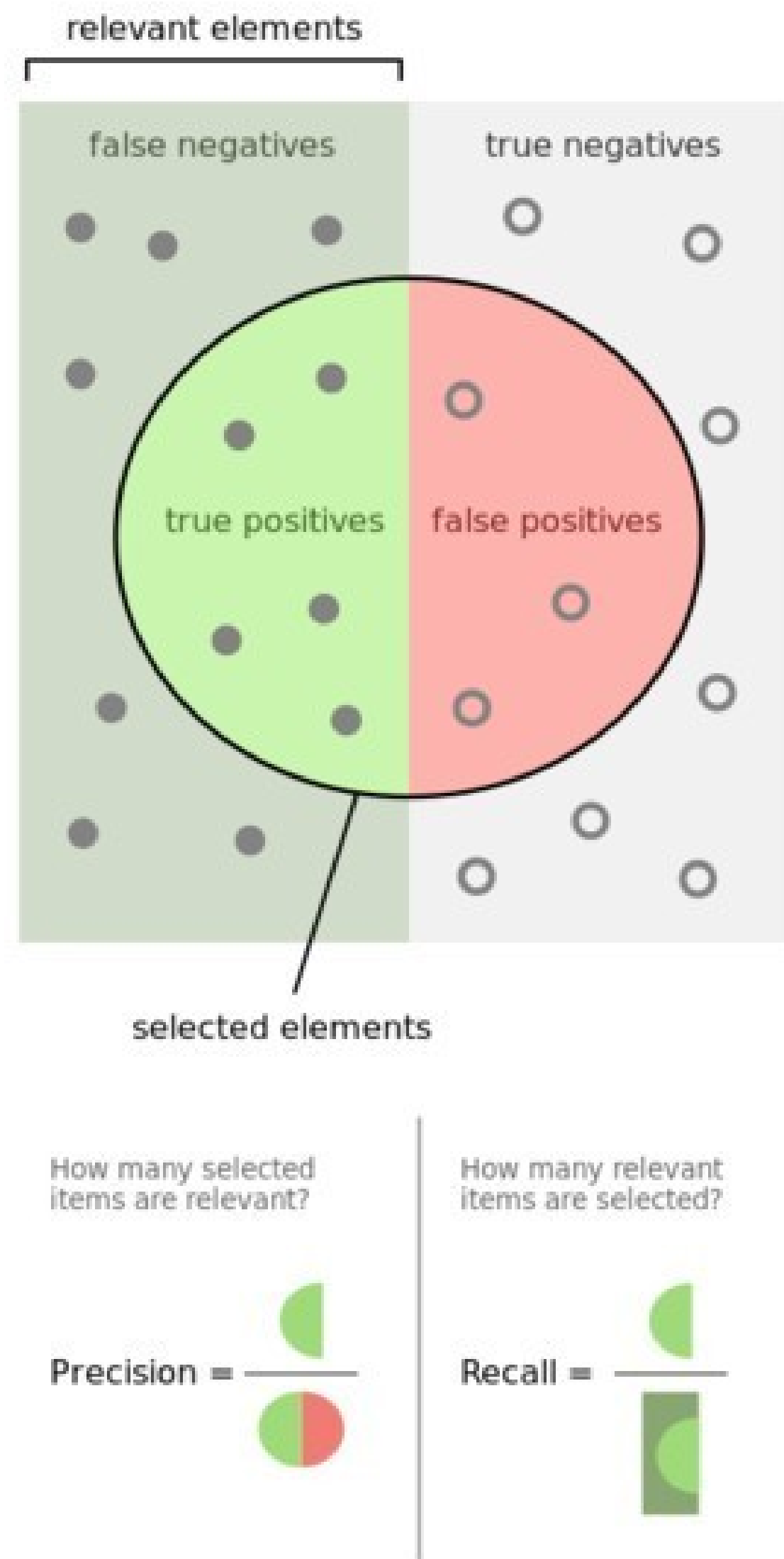


Дерево Решений

Случайный лес.



Метрики классификации



Precision

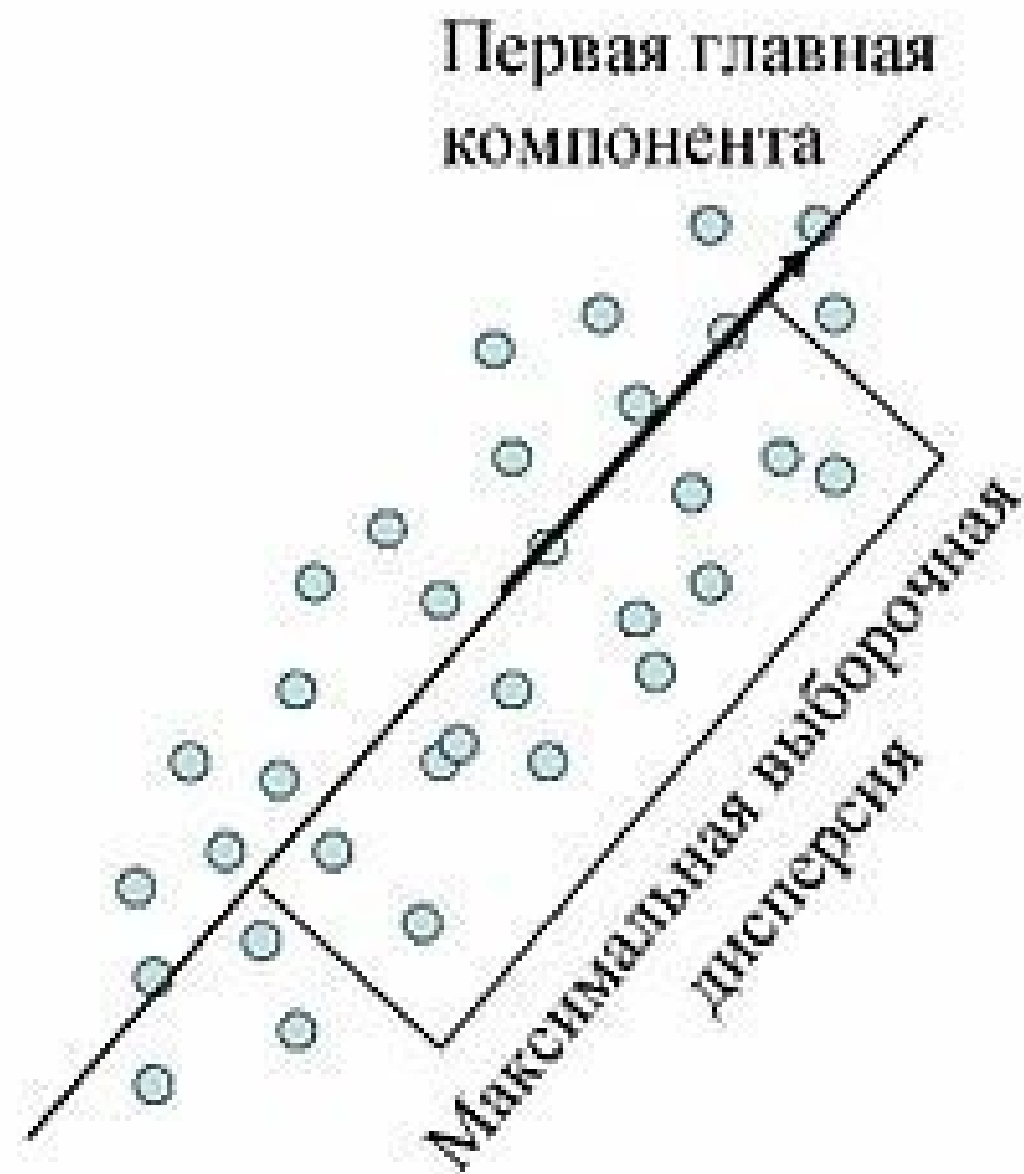
Recall

F1-мера

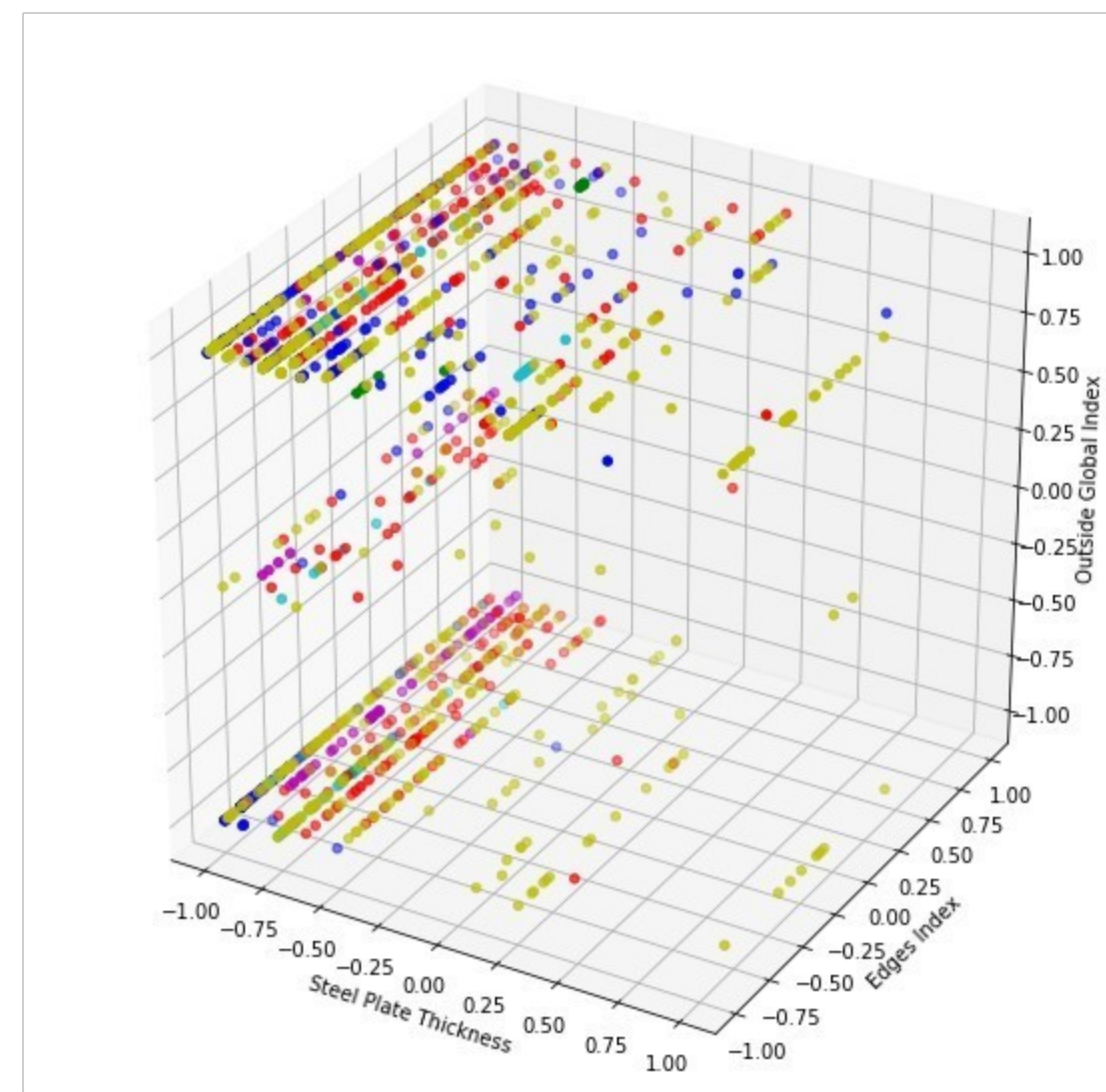
$$F_{\beta} = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}$$

Принцип минимальных компонент.

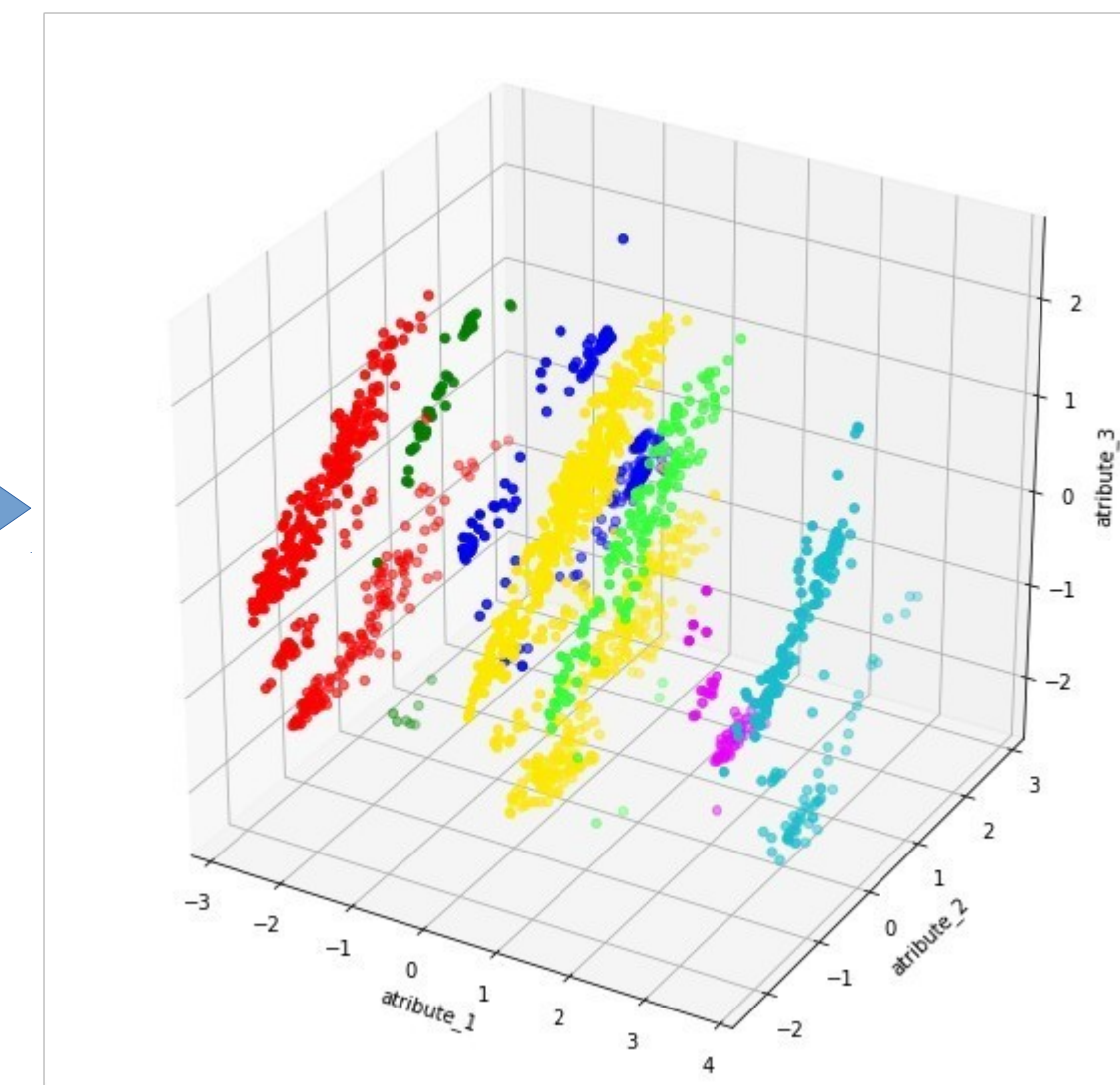
Поиск ортогональных проекций с наибольшим рассеянием



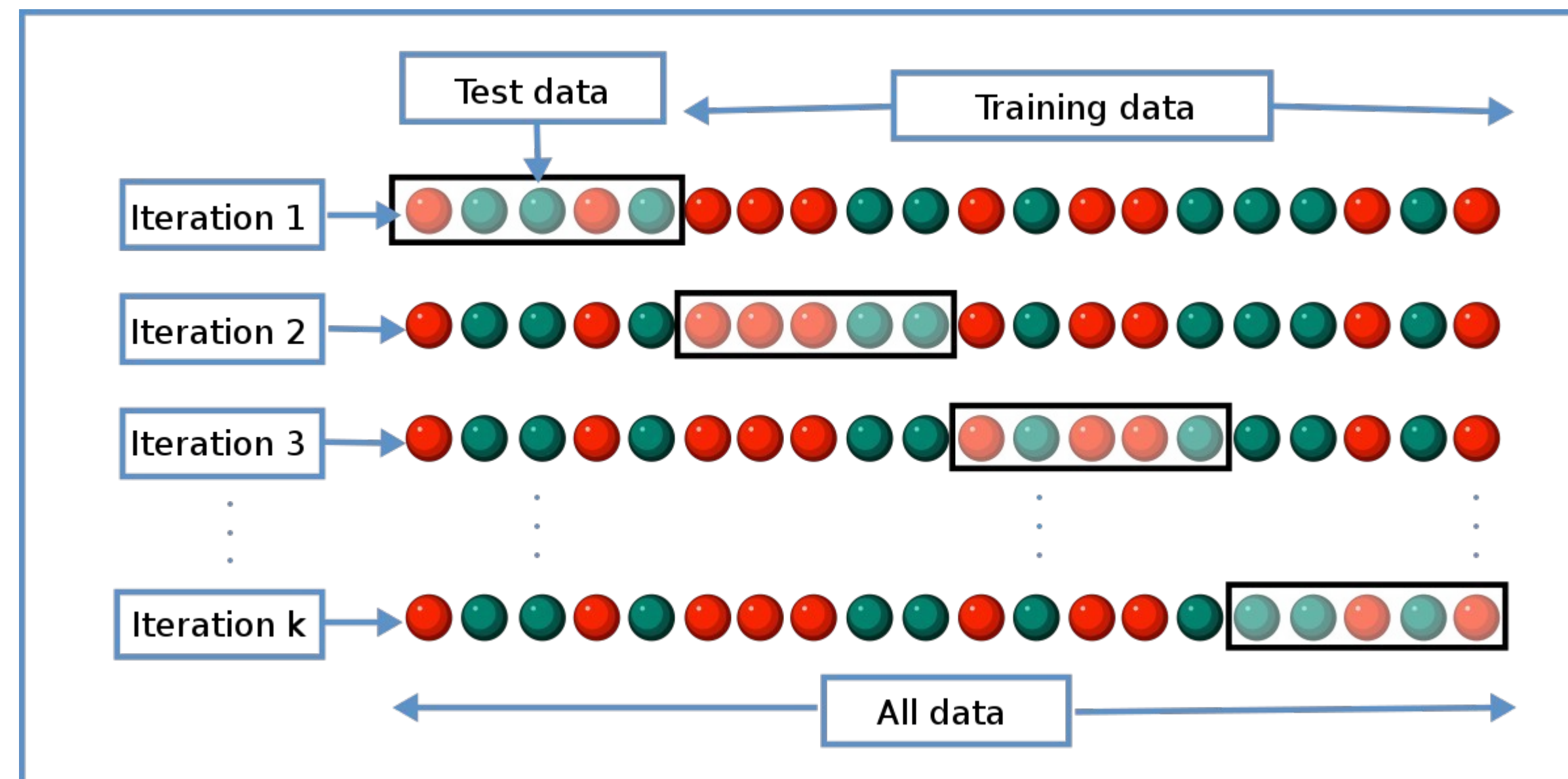
Было



Стало

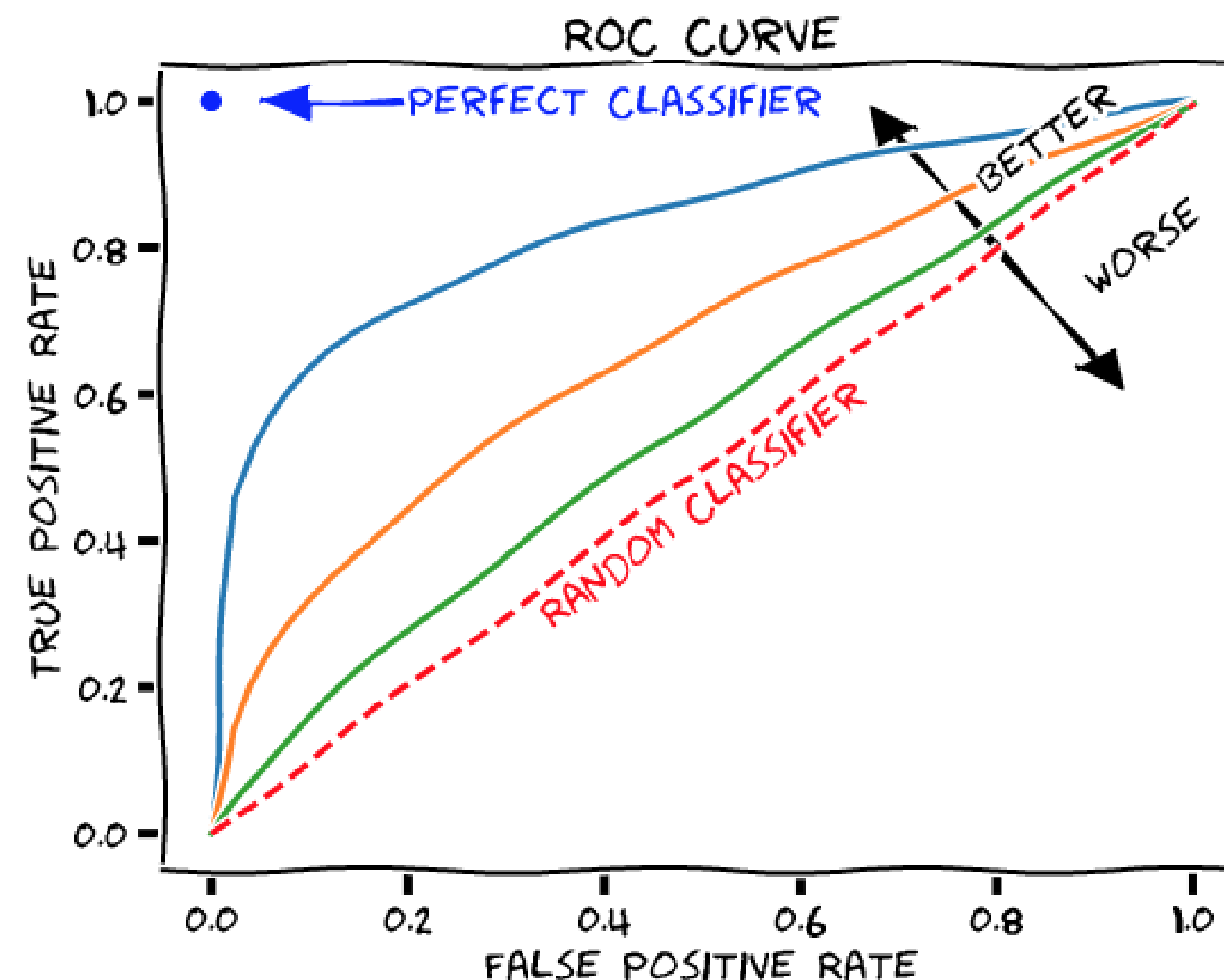


Кросс-валидация



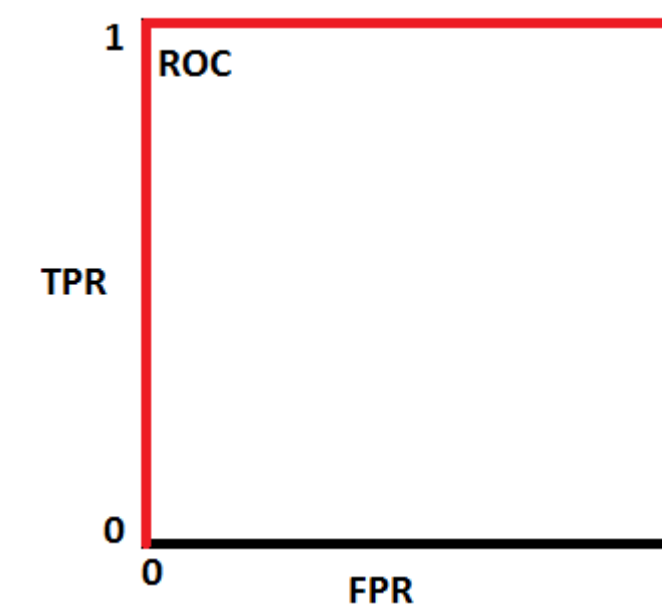
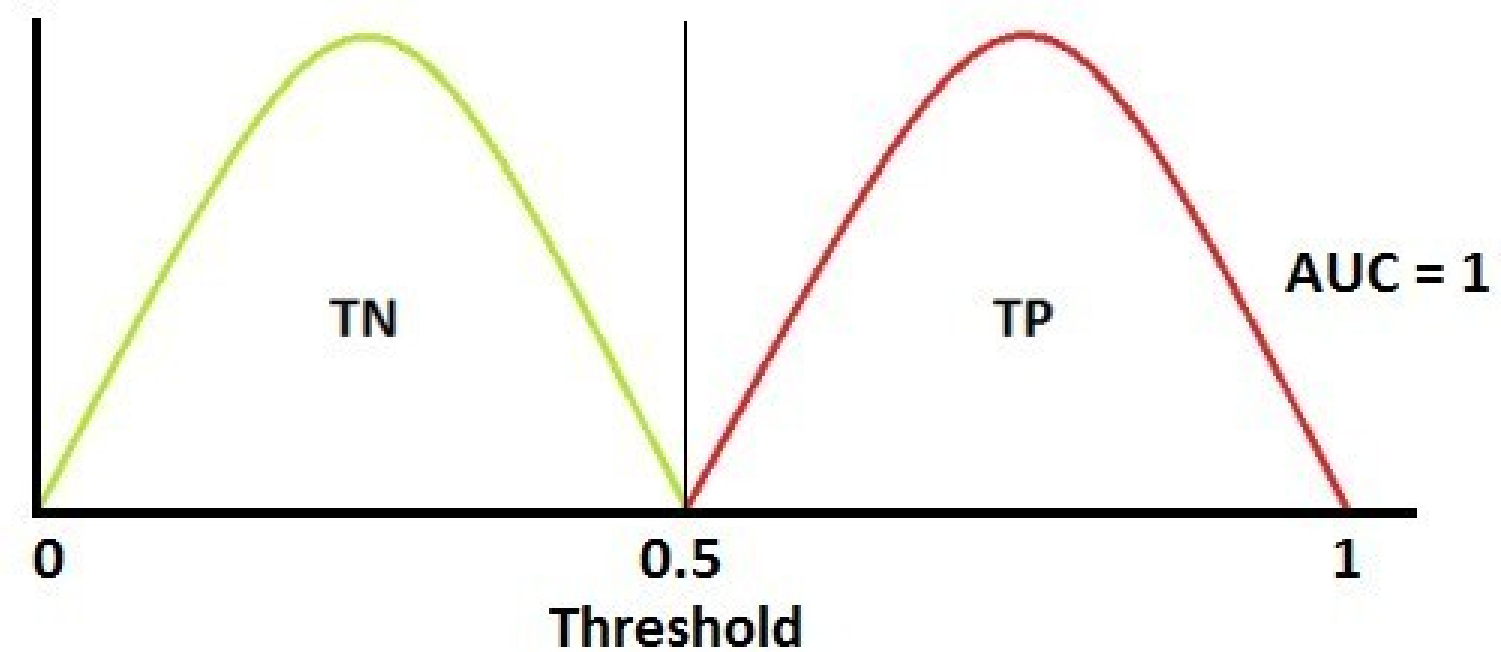
Оцениваем модель на нескольких тестовых данных

Метрики классификации: ROC-кривая

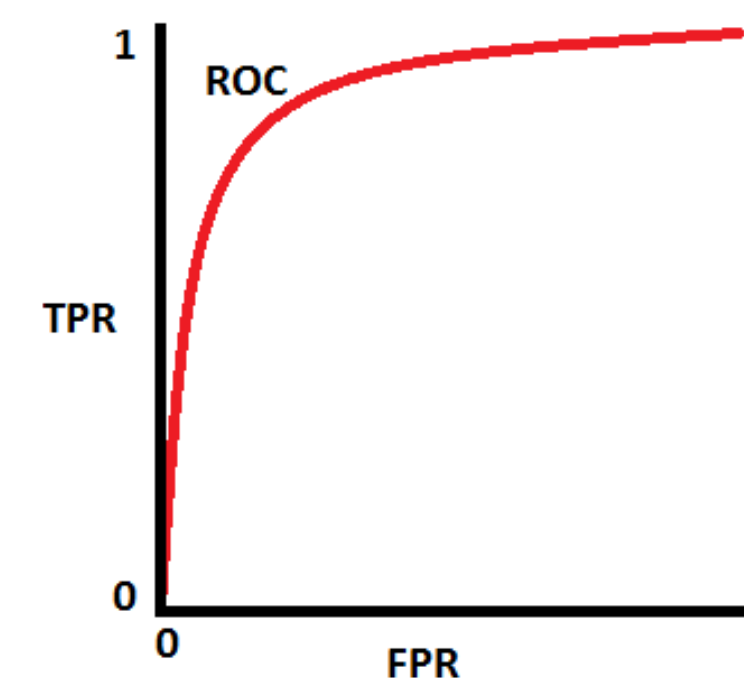
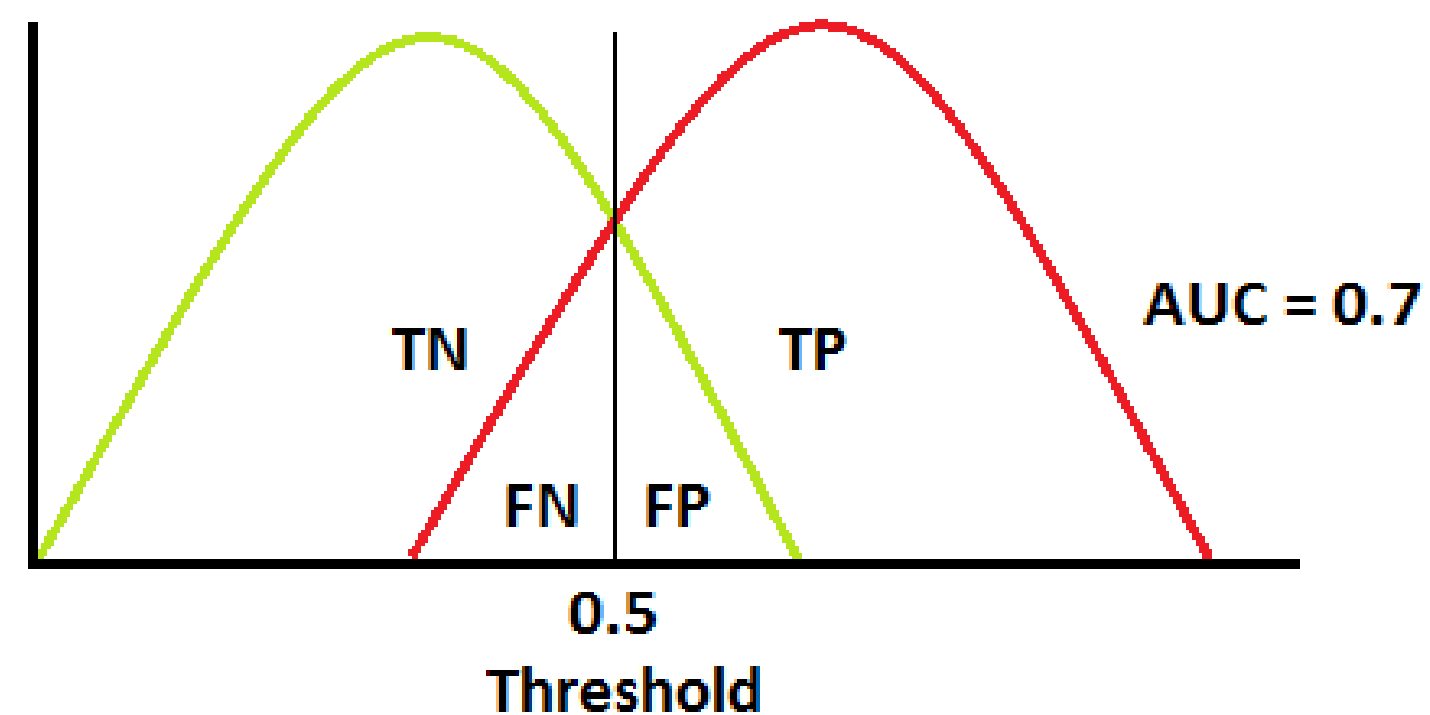


Позволяет определить порог,
при котором мы будем отделять один класс от другого

ROC-кривая



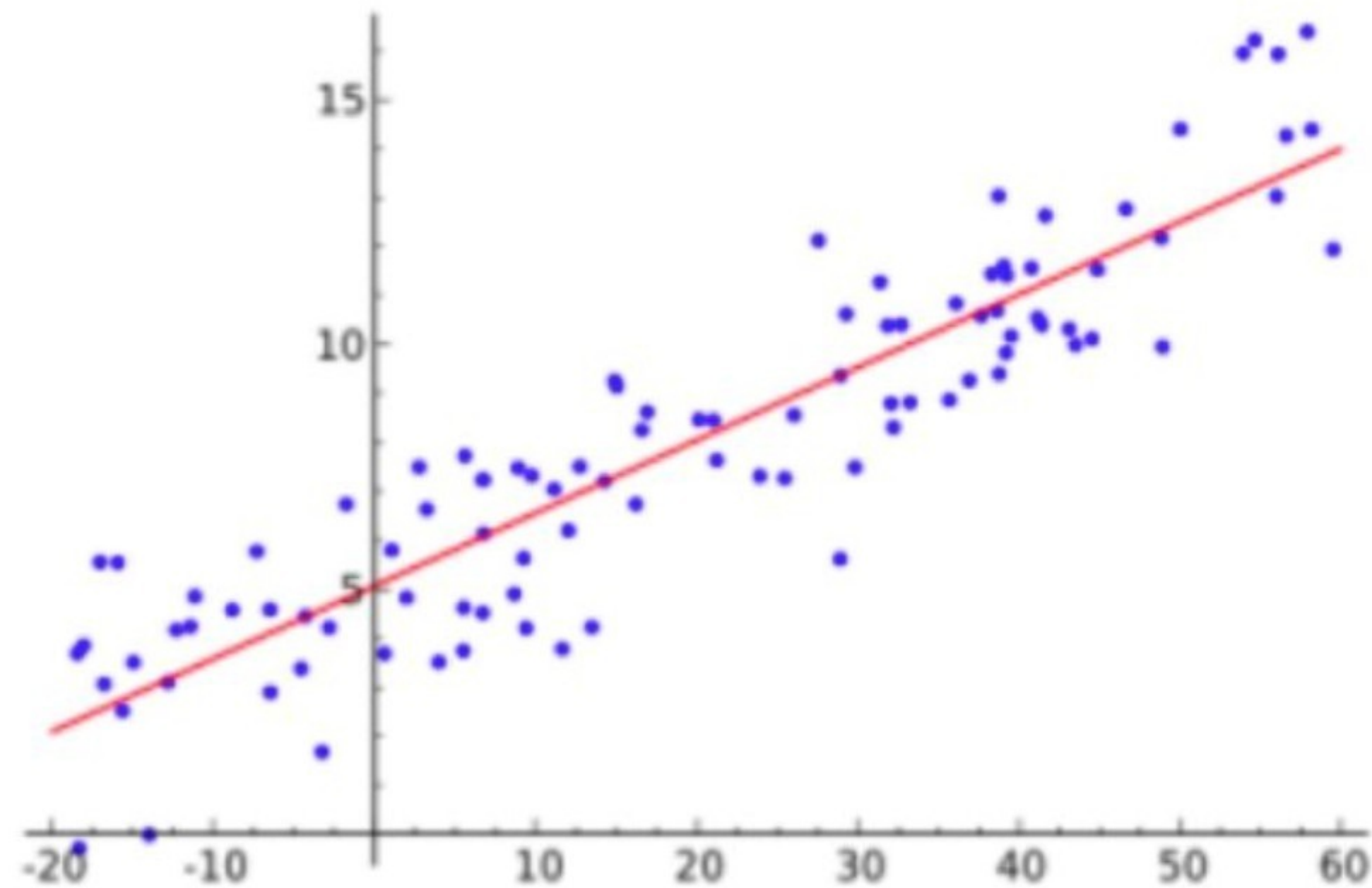
Идеальная модель — порог 50%



Модель с некоторыми ошибками — порог выбирается в зависимости от допускаемых ошибок

Регрессия

Отличается тем, что допустимым ответом является действительное число или числовой вектор.



Спасибо за внимание!