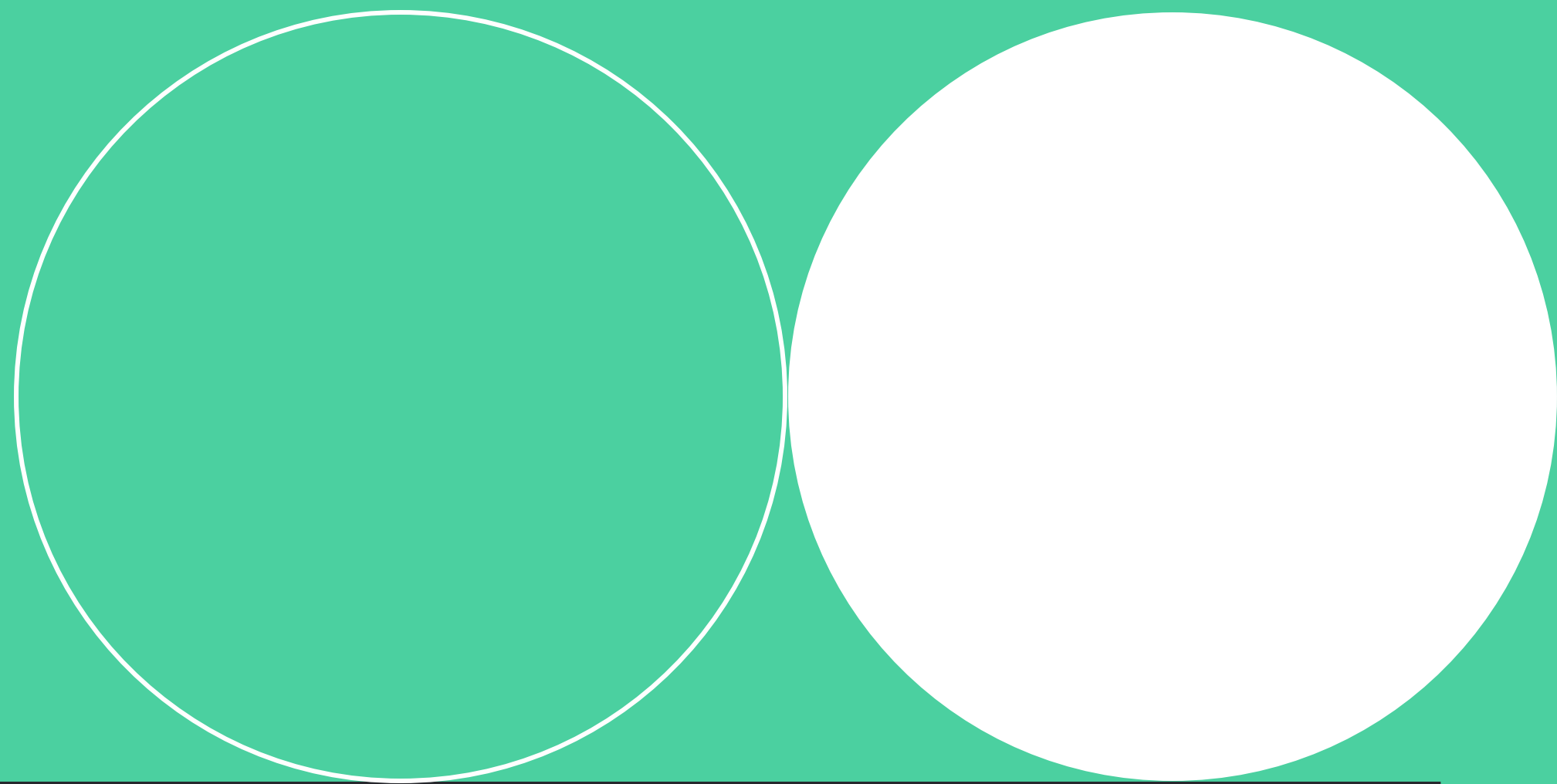

Линейный классификатор и логистическая регрессия



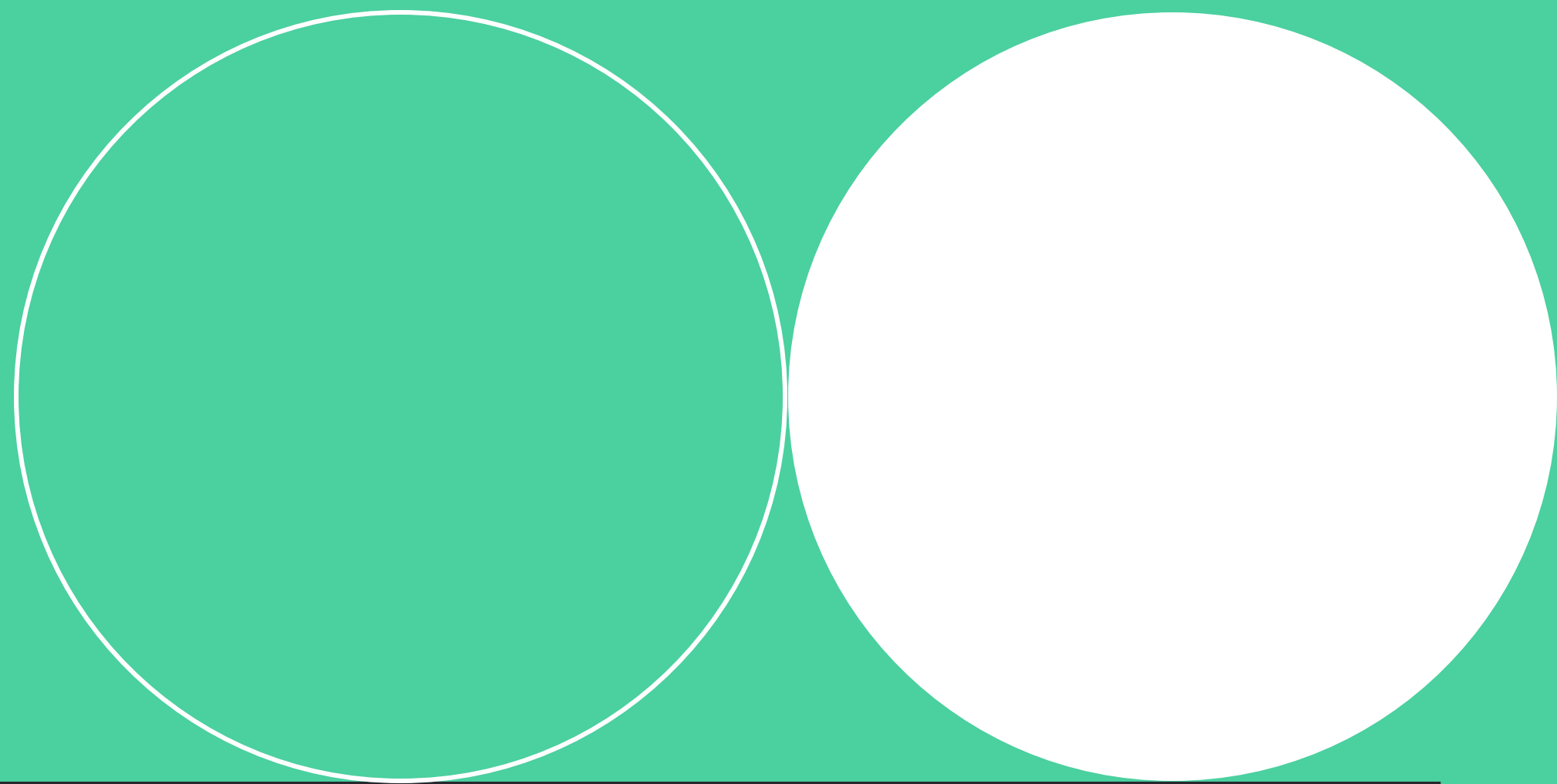
Цели занятия



- Рассмотреть задачу классификации
- Рассмотреть линейный классификатор
- Рассмотреть варианты линейного классификатора.
Логистическую регрессию и SVM.



О чем поговорим и что сделаем



План занятия

1

Задача классификации

2

Логистическая регрессия: практическое задание

3

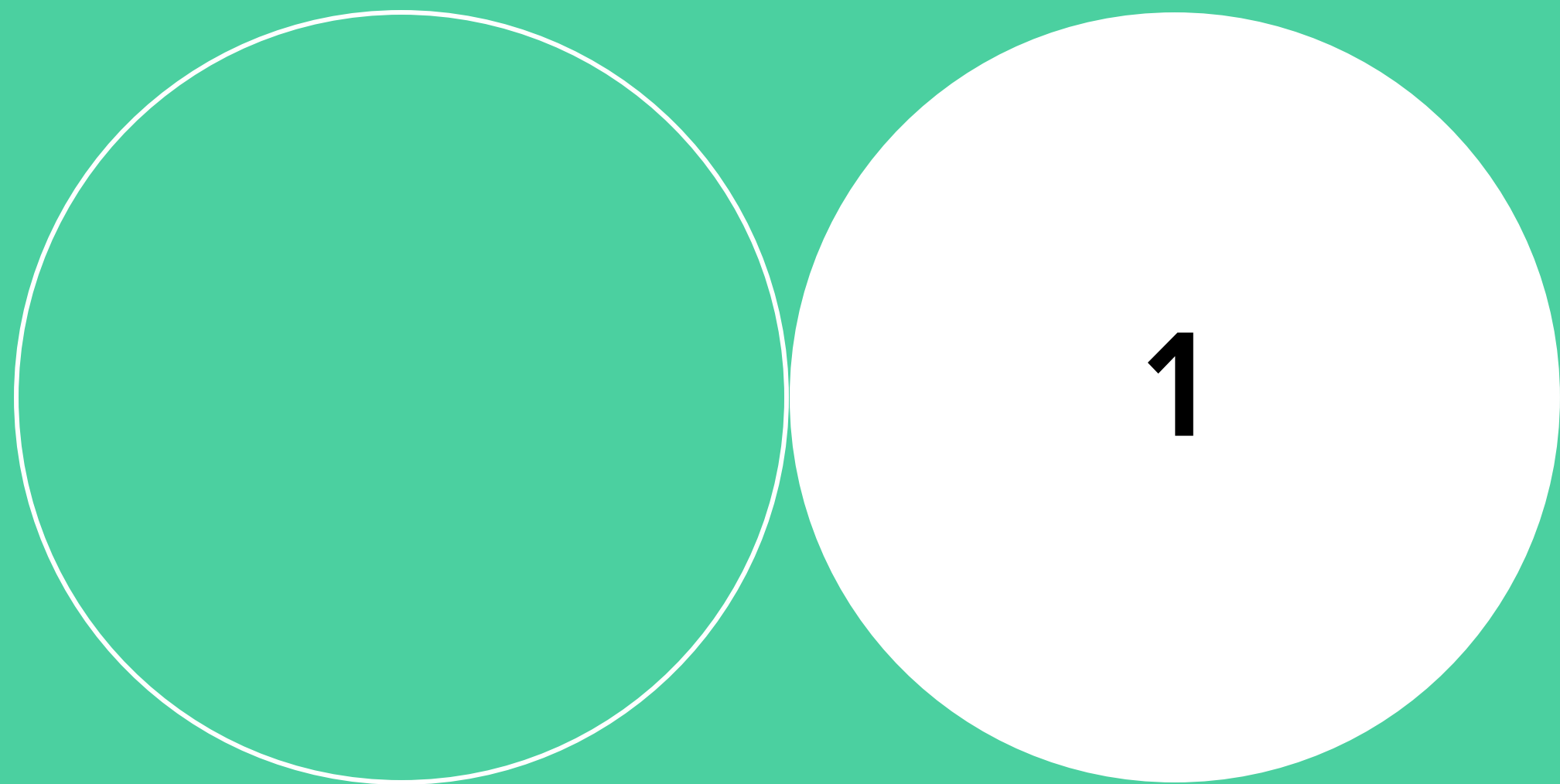
SVM

4

Практика

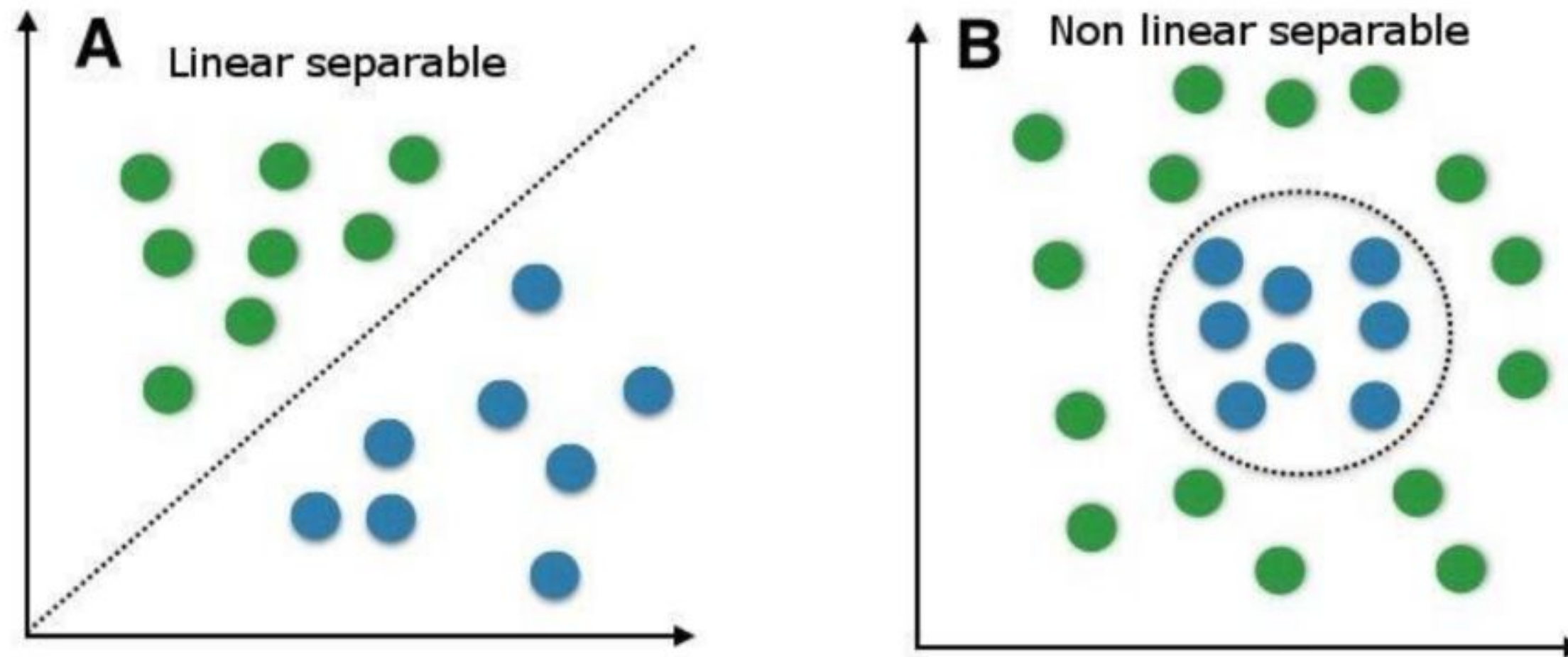


Классификация



Классификации – задача *предсказания ответа* из **конечного** множества вариантов.

Линейный классификатор - решает задачу разделения признакового пространства на две части, в каждом из которых находится свой класс.



Линейная делимость данных

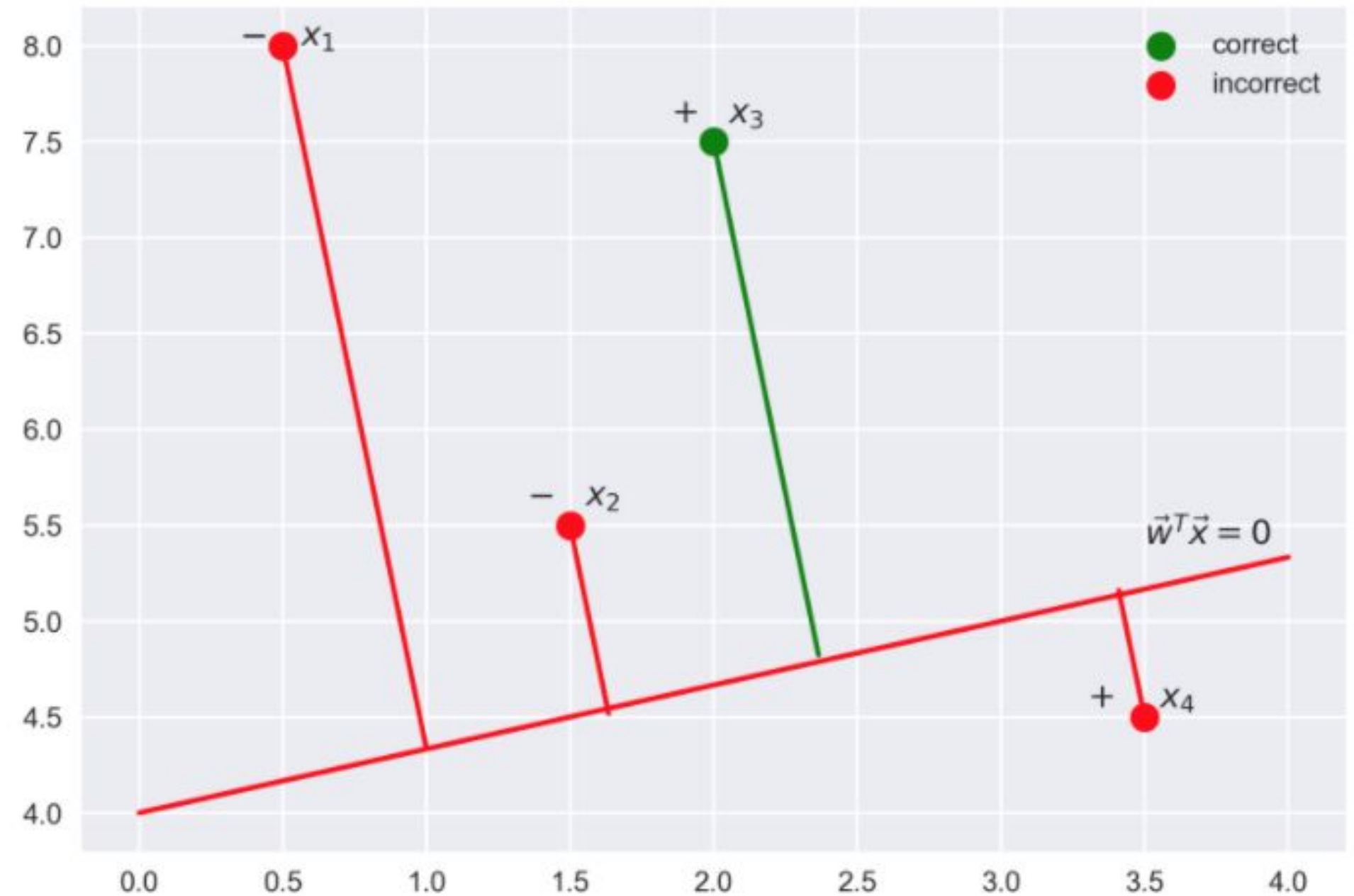


Отступ - расстояние от разделяющей гиперповерхности до объекта

$$M(\vec{x}_i) = y_i \vec{w}^T \vec{x}_i$$

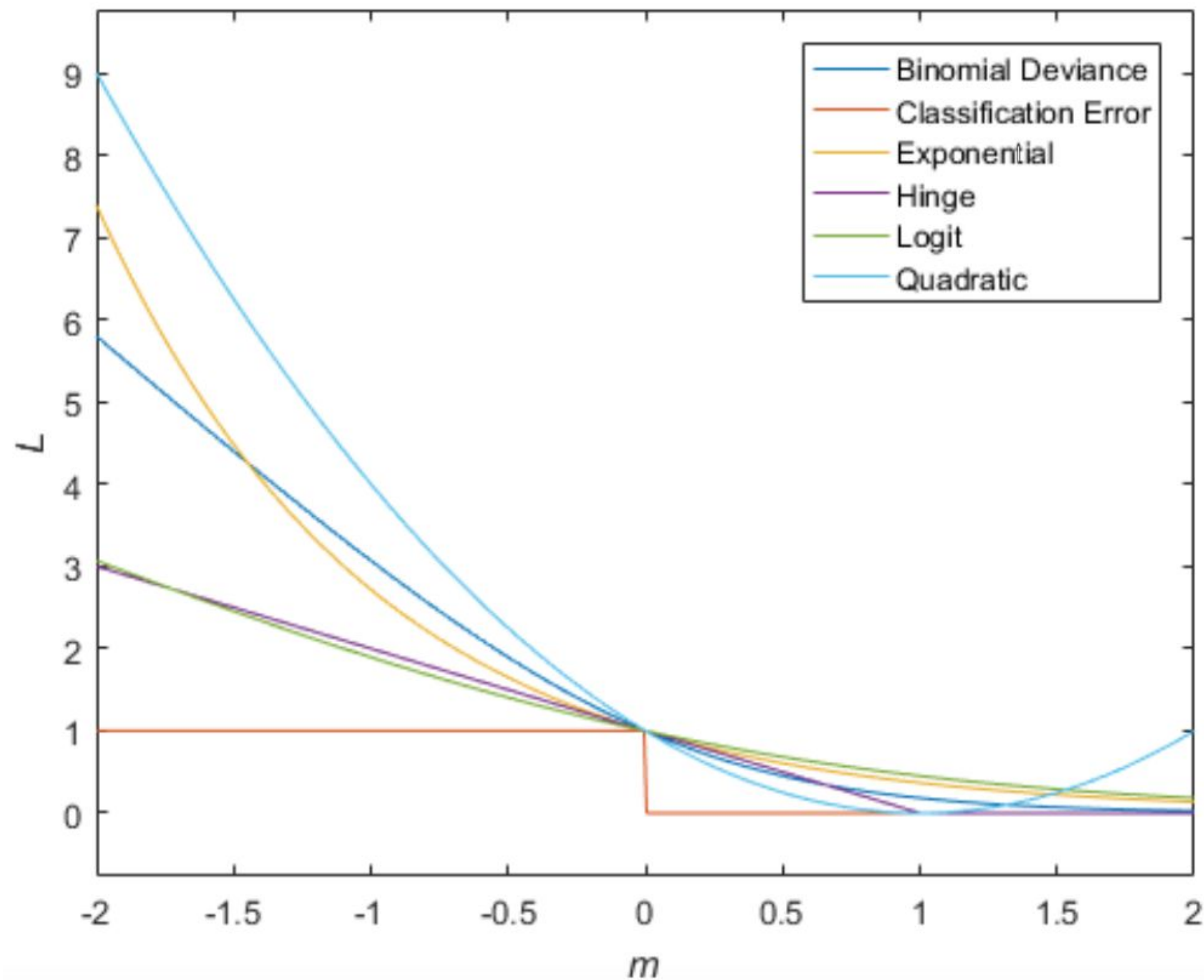
Отступ можно понимать как «степень погруженности» объекта в свой класс.

<https://habr.com/ru/company/ods/blog/323890/>
<https://dyakonov.org/логистическая-функция-ошибки/#more-6139>



Функция потерь L как функция от отступа M

$$L = f(M)$$



Логистическая регрессия



4

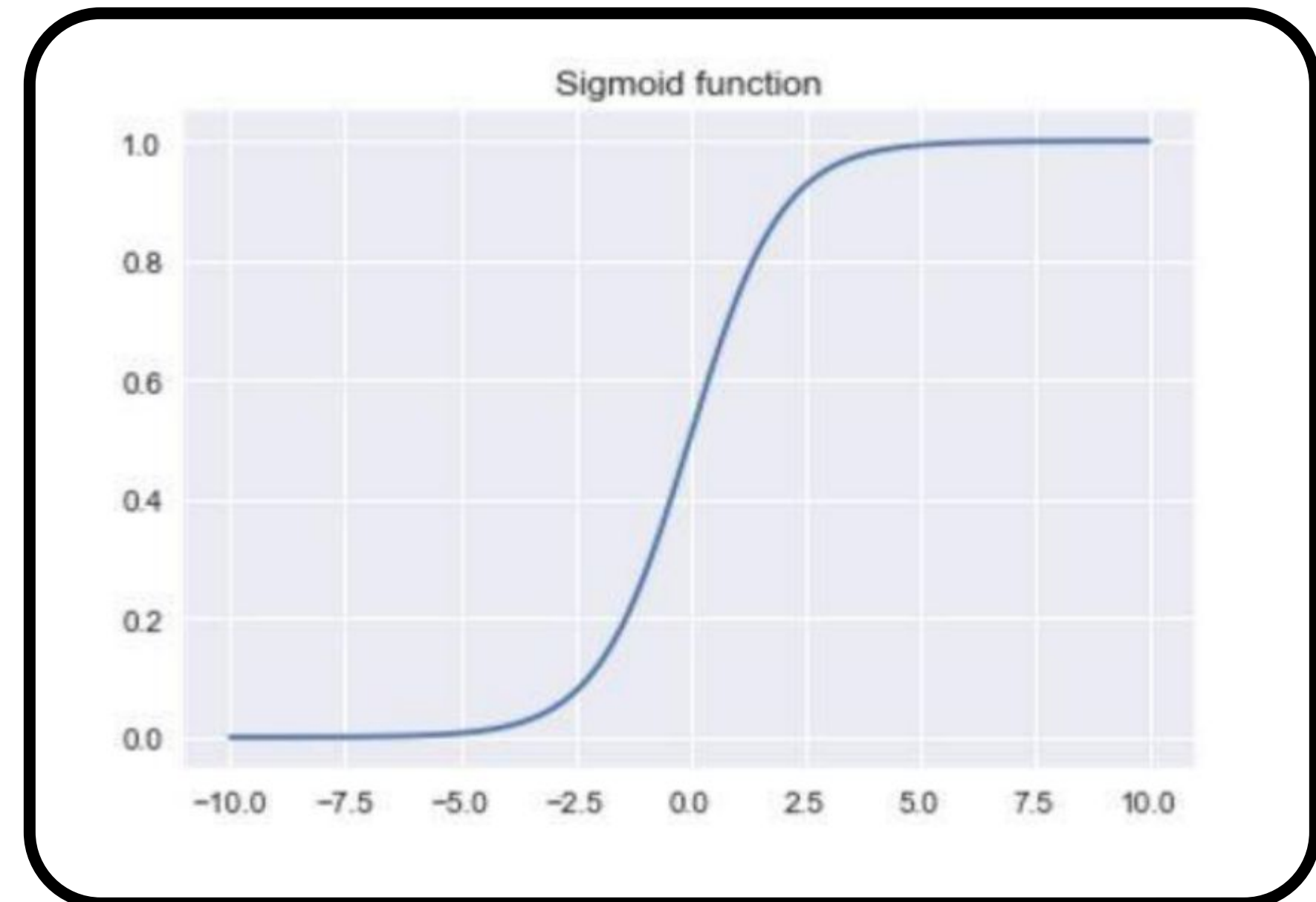
Логистическая регрессия

Логистическая регрессия - линейный классификатор, позволяющий **оценивать вероятности** принадлежности объектов классам.

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

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

<https://habr.com/ru/post/485872/>
<https://habr.com/ru/company/ods/blog/323890/>
<https://dyakonov.org/логистическая-функция-ошибки/#more-6139>



Функция потерь

Модель предсказывает
вероятность классов $\{0, +1\}$

$$p(y_i | x_i, w) = a_i^{y_i} (1 - a_i)^{1 - y_i}$$

Максимизировать
правдоподобие

$$p(y | X, w) = \prod_i p(y_i | x_i, w)$$

Функция потерь

$$\mathcal{L}_{\log}(X, \vec{y}, \vec{w}) = \sum_i (-y_i \log a_i - (1 - y_i) \log(1 - a_i))$$



Функция потерь

Модель предсказывает
вероятность классов $\{-1, +1\}$

$$P(y = y_i \mid \vec{x}_i, \vec{w}) = \sigma(y_i \vec{w}^T \vec{x}_i)$$

Максимизировать
правдоподобие

$$P(\vec{y} \mid X, \vec{w}) = \prod_{i=1}^{\ell} P(y = y_i \mid \vec{x}_i, \vec{w})$$

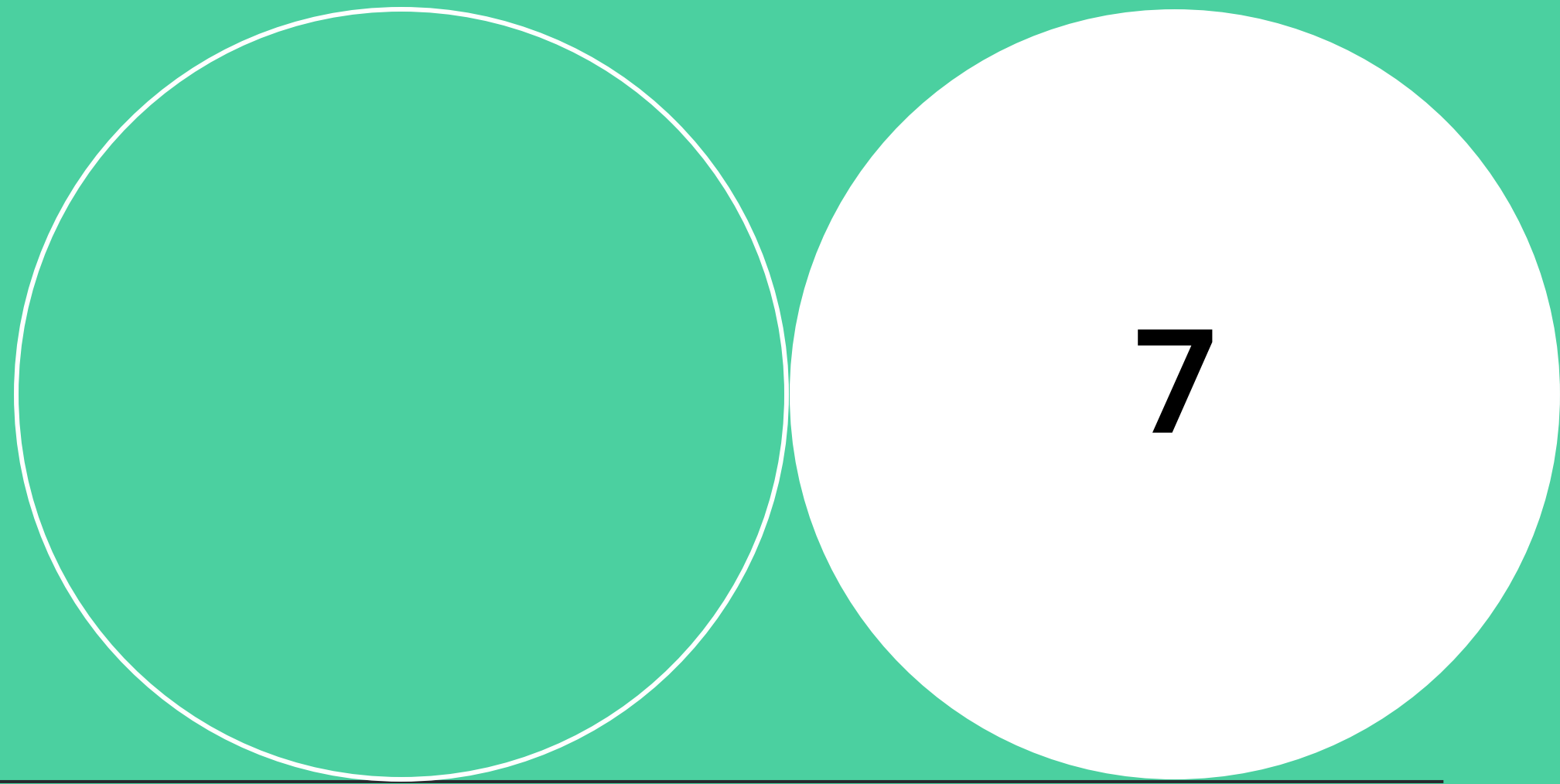
$$\log P(\vec{y} \mid X, \vec{w}) = - \sum_{i=1}^{\ell} \log(1 + \exp^{-y_i \vec{w}^T \vec{x}_i})$$

Функция потерь

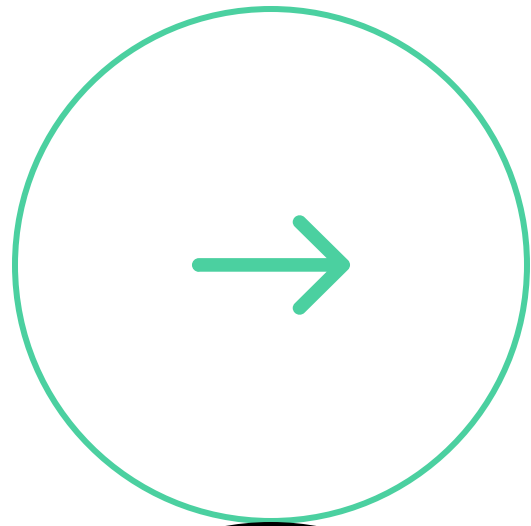
$$\mathcal{L}_{\log}(X, \vec{y}, \vec{w}) = \sum_{i=1}^{\ell} \log(1 + \exp^{-y_i \vec{w}^T \vec{x}_i})$$



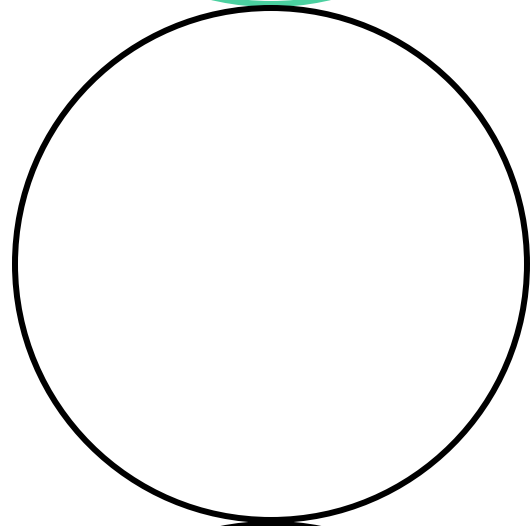
SVM



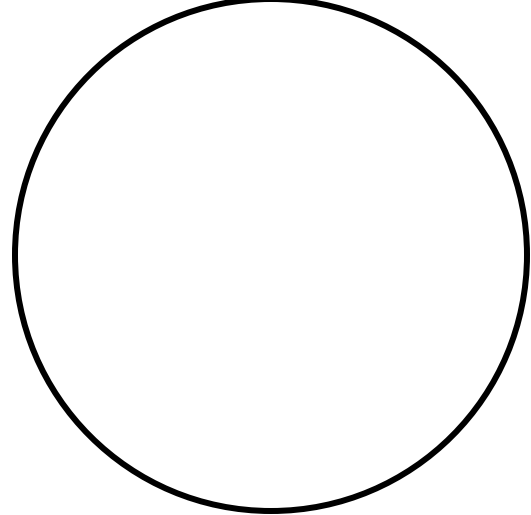
Множество гиперплоскостей



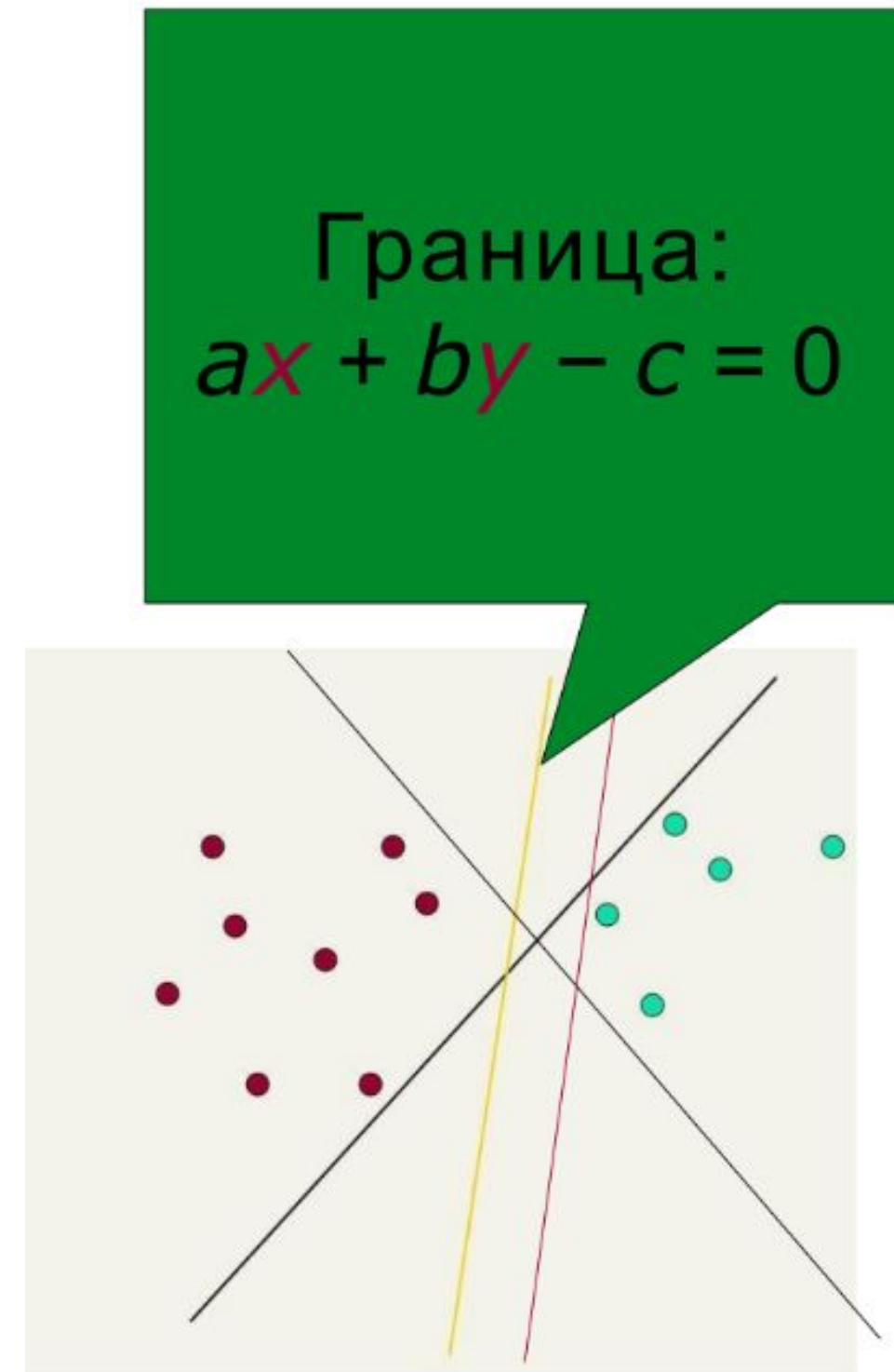
Множество решений для a, b, c .



SVM находит оптимальную разделяющую поверхность

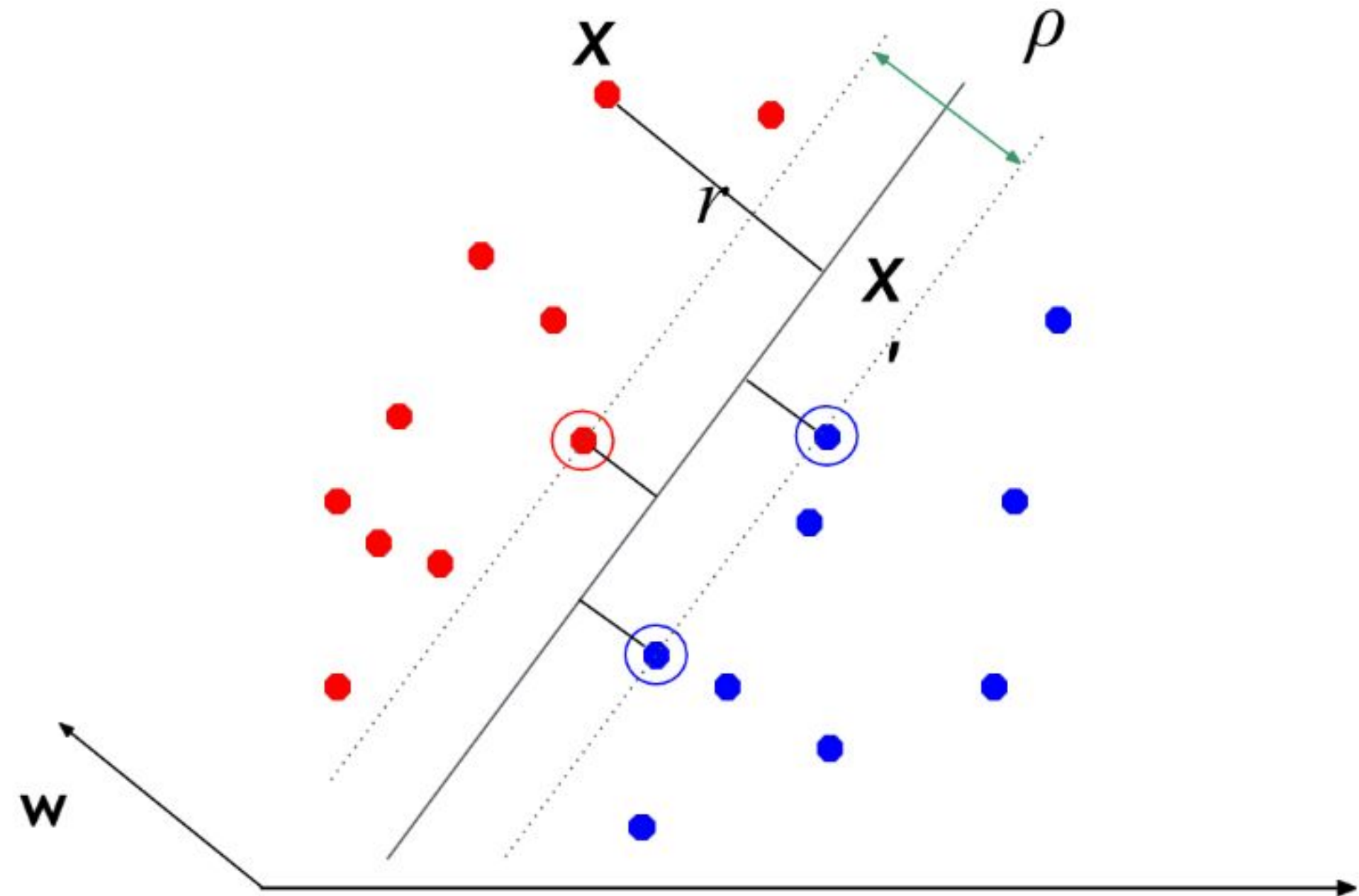


Максимизирует «зазор»



Максимальный зазор

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

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

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



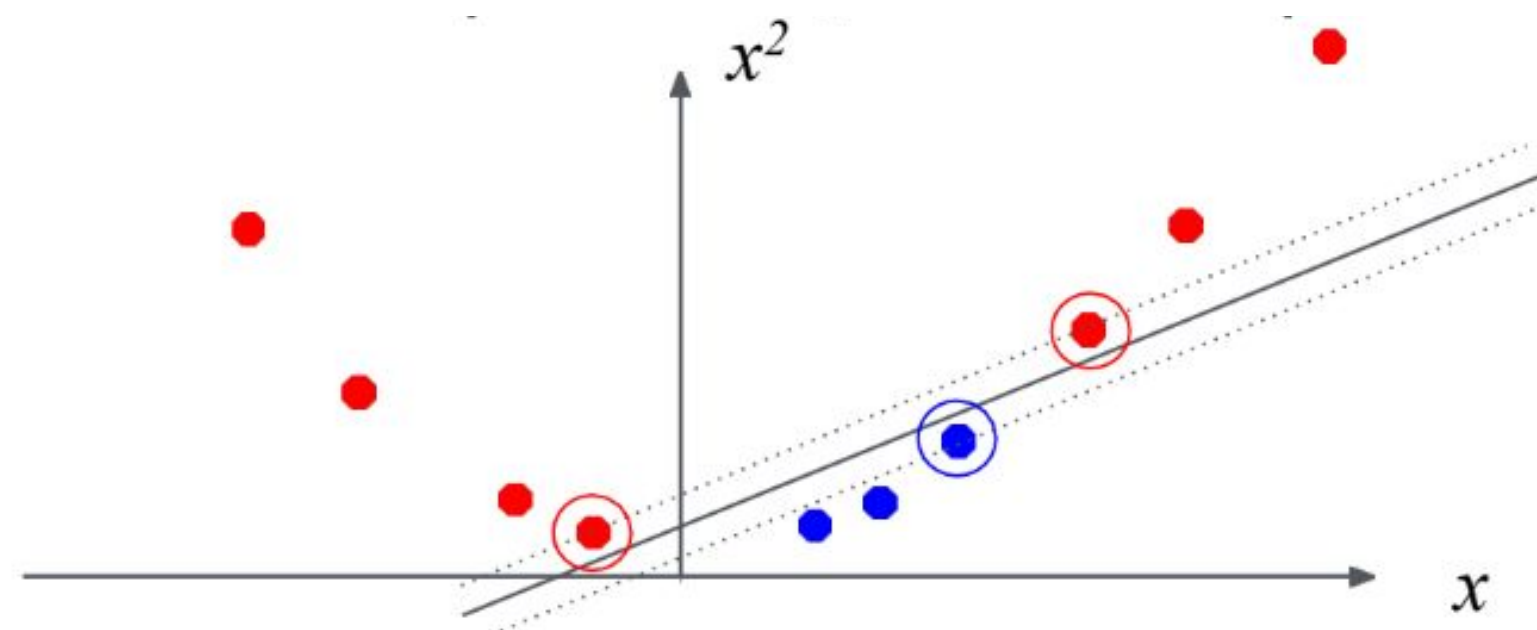
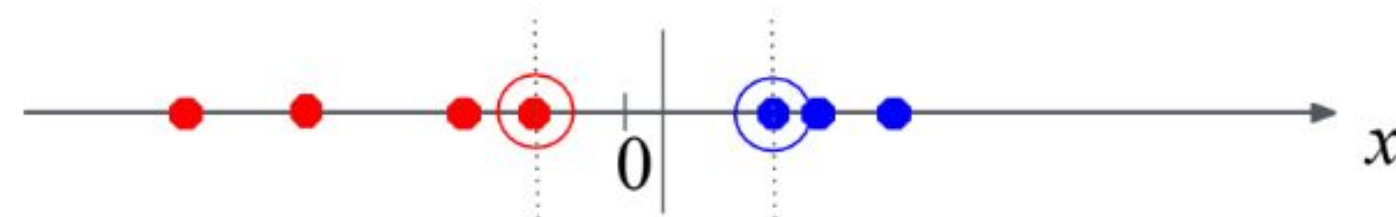
Non-linear SVMs



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

- Полиномиальное
- Полиномиальное со смещением

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}')^d$$

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + 1)^d$$



Kernels

- Радиальная базисная функция $k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$, для $\gamma > 0$

- Радиальная базисная функция Гаусса

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$



Практика

Что мы сегодня узнали



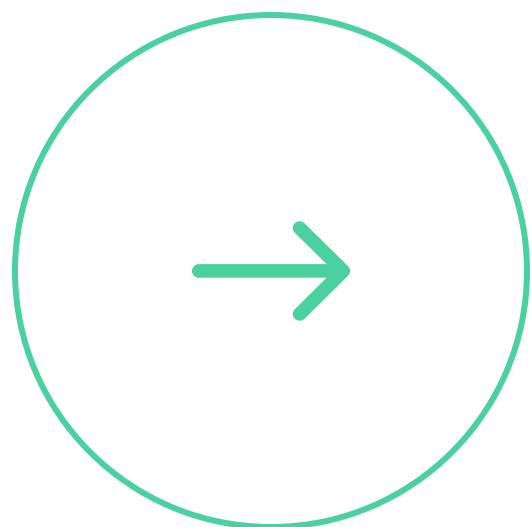
Итоги занятия

- 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



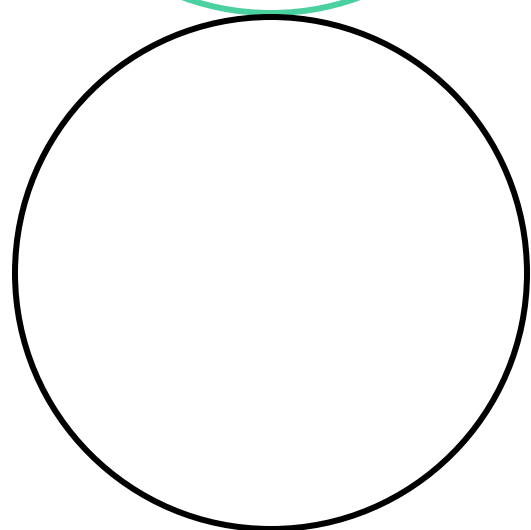
Полезные материалы





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!

