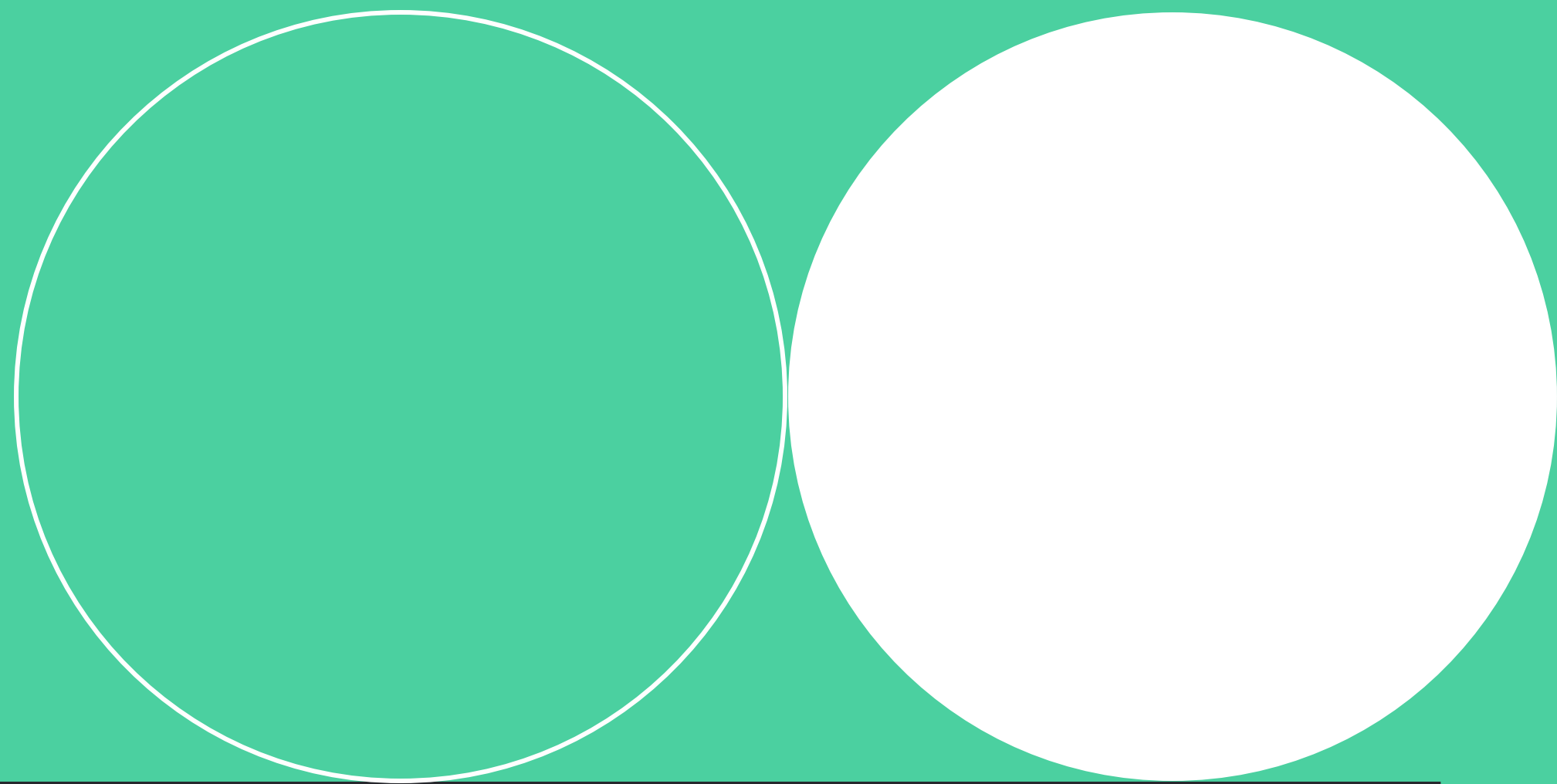

Линейный классификатор и логистическая регрессия

Занятие 1.2



Цели занятия

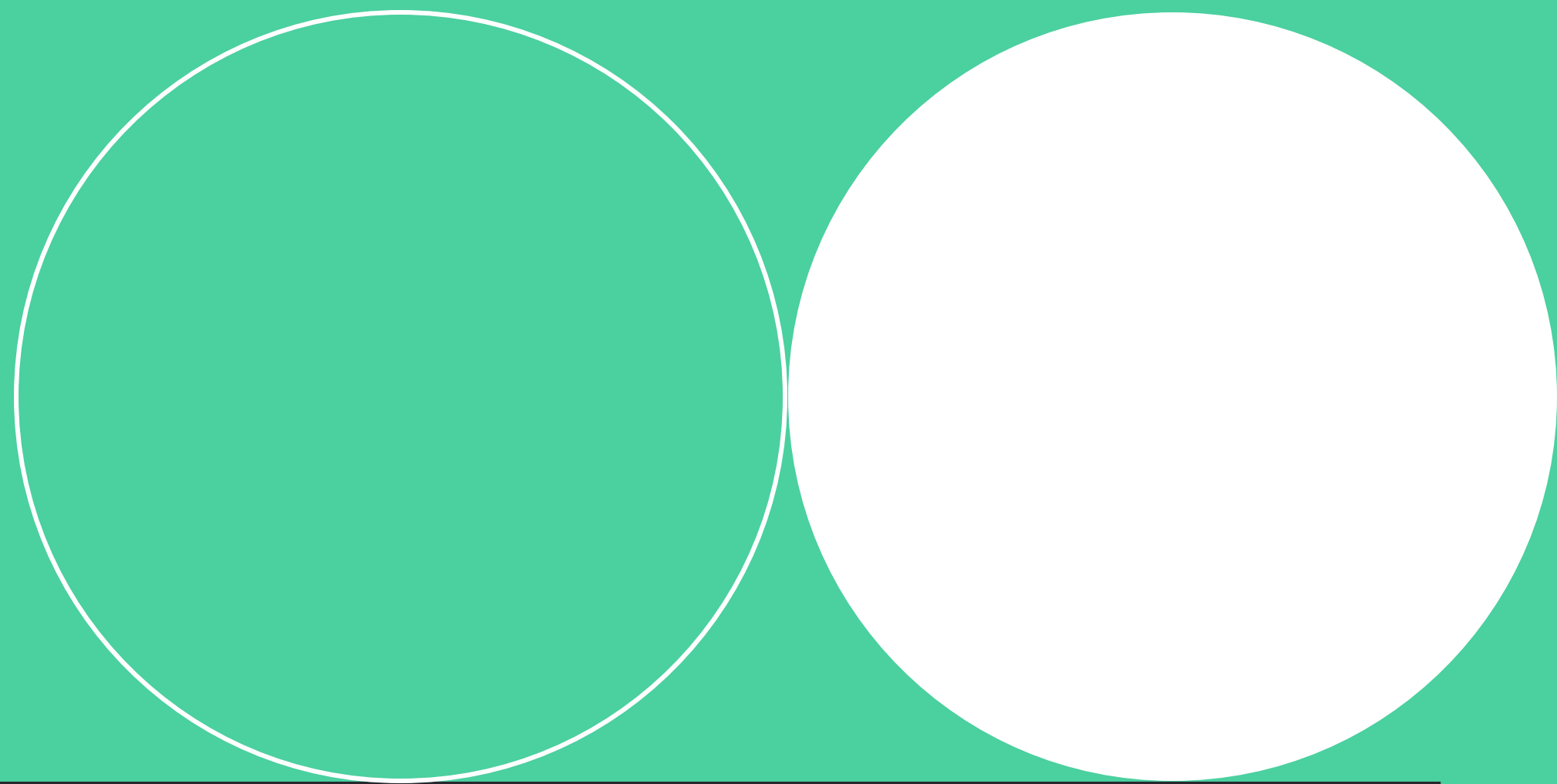


В конце занятия вы:

- будете знать преимущества и недостатки линейных моделей, а также требования к данным;
- научитесь реализовывать алгоритм градиентного спуска и логистическую регрессию;
- повторите понятие условной вероятности.



О чем поговорим и что сделаем



План занятия

1

Линейные модели: требования к данным и практика

2

Логистическая регрессия: практическое задание

3

Градиентный спуск: теория и практическое задание

4

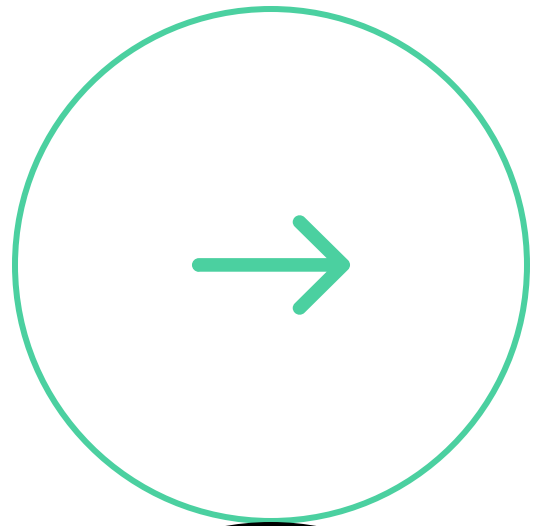
Немного про условную вероятность.



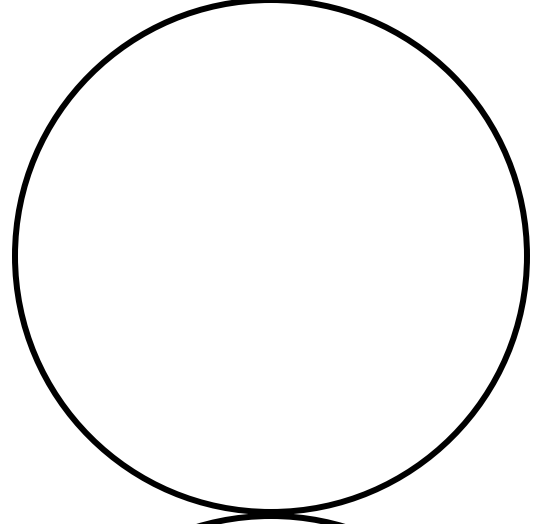
Линейные модели



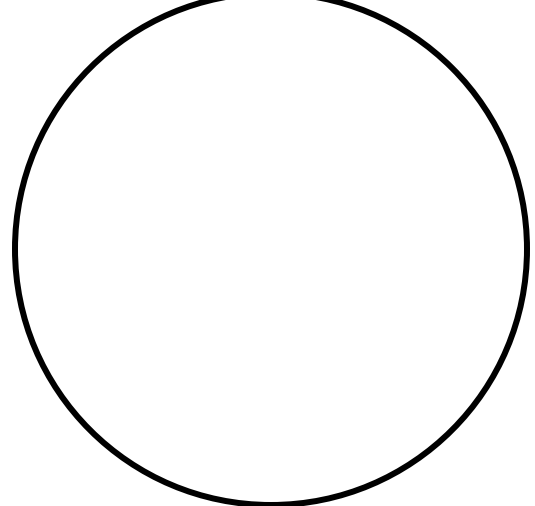
Причины популярности



Линейные модели подходят для описания многих процессов



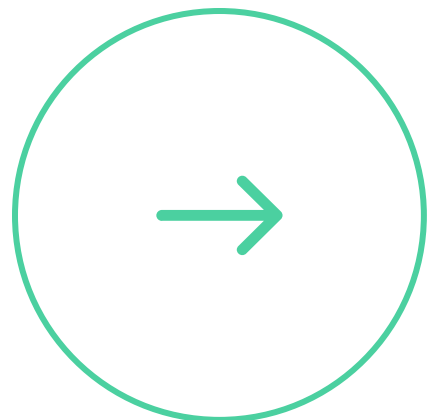
Относительная простота вычислений и интерпретации результатов



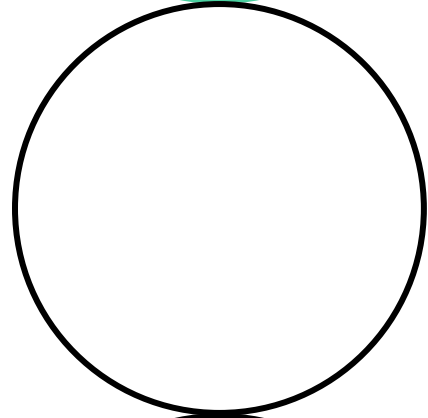
Вклад нескольких факторов часто можно разбить на сумму влияния каждого фактора в отдельности



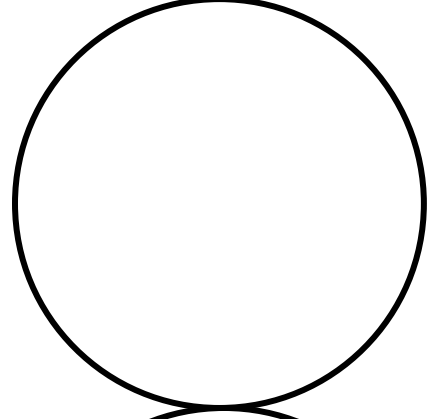
Примеры использования



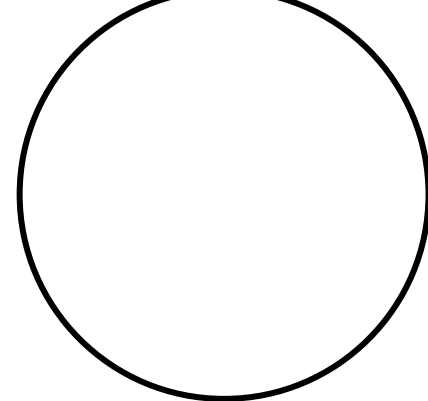
Прогноз продаж по объему инвентаря, загрузке, площади и другим «линейным» характеристикам



Построение вероятностных моделей в страховании, кредитном скоринге, инвестиционных проектах



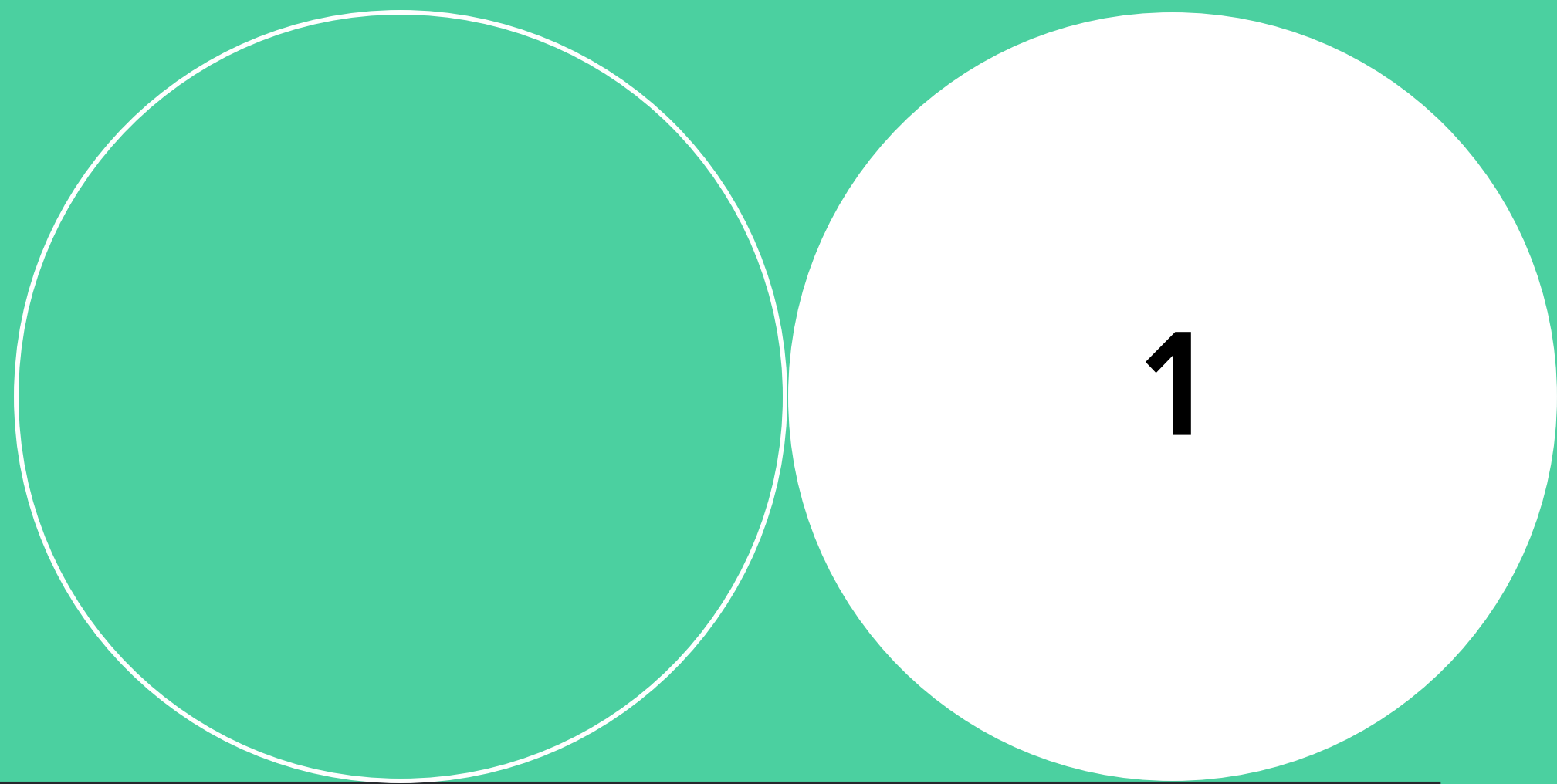
Предсказание цены товара на основании его характеристик



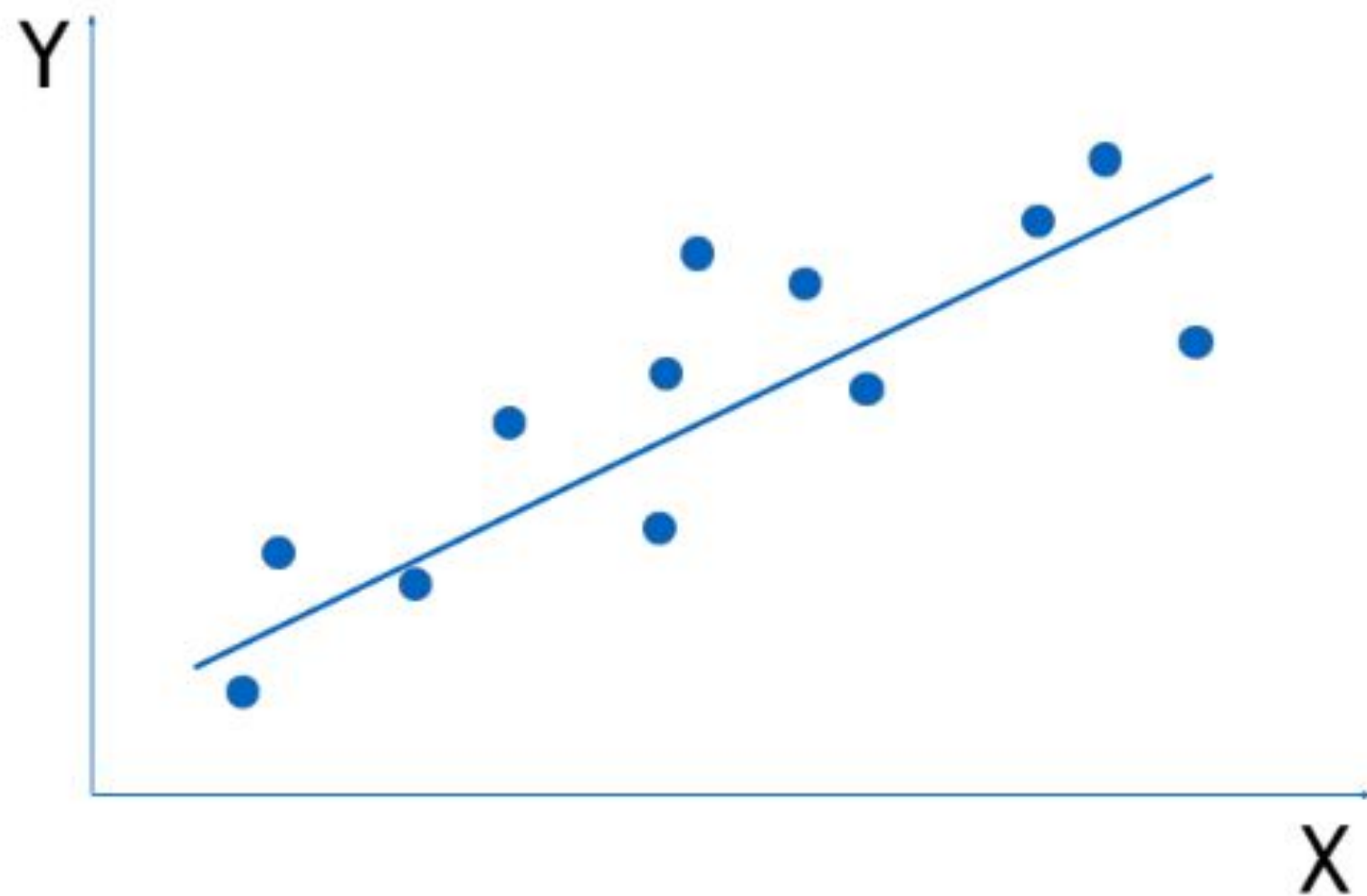
Построение трендов



Определение и код



Определение

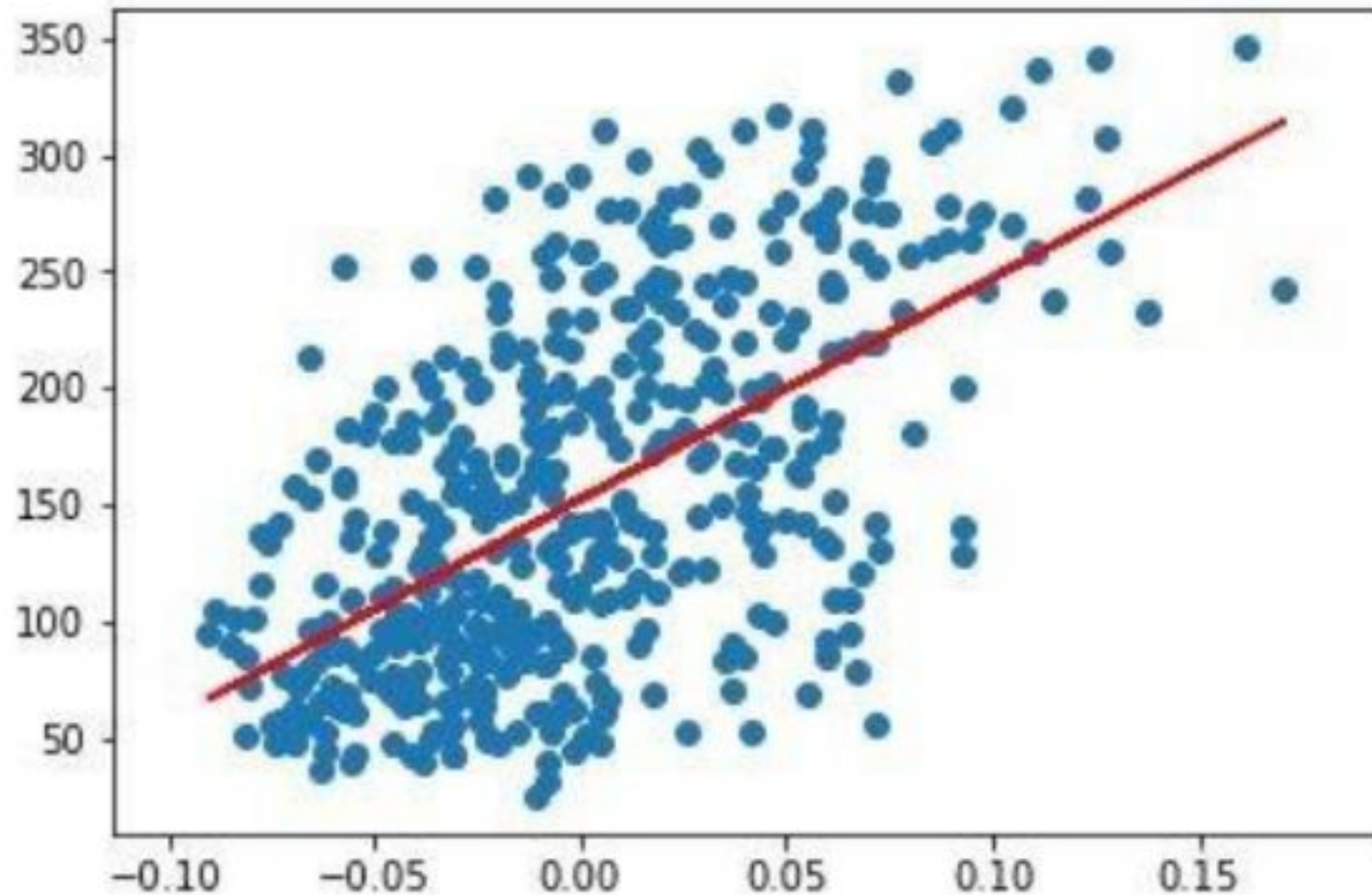


$$y_i = \sum_{j=1}^m w_j X_{ij} + e_i$$

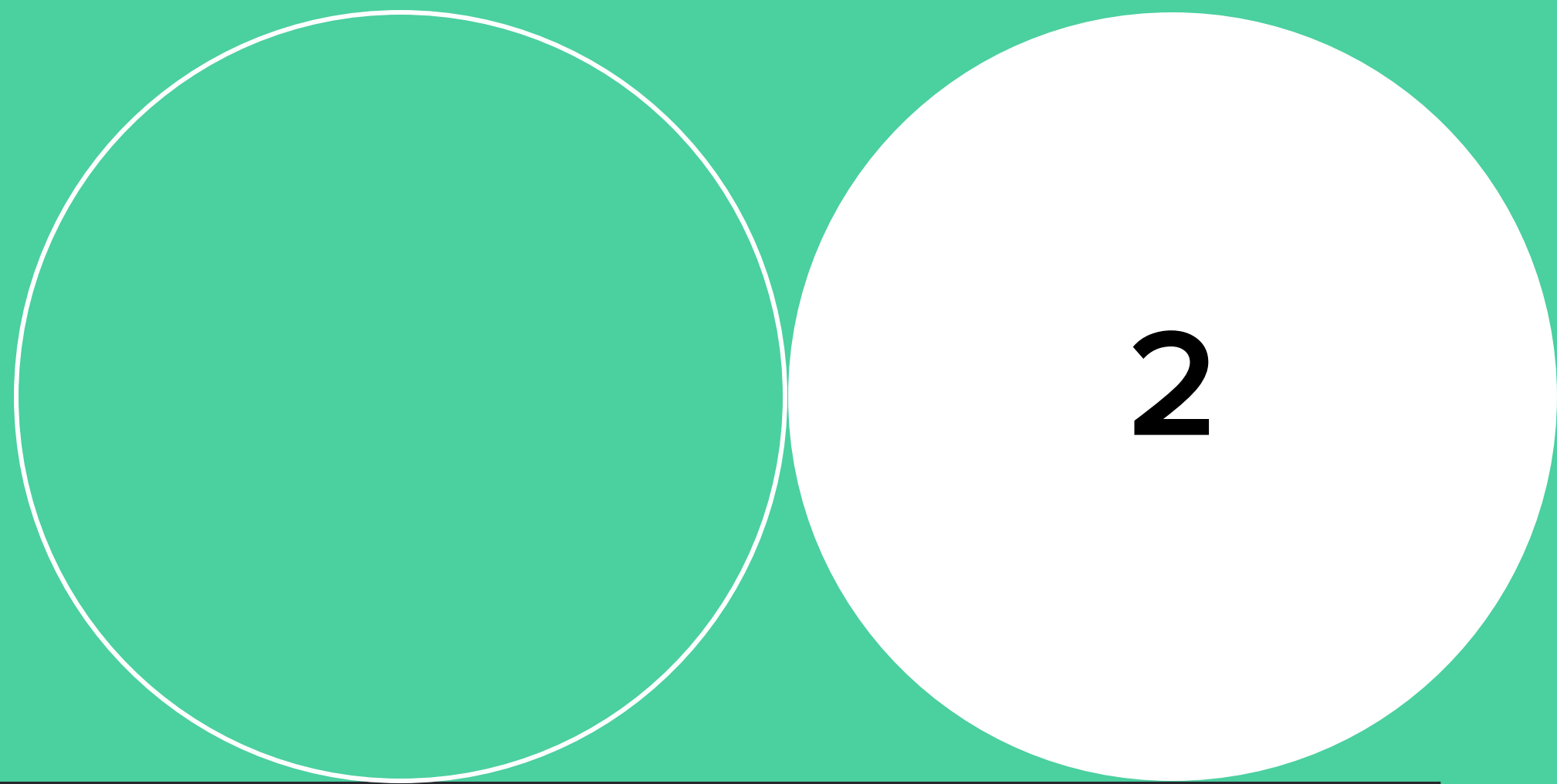
Y - целевая переменная
W - вектор весов модели
X - матрица наблюдений
e - ошибка модели



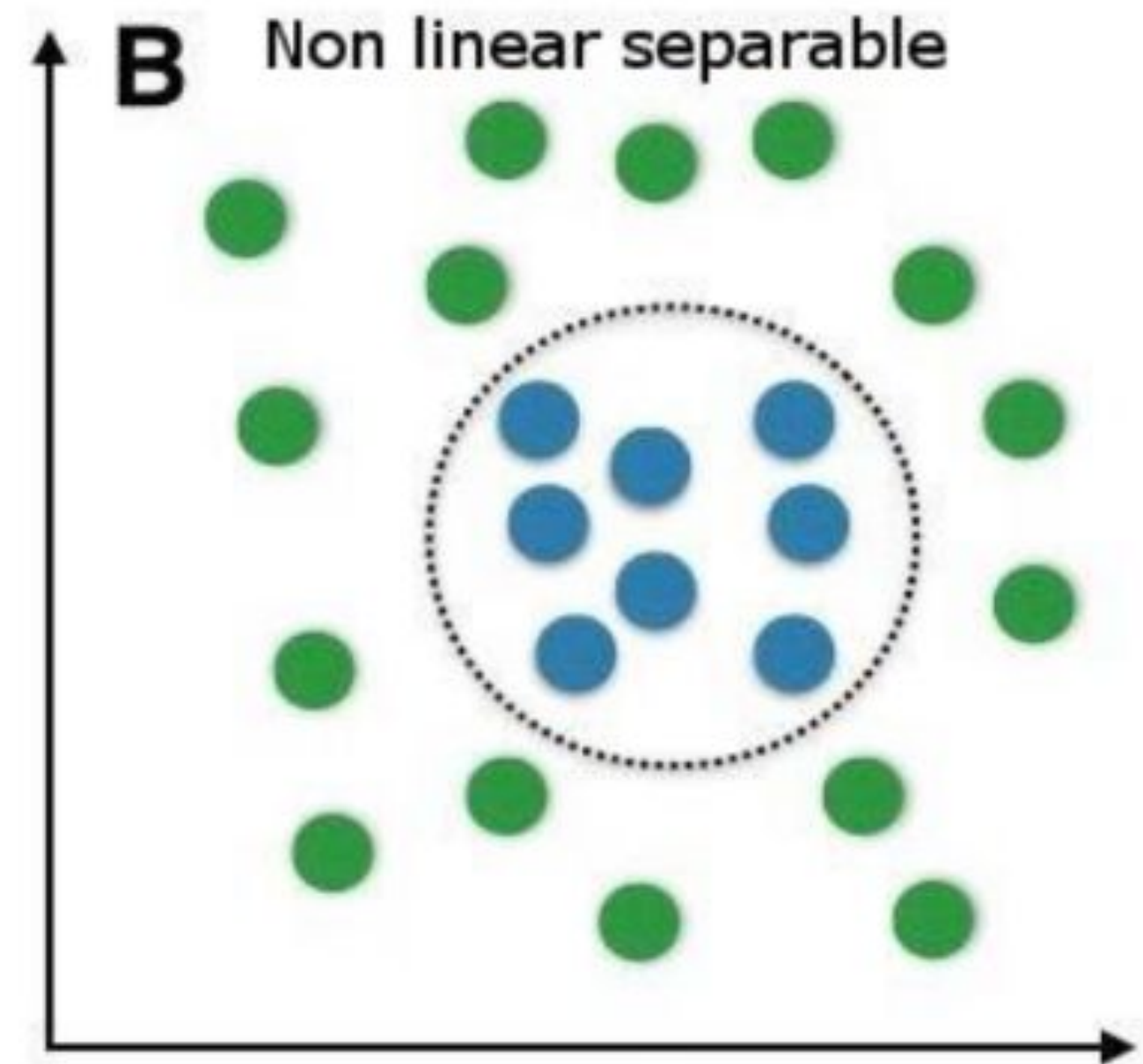
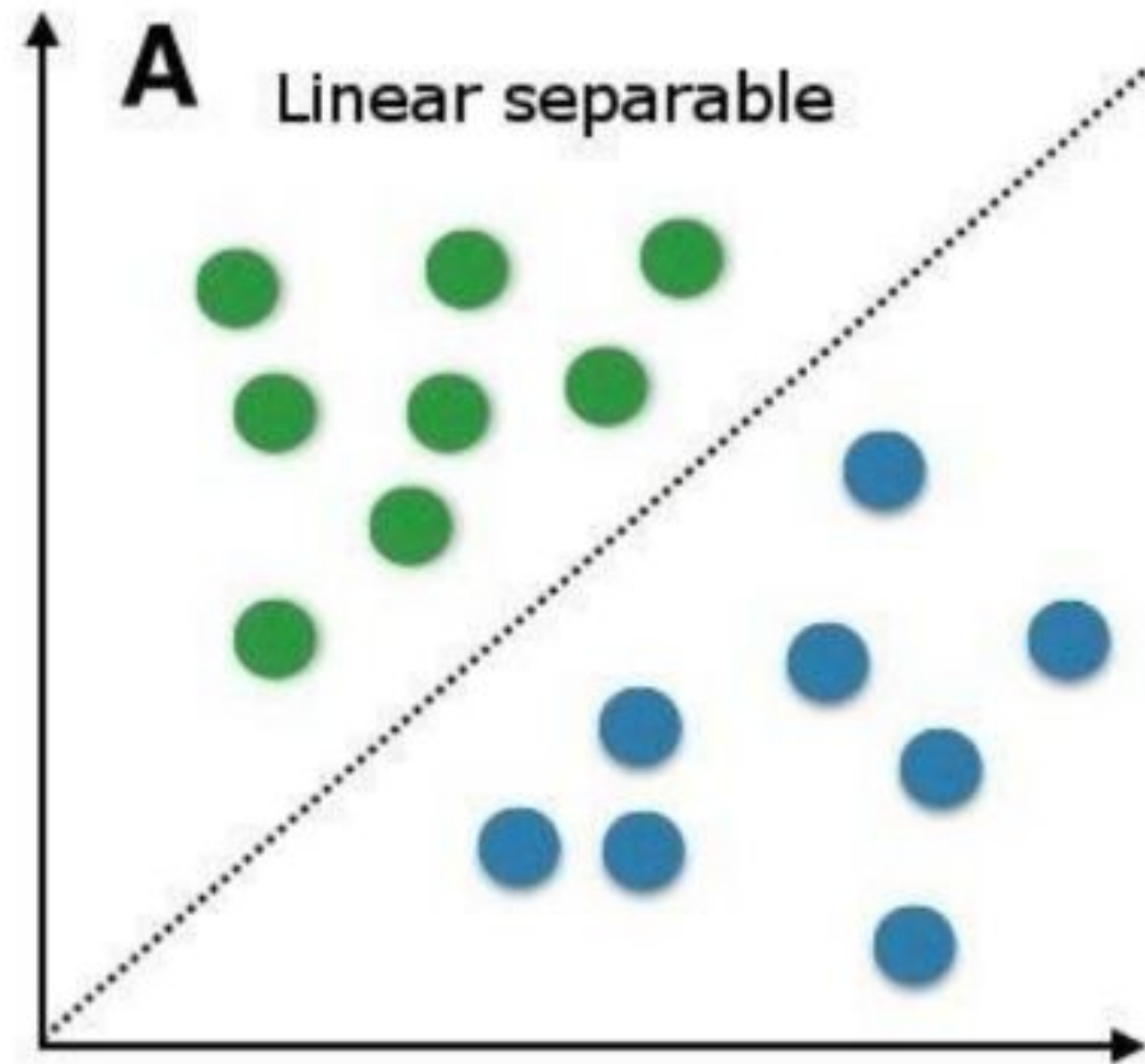
Пример из кода LINEAR REGRESSION.IPYNB



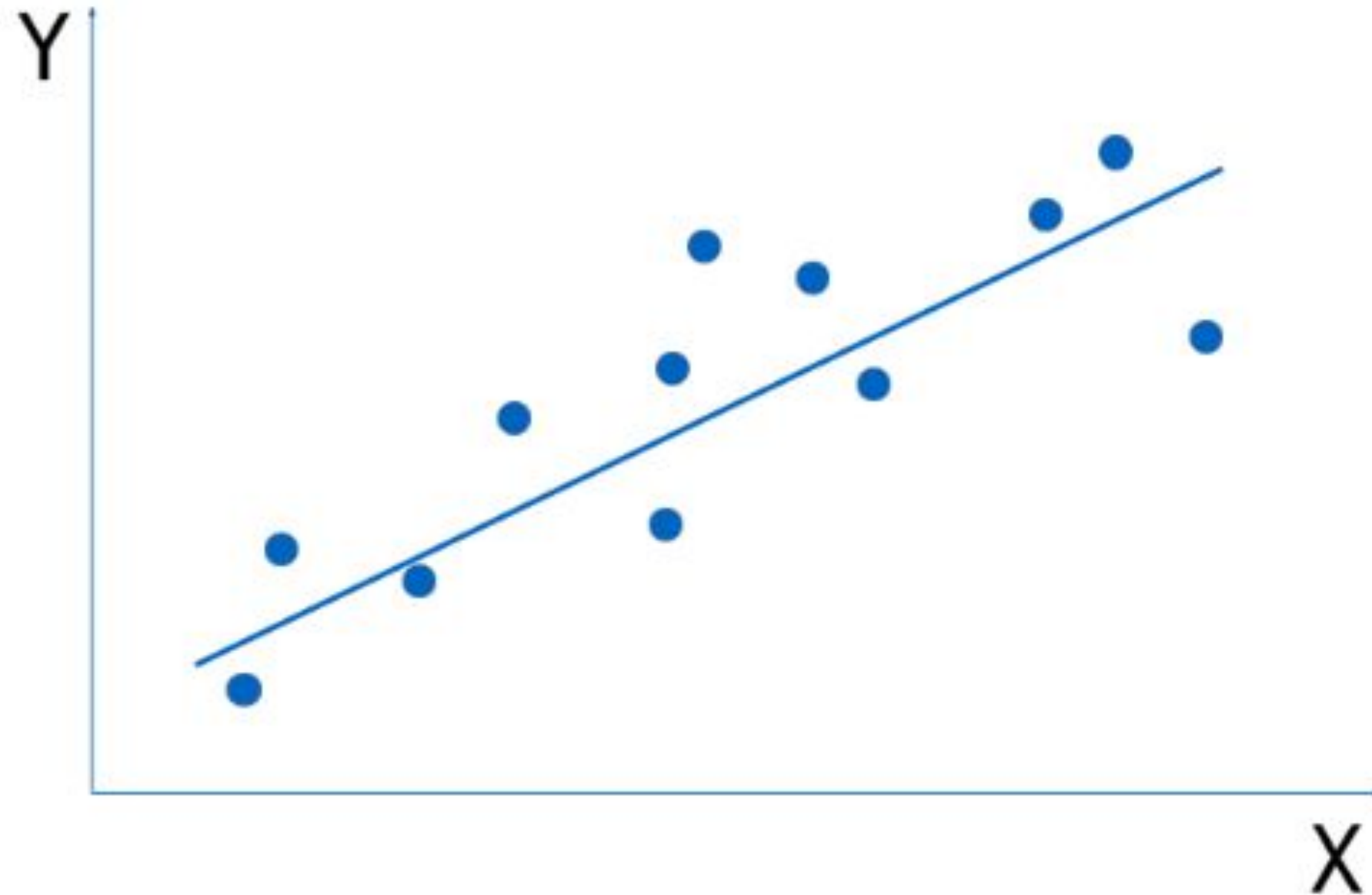
Построение линейной модели



Как строим линейную модель



Метод максимального правдоподобия

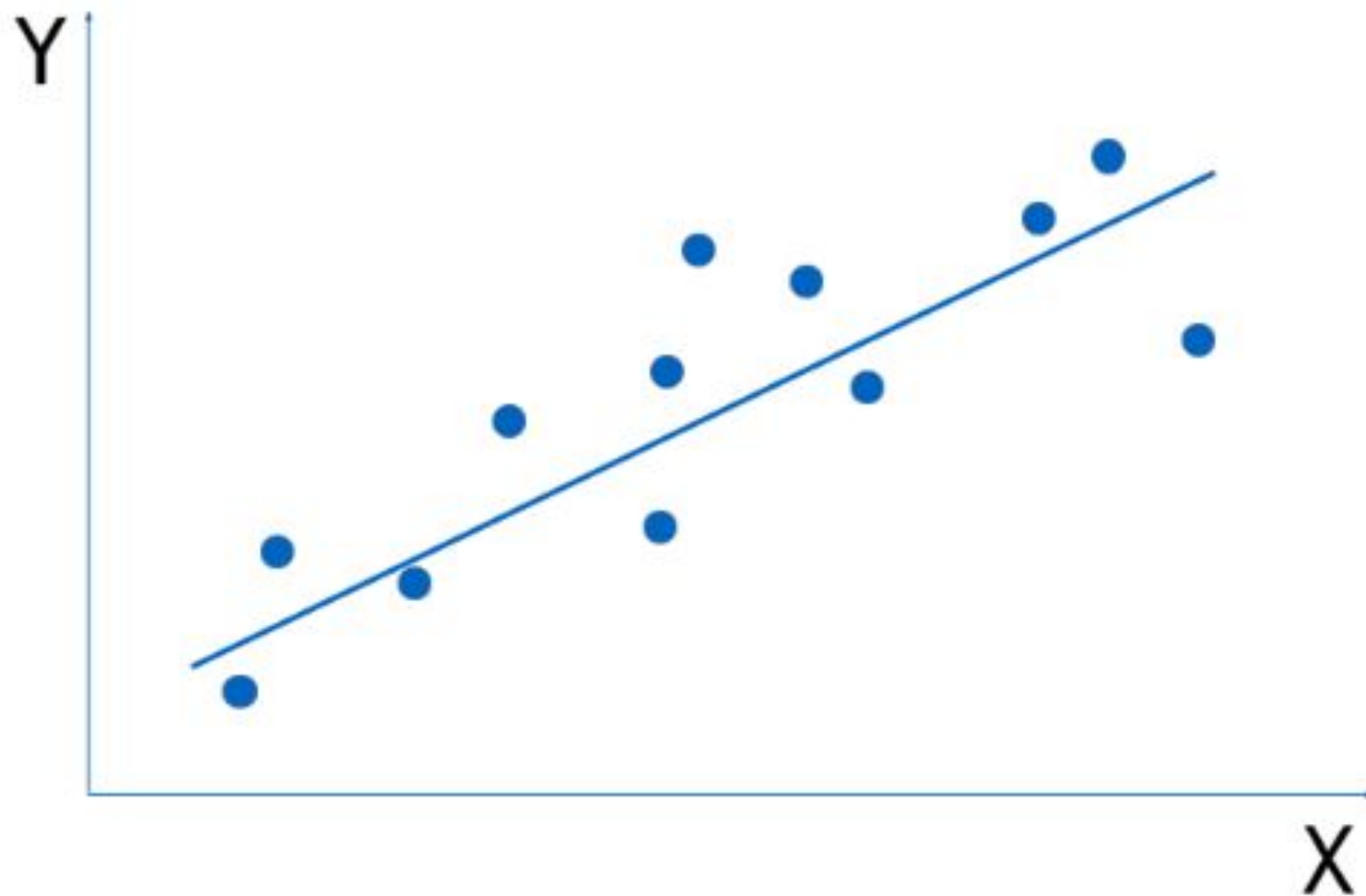


Как можно получить эту прямую?

$p(y | x, a)$ — вероятность получить y при входных данных x . a — параметр модели



Метод максимального правдоподобия



Как можно получить эту прямую?

$p(y | x, \alpha)$ — вероятность получить y при входных данных x . α — параметр модели

Введем функцию:

$$W(\alpha) = \prod_i p(x_i, \alpha)$$



Метод максимального правдоподобия

Функция максимального правдоподобия:

$$L(\alpha) = \sum_i \log p(x_i, \alpha)$$



Метод максимального правдоподобия

Функция максимального правдоподобия:

$$L(\alpha) = \sum_i \log p(x_i, \alpha)$$

Как подобрать значение α , чтобы максимизировать $L(\alpha)$?

Необходимо минимизировать среднеквадратичную ошибку между прогнозными и фактическими значениями



Доказательство
<https://habrahabr.ru/company/ods/blog/323890/#metod-maksimalnogo-pravdopodobiya>



Время кода

REGRESSION_CARS.IPYNB



Практическое задание 1



Время практики

SAT_MODEL.IPYNB



Логистическая регрессия



4

Прогноз вероятности

Прогнозирует вероятность отнесения наблюдения к определенному классу

Модель:

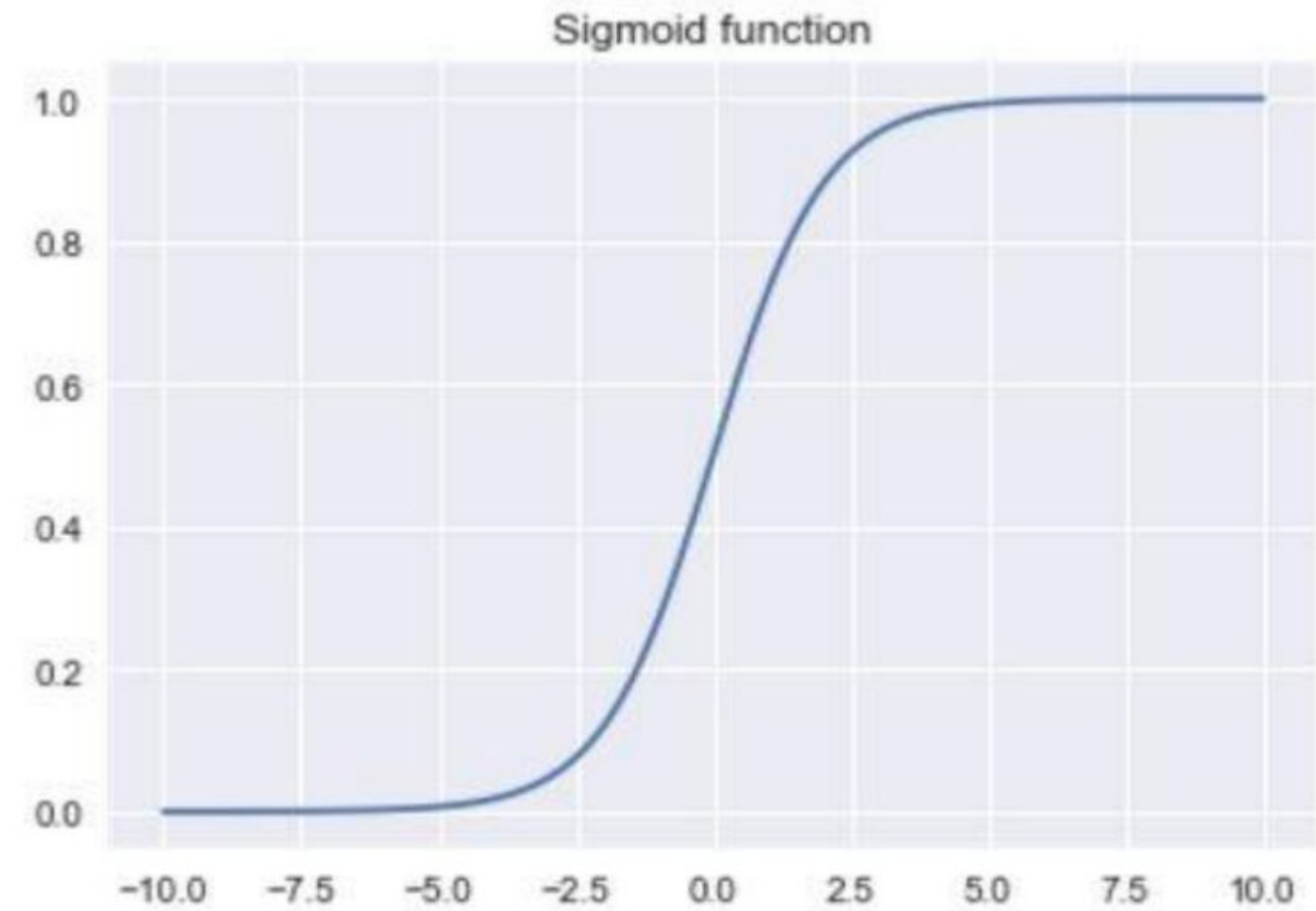
$$L = a_0 + a_1X_1 + a_2X_2 + \dots + a_nX_n$$



Прогноз вероятности

Вероятность:

$$p = \frac{1}{1 + e^{-L}}$$

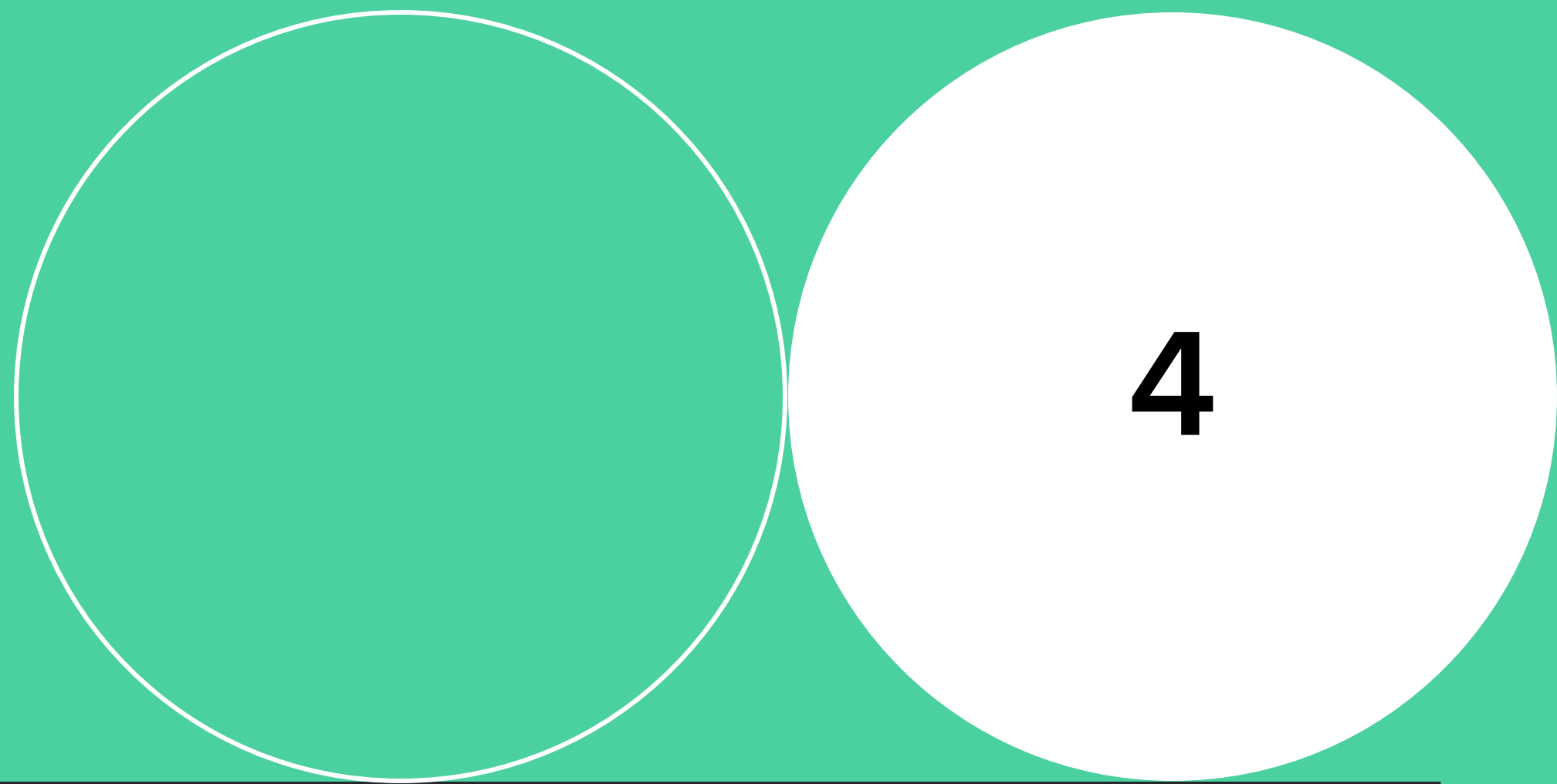


Основа практики

LOGISTIC_REGRESSION_ATHLETES_CLASSIFIER.IPYNB



Практическое задание 2

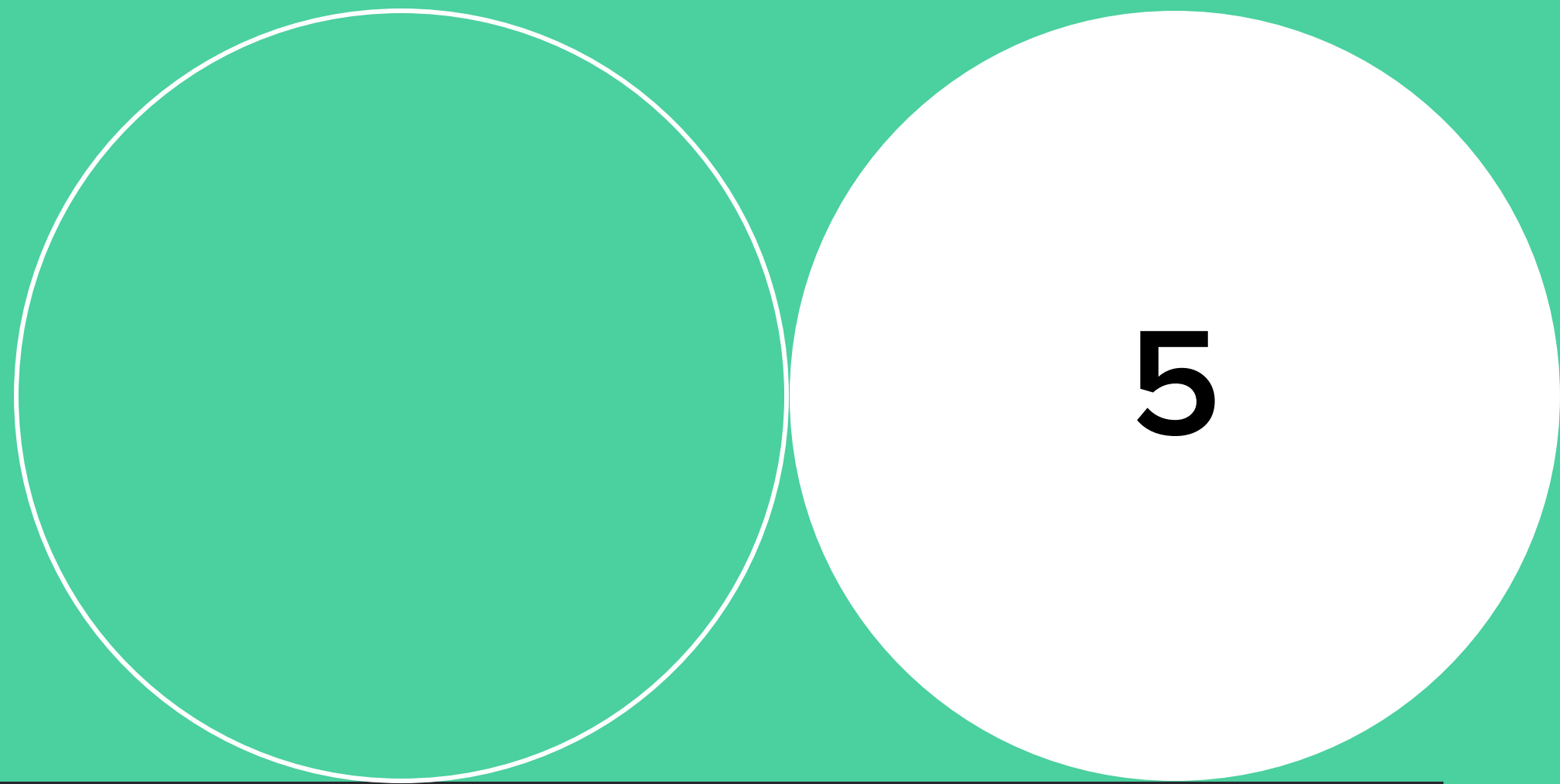


Улучшаем точность модели

С новыми признаками



Если классов больше двух

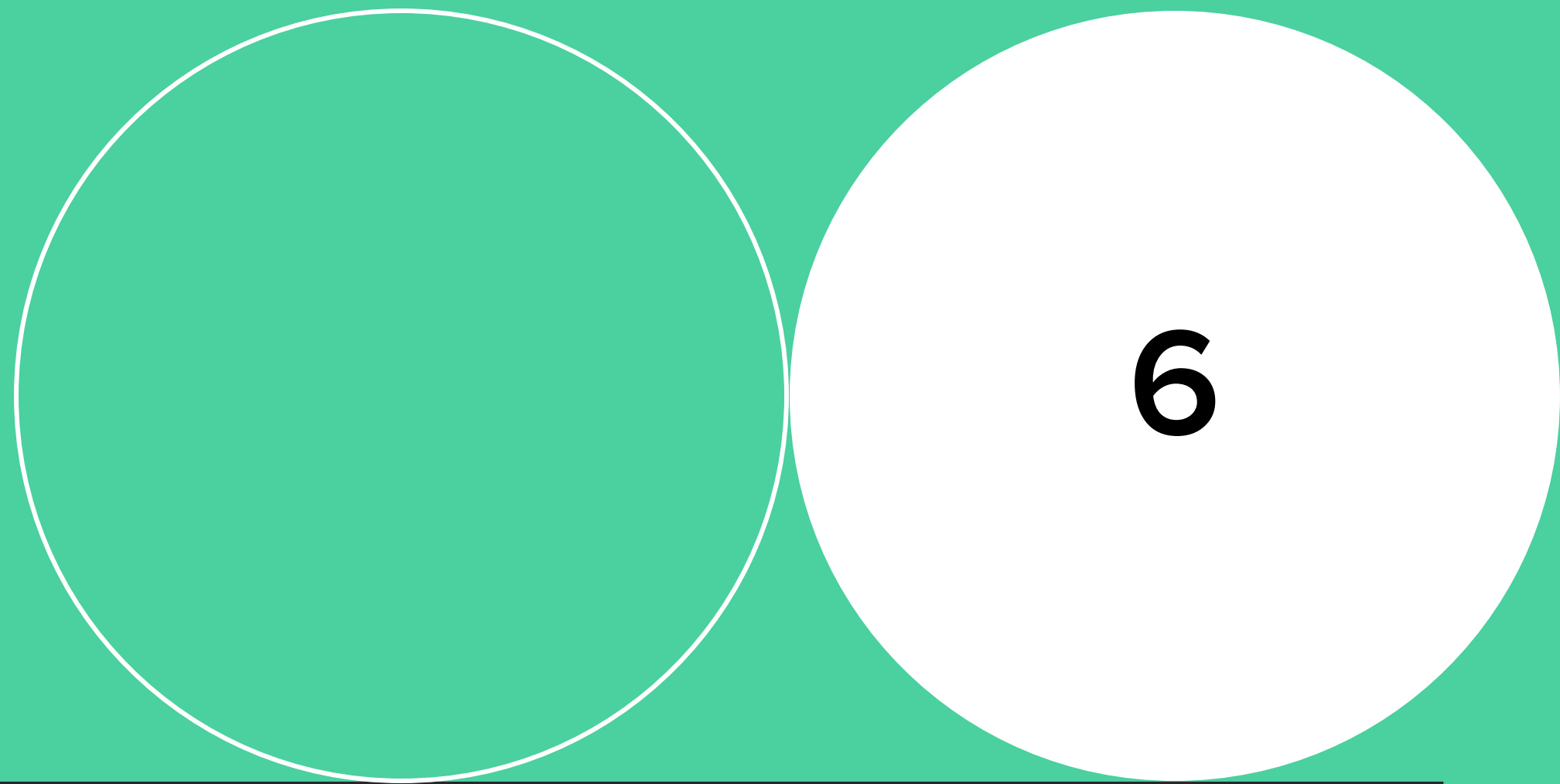


Пример

IRIS_DATASET.IPYNB



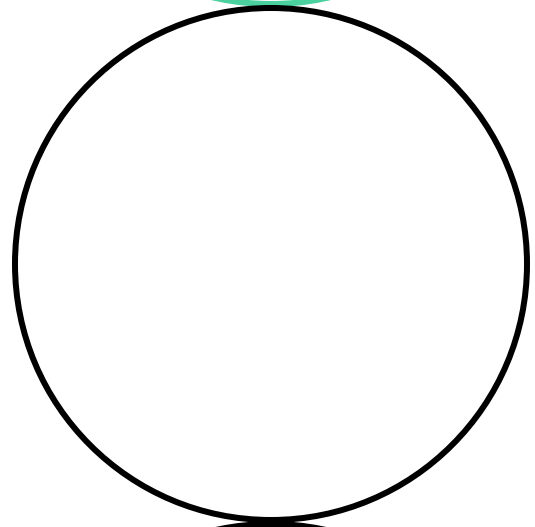
Для каких данных это работает?



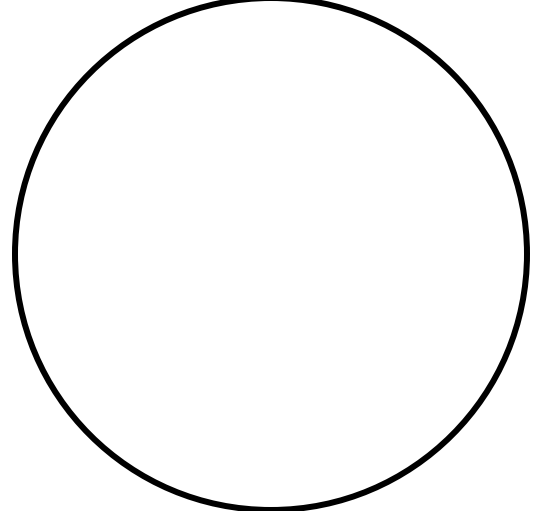
Требования к данным



**Линейная зависимость
целевой переменной**



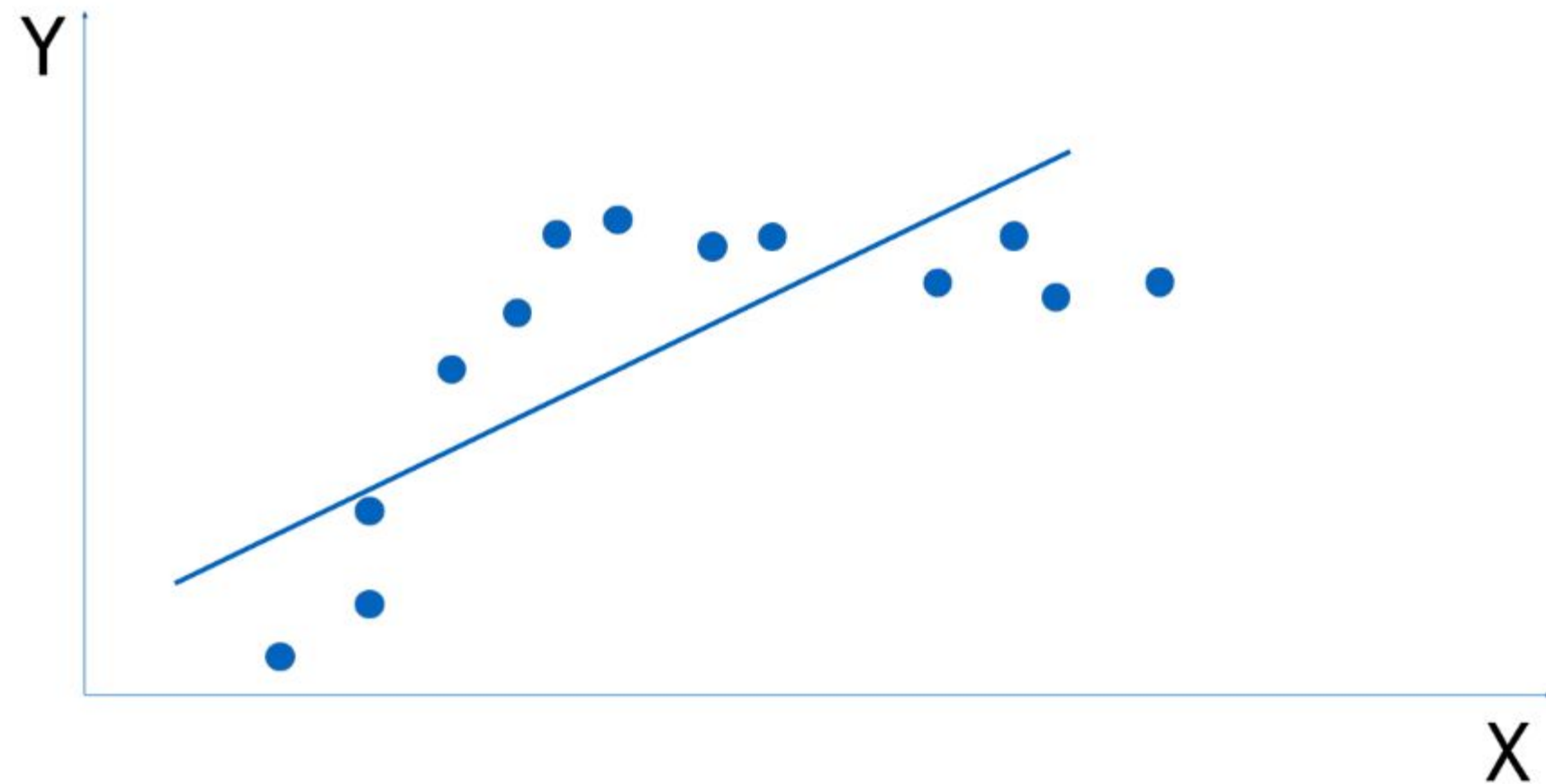
**Нормальное
распределение остатков**



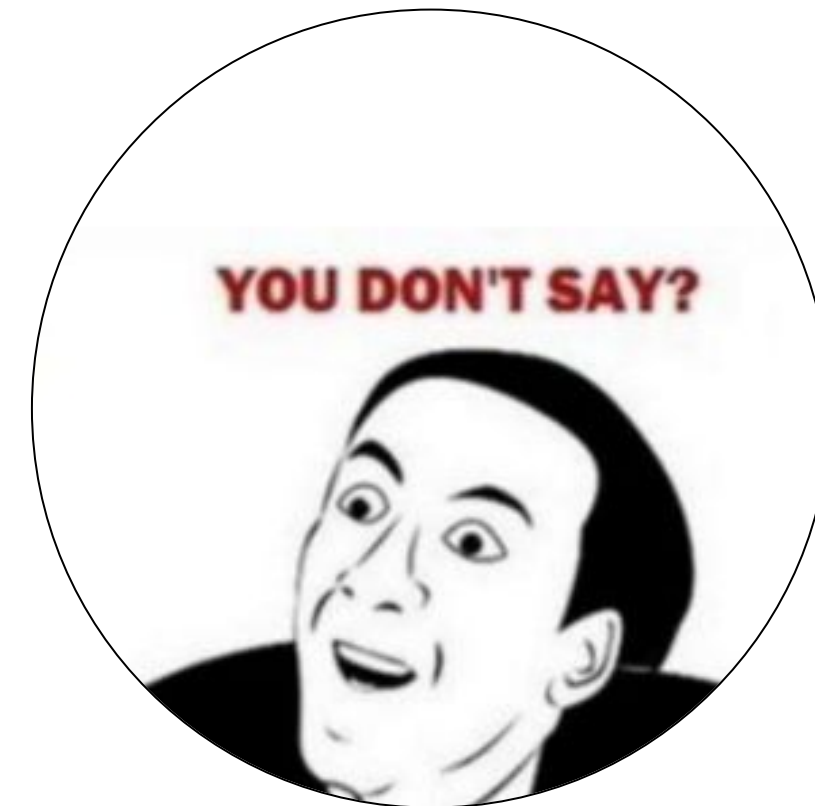
**Постоянная
изменчивость остатков**



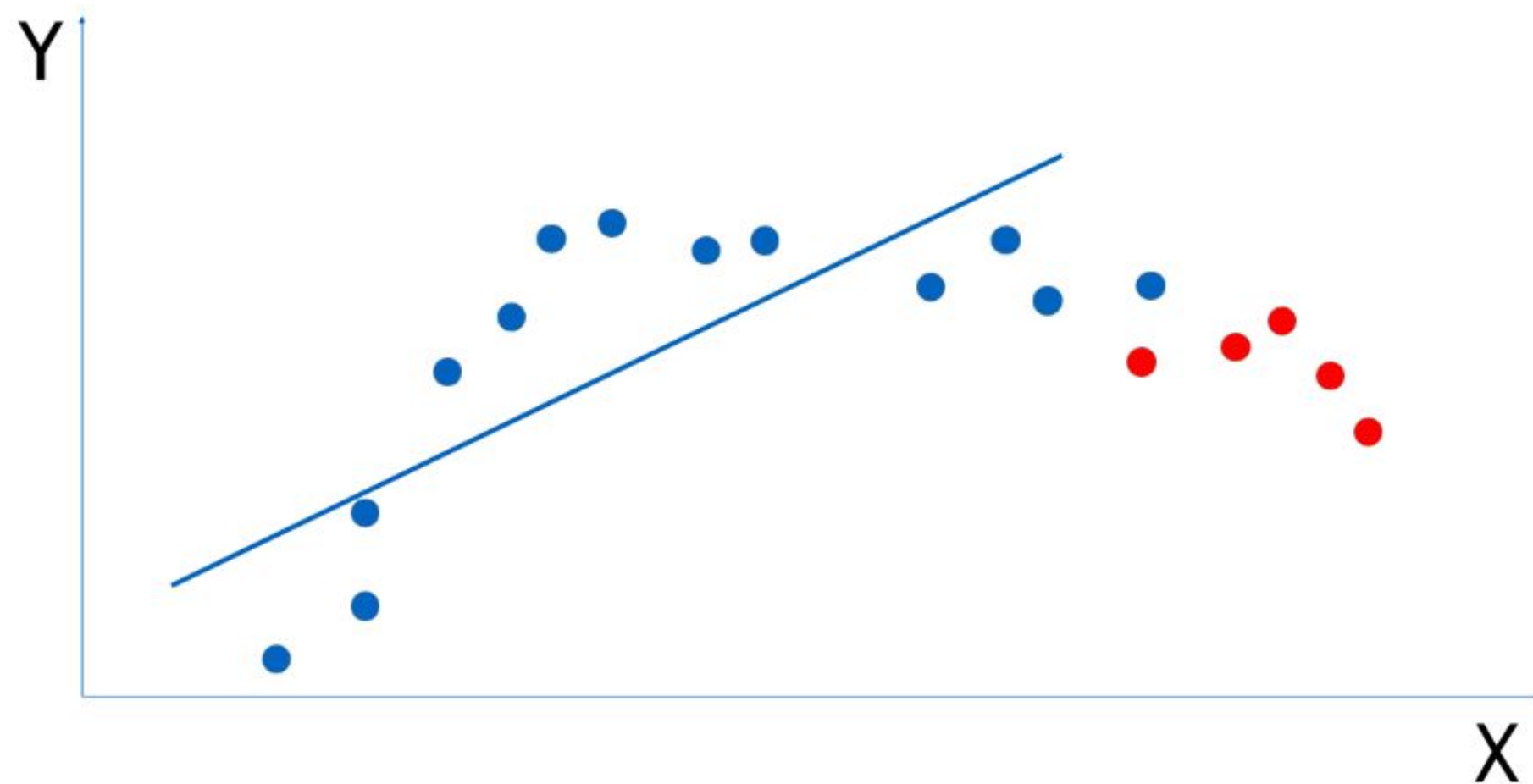
Требования к данным



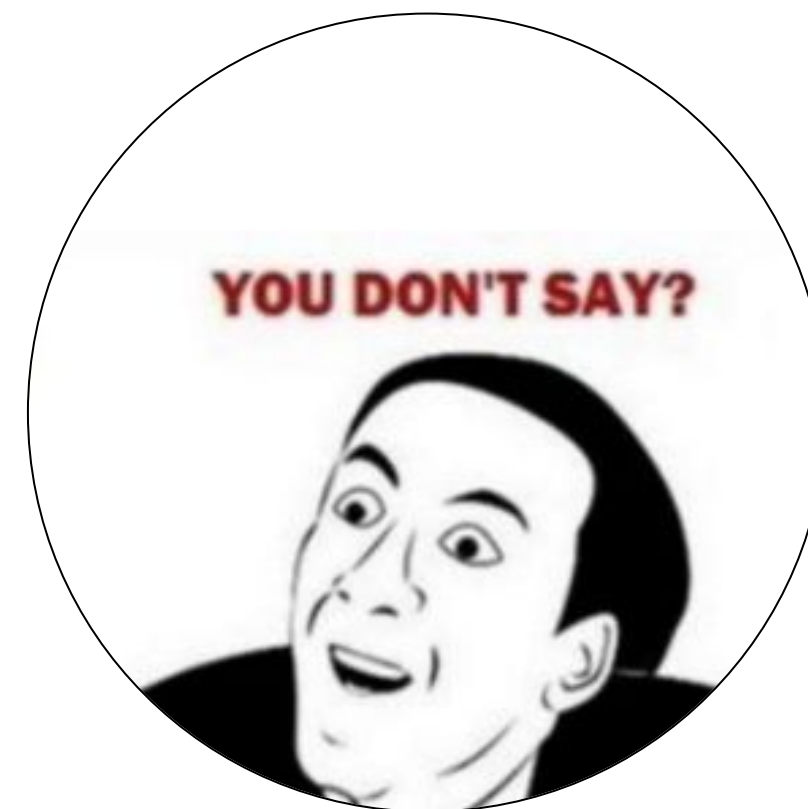
Линейная взаимосвязь X и Y



Требования к данным

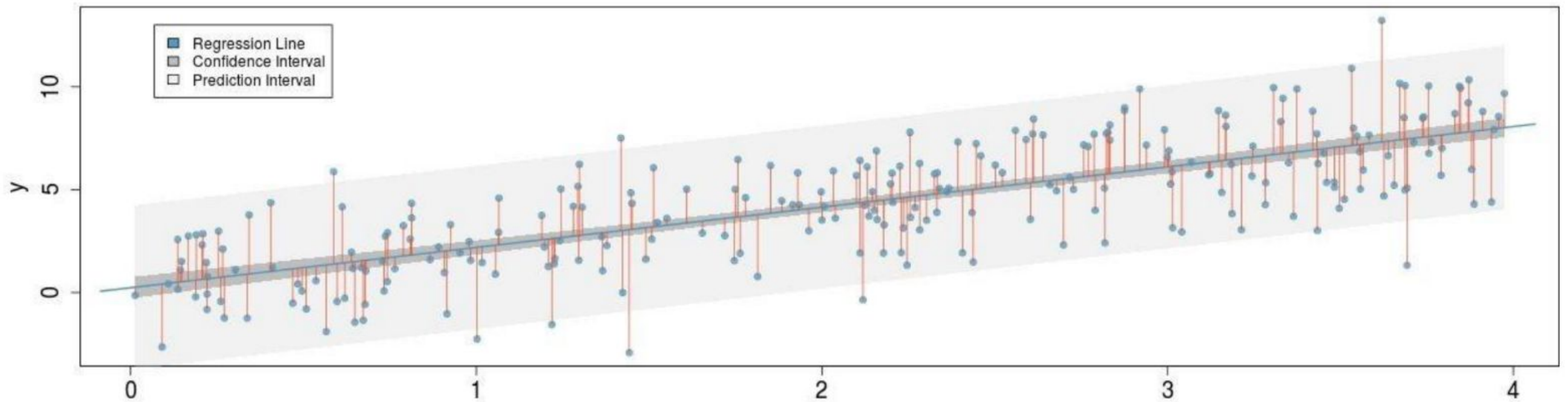


Линейная взаимосвязь X и Y



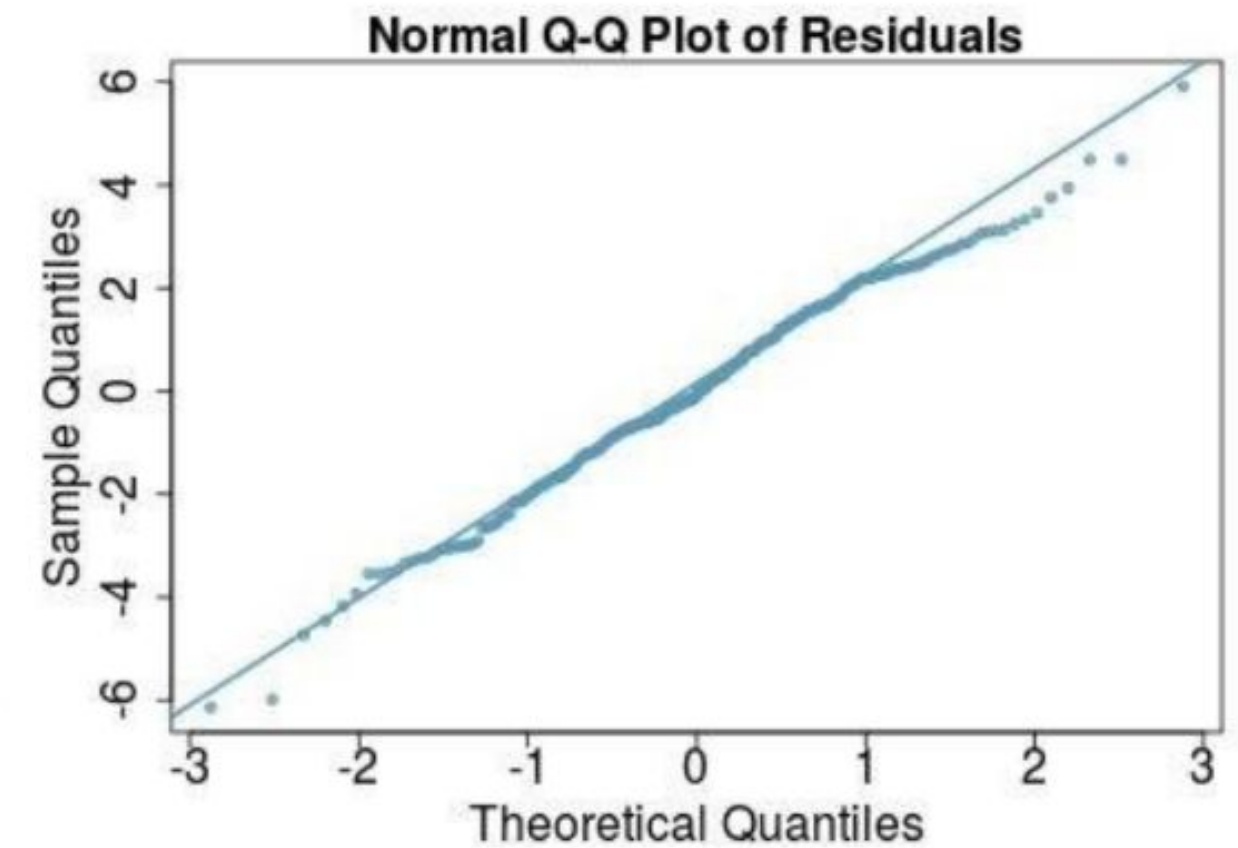
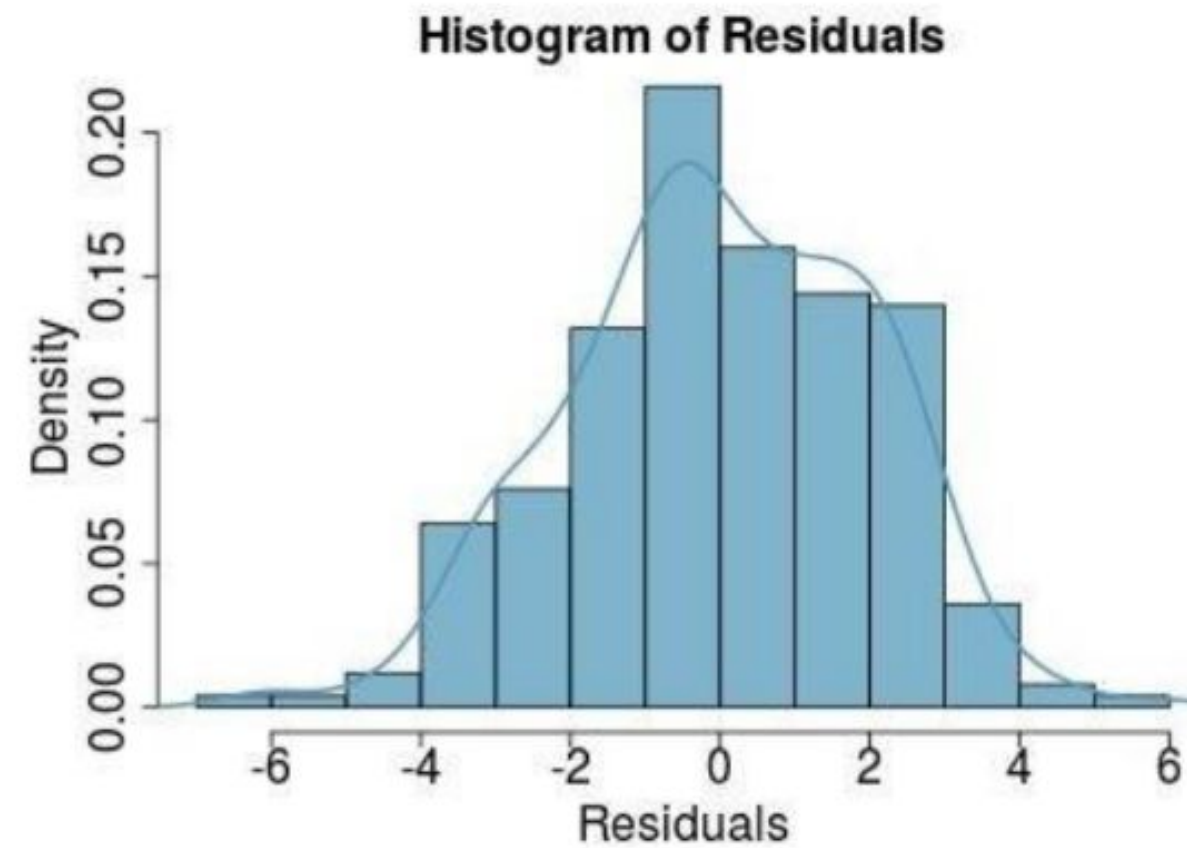
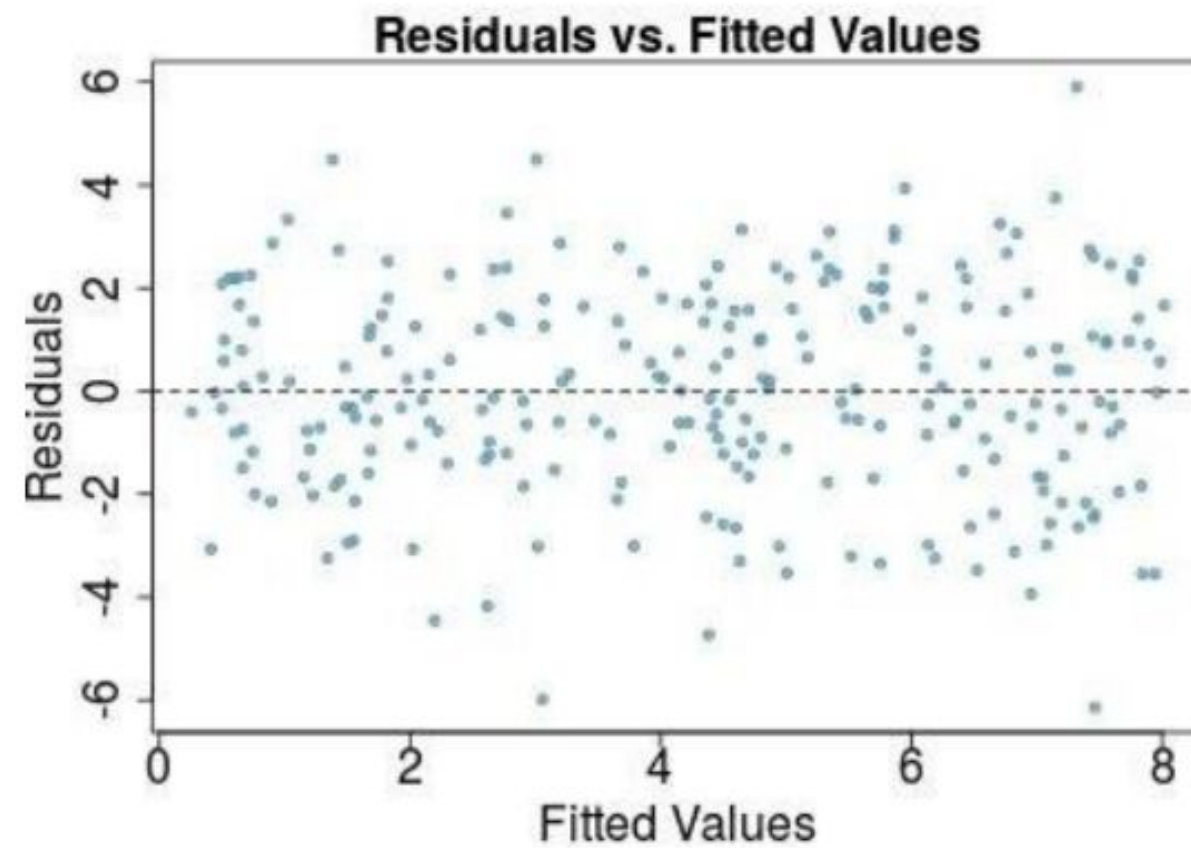
Нормальное распределение остатков

[HTTPS:// GALLERY.SHINYAPPS.IO/SLR_DIAG/](https://gallery.shinyapps.io/slr_diag/)



Нормальное распределение остатков

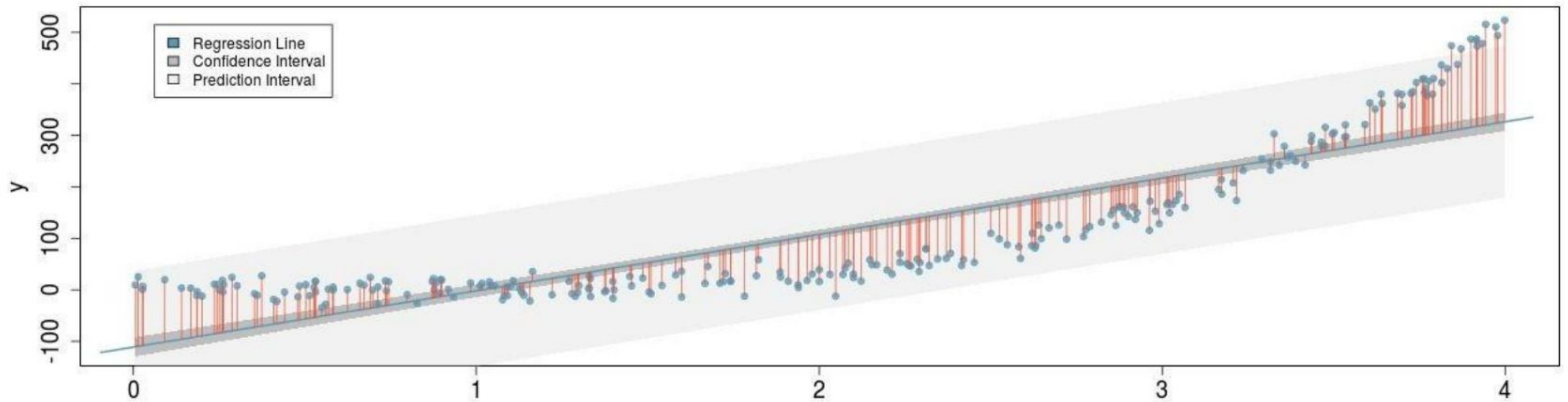
[HTTPS:// GALLERY.SHINYAPPS.IO/SLR_DIAG/](https://gallery.shinyapps.io/slr_diag/)



Гомоскедастичность

Постоянная изменчивость остатков

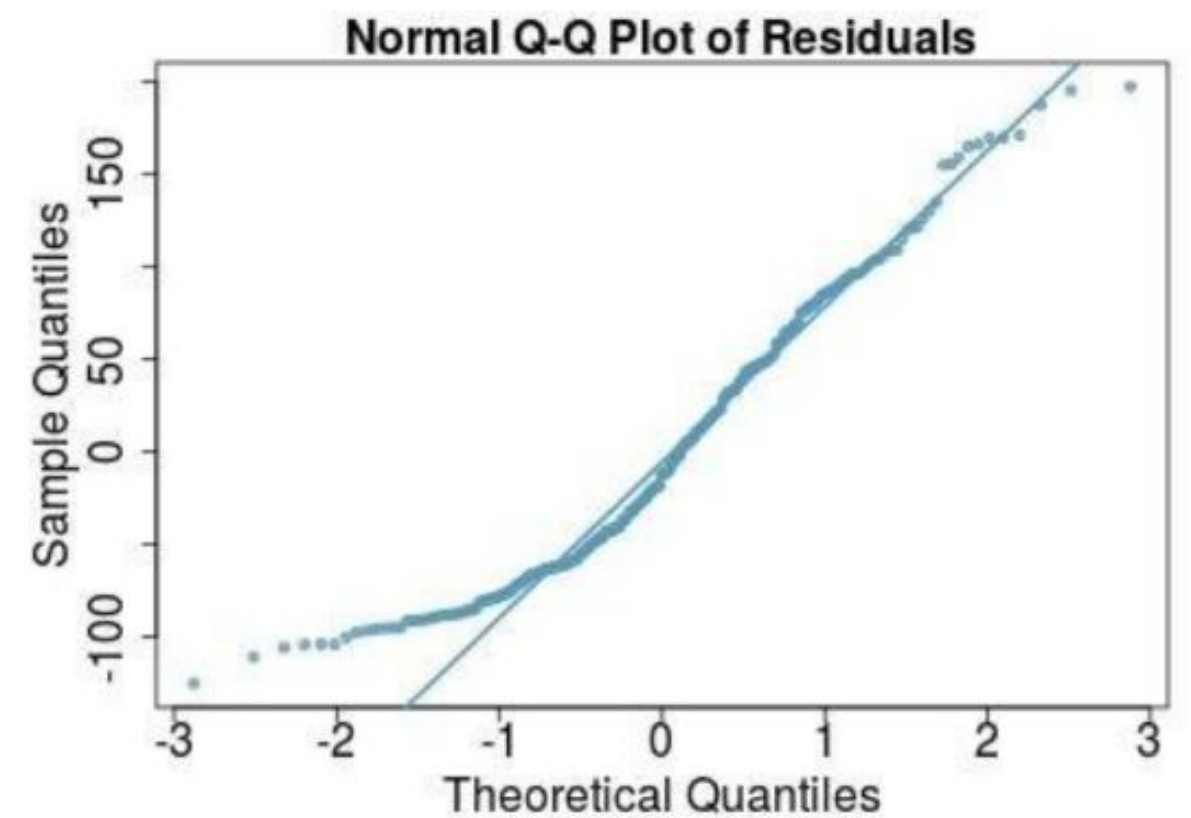
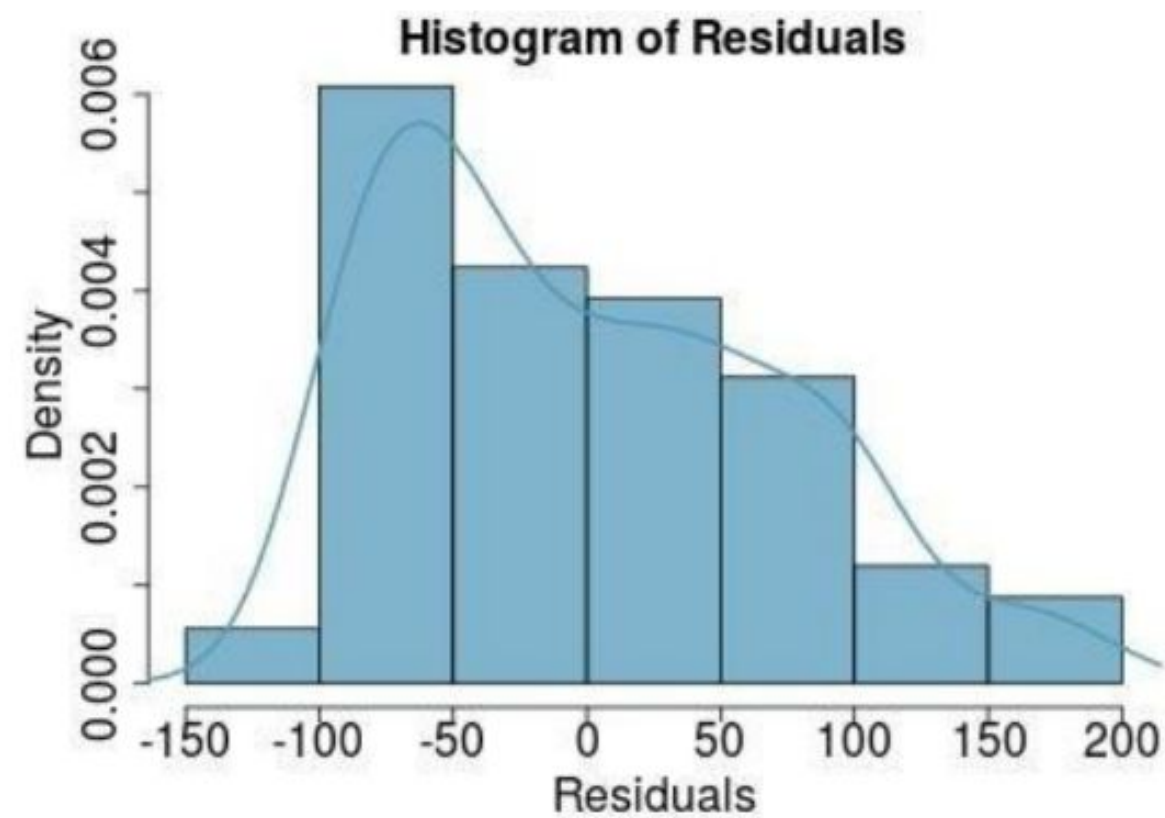
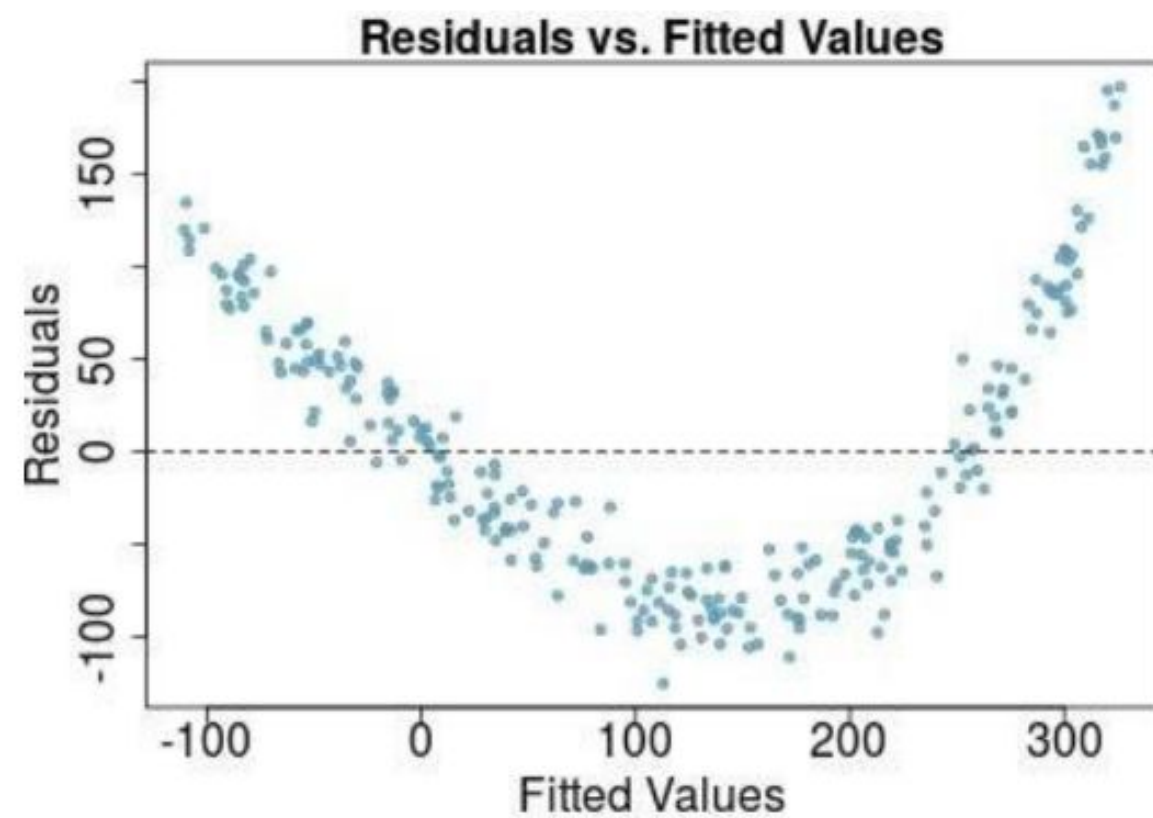
Пример постоянной гетероскедастичной последовательности



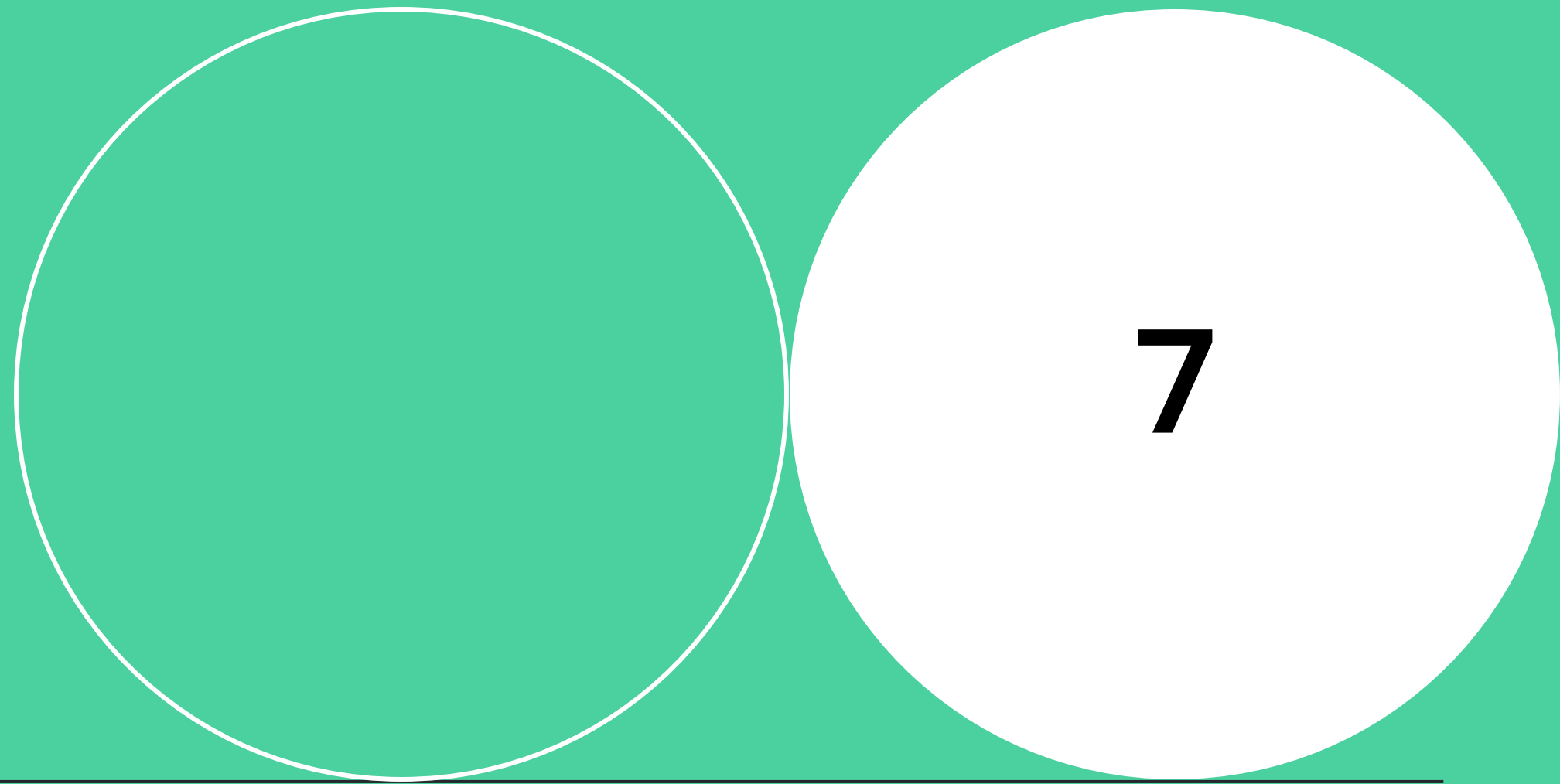
Гомоскедастичность

Постоянная изменчивость остатков

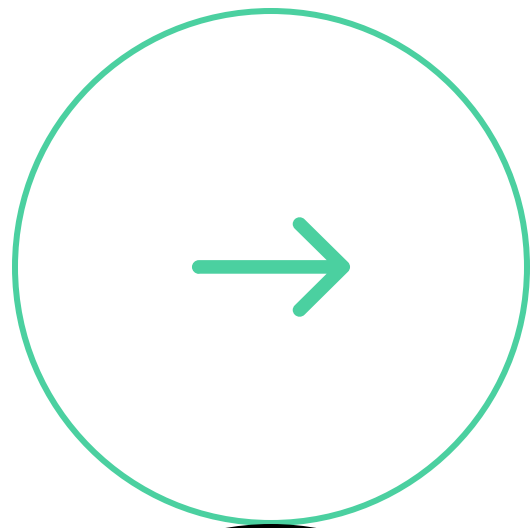
Пример постоянной гетероскедастичной последовательности



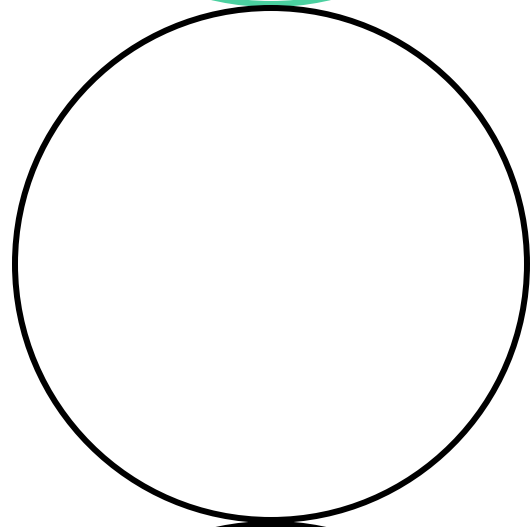
SVM



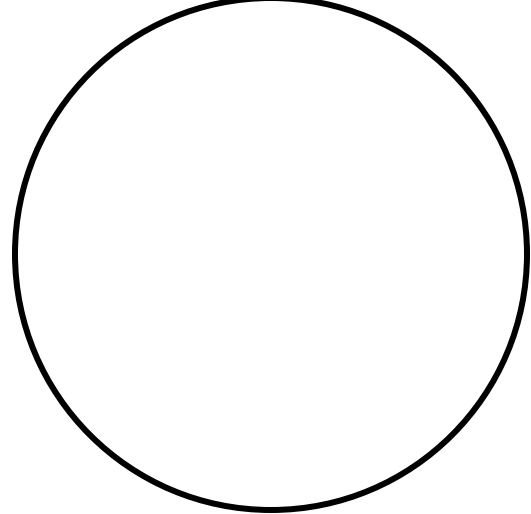
Множество гиперплоскостей



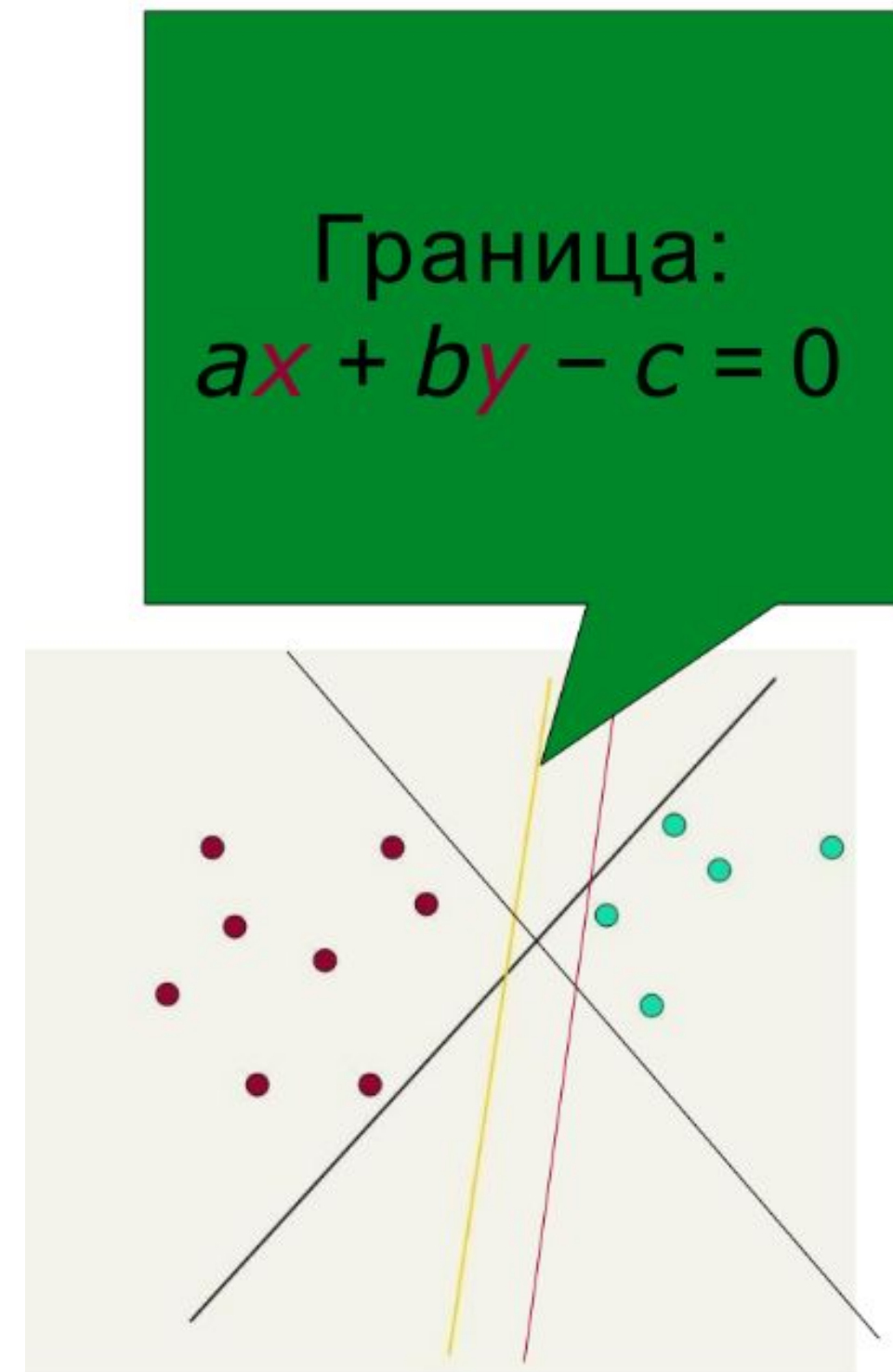
Множество решений для a, b, c .



SVM находит оптимальную разделяющую поверхность

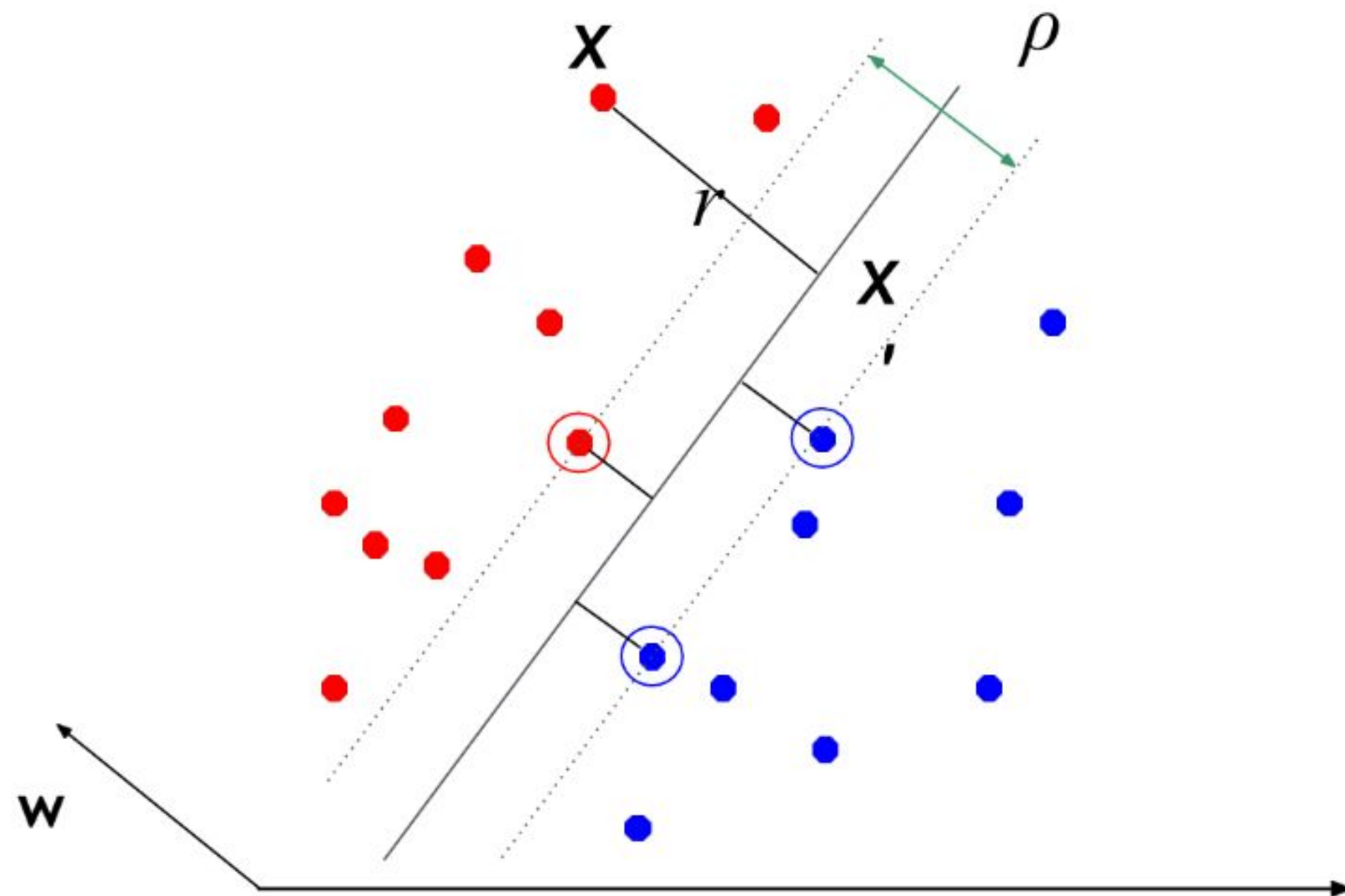


Максимизирует «зазор»



Максимальный зазор

- w – нормаль к разделяющей плоскости
- x_i - sample
- y_i : - класс sample i (+1 or -1)
(важно, не 1 и 0)
- Классификатор: $f(x_i) = \text{sign}(w^T x_i + b)$
- Зазор для точки x $r = y \frac{w^T x + b}{\|w\|}$
- Зазор всего датасета – минимум зазора для всех точек



Формула

Итого получаем задачу оптимизации:

Найти \mathbf{w} и b такие что

максимально; и для всех $\{(\mathbf{x}_i, y_i)\}$

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \text{ если } y_i = 1; \quad \mathbf{w}^T \mathbf{x}_i + b \leq -1 \text{ если } y_i = -1$$

Перепишем в более понятном виде:

Найти \mathbf{w} и b такие что

$$\Phi(\mathbf{w}) = 0.5 \mathbf{w}^T \mathbf{w} \text{ максимально}$$

$$\text{И для всех } \{(\mathbf{x}_i, y_i)\}: \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$



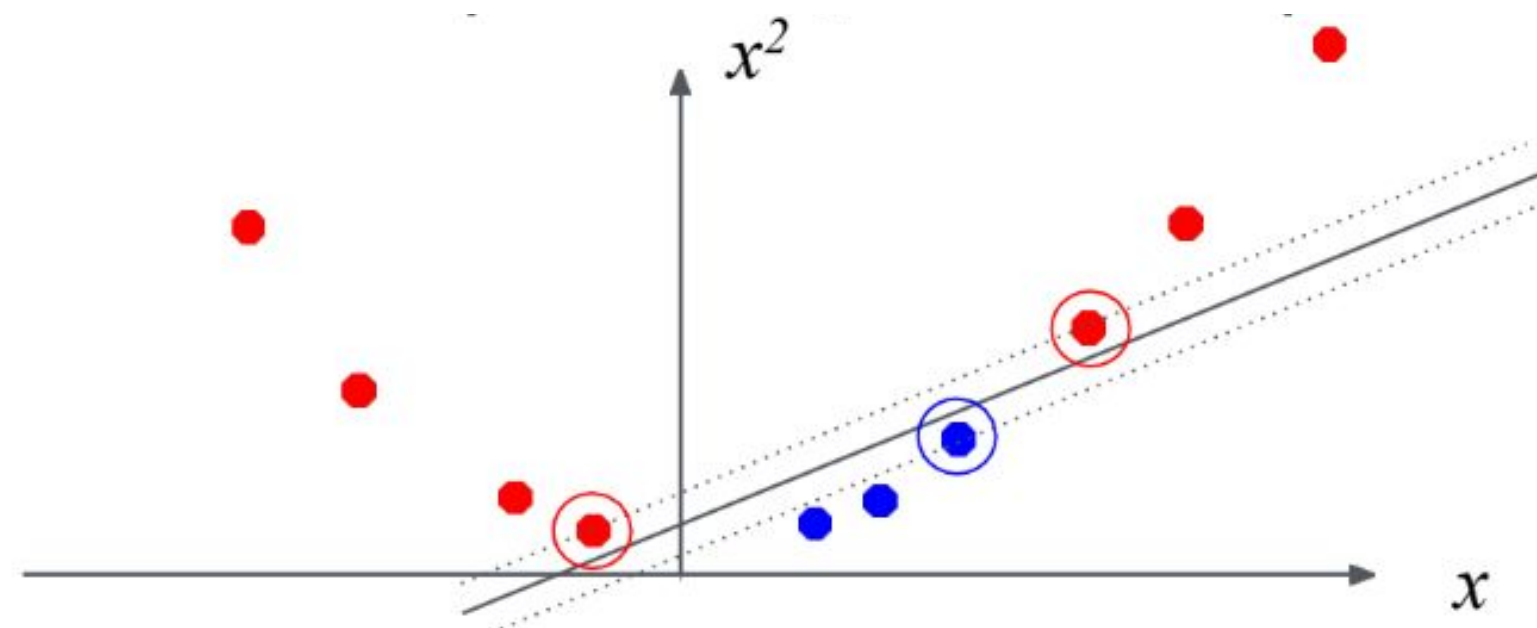
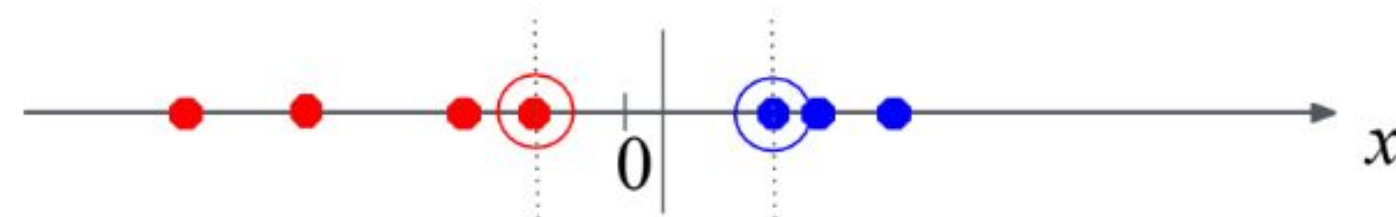
Non-linear SVMs



Линейно разделимые датасеты хорошо классифицируются

Но что делать, если они не линейно разделимы?

Можно попробовать отобразить данные в пр-во более высокой размерности



The «Kernel Trick»

- SVM зависит от скалярного произведения $K(x_i, x_j) = x_i^T x_j$
- Если каждая точка отображается в пр-во более высокой размерности при помощи $\Phi: x \rightarrow \phi(x)$, тогда скалярное произведение становится:
- $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$
- Функция ядра - это функция, соответствующая скалярному произведению в пр-ве более высокой размерности



Kernels

- Примеры
- Линейное
- Полиноминое $K(x,z) = (1+x^Tz)^d$
- RBF

$$K(x_i, x_j) = e^{-\|x_i - x_j\|^2 / 2\sigma^2}$$



Что мы сегодня узнали

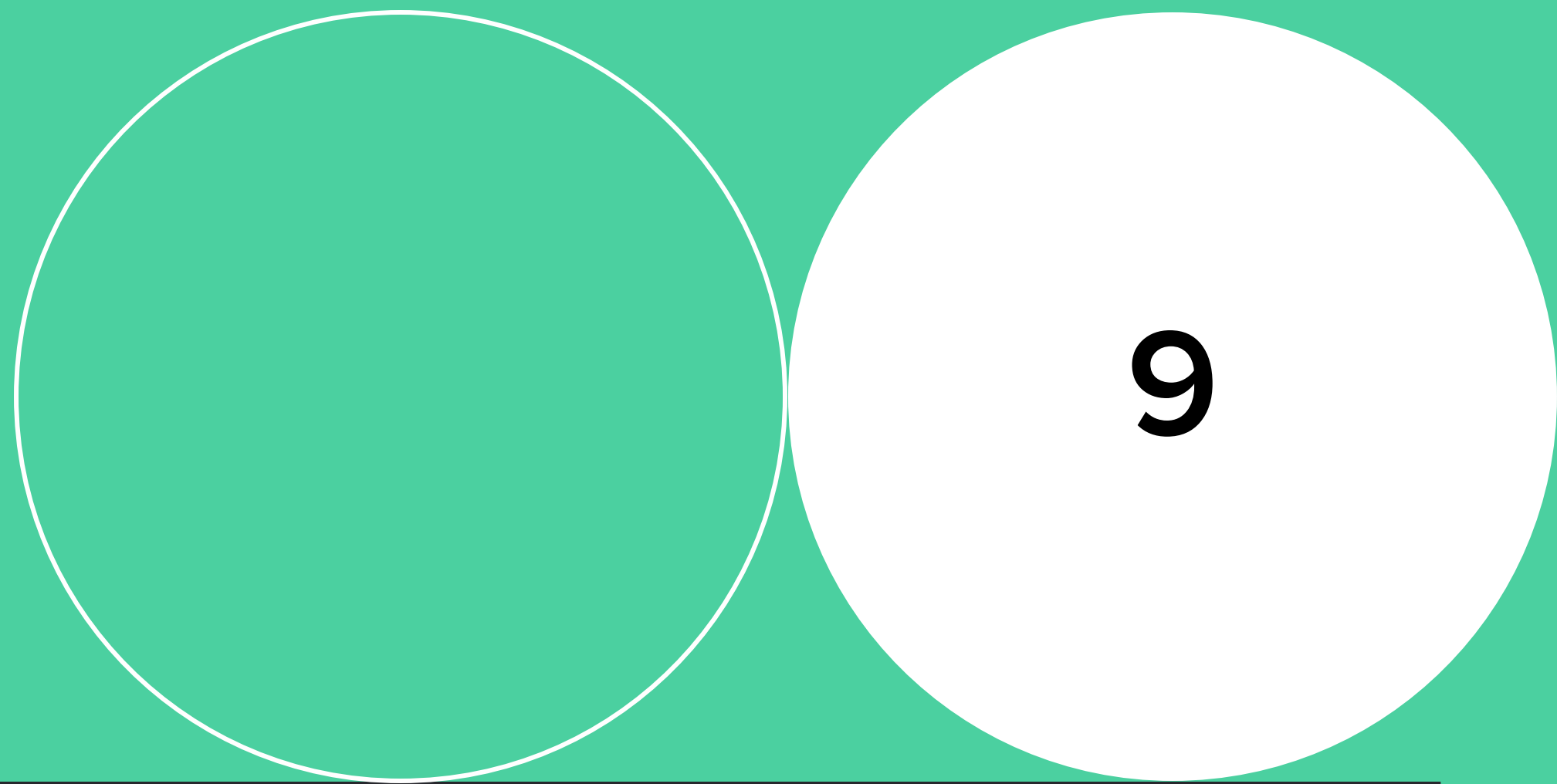


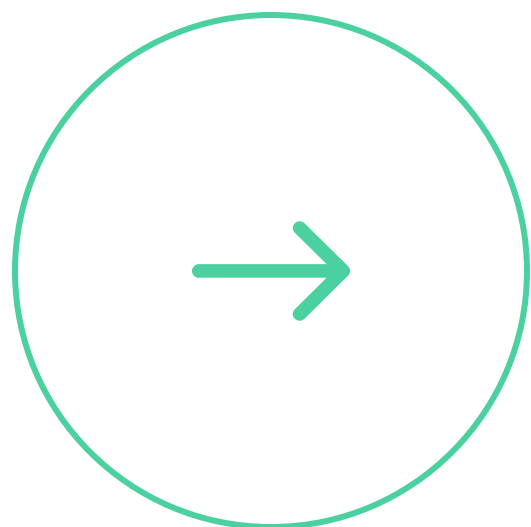
Итоги занятия

- Вспомнили основы теории вероятностей
- Изучили линейные модели и требования к ним на основе функции правдоподобия
- Реализовали логистическую регрессию
- Изучили алгоритм градиентного спуска и потренировались в его реализации



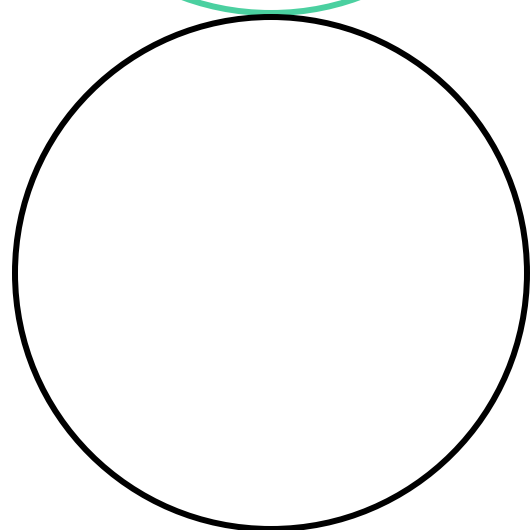
Полезные материалы





Статья о линейных моделях в ODS

<https://habrahabr.ru/company/ods/blog/323890/>



Курс «Основы статистики» на Stepik.org

<https://stepik.org/course/Основы-статистики-76>



Спасибо за внимание!

