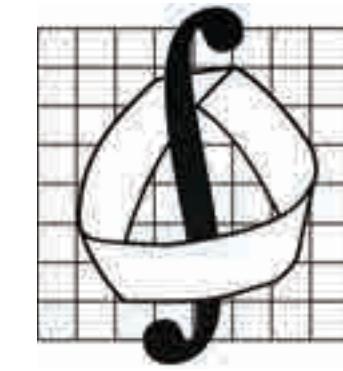


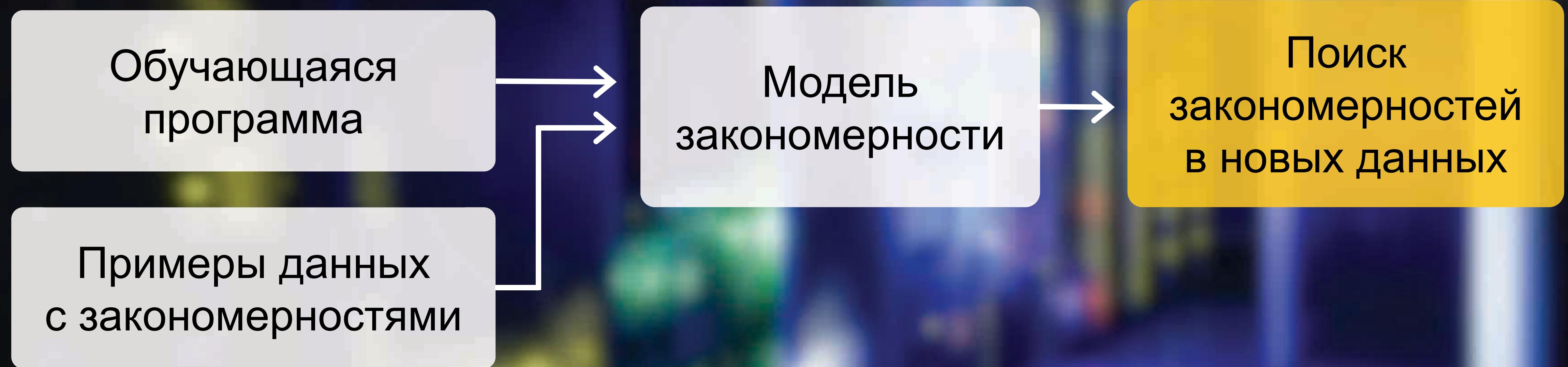
Яндекс



Машинное обучение и математика: что общего?

Елена Игоревна Бунина, профессор механико-математического факультета МГУ, директор отделения компьютерных наук Школы анализа данных Яндекса

Главная идея машинного обучения



Основные понятия машинного обучения

Дано:

X

множество объектов (заданное своими признаками,
чаще всего как точки n -мерного пространства)

Y

множество ответов (чаще всего $\{0, 1\}$ или $\{-1, 1\}$,
или $\{1, 2, \dots, n\}$ или \mathbb{R})

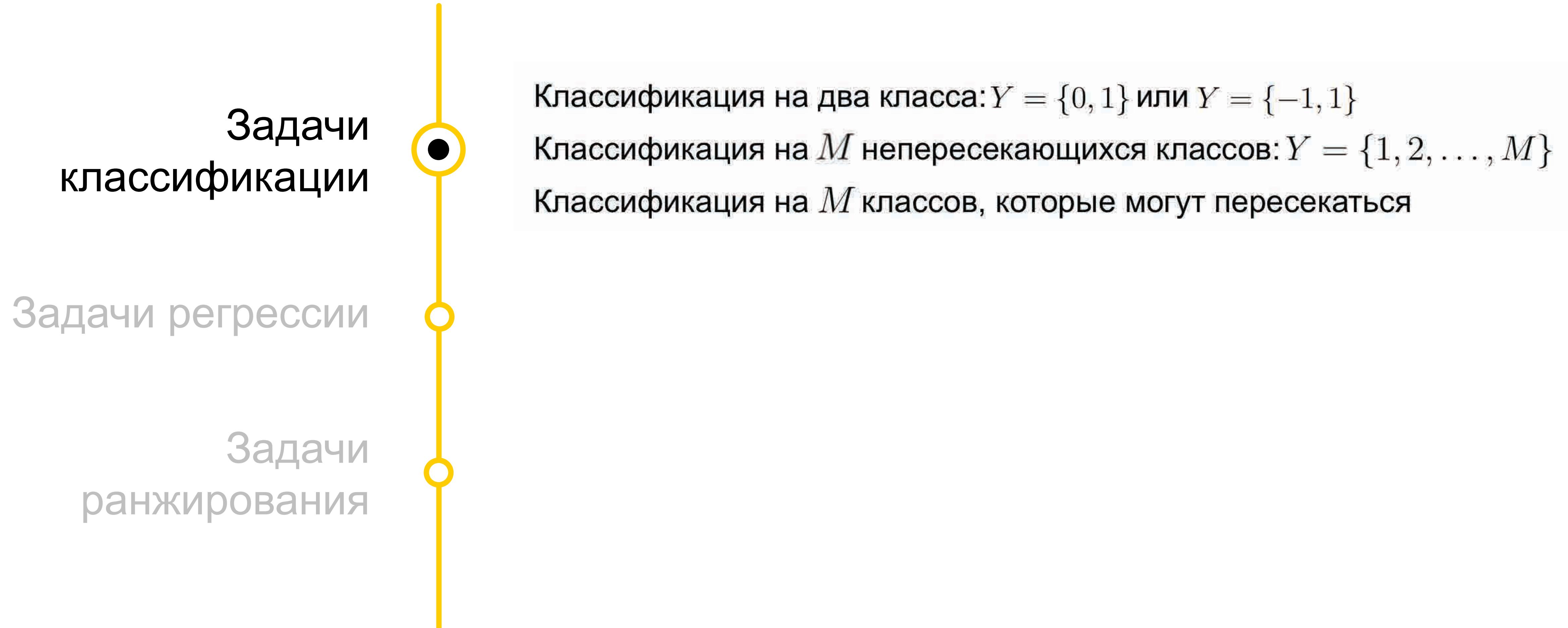
множество $\{x_1, x_2, \dots, x_l\} \subset X$ – обучающая
выборка

$y_i = y(x_i) \in Y, i = 1, \dots, l,$ – известные ответы

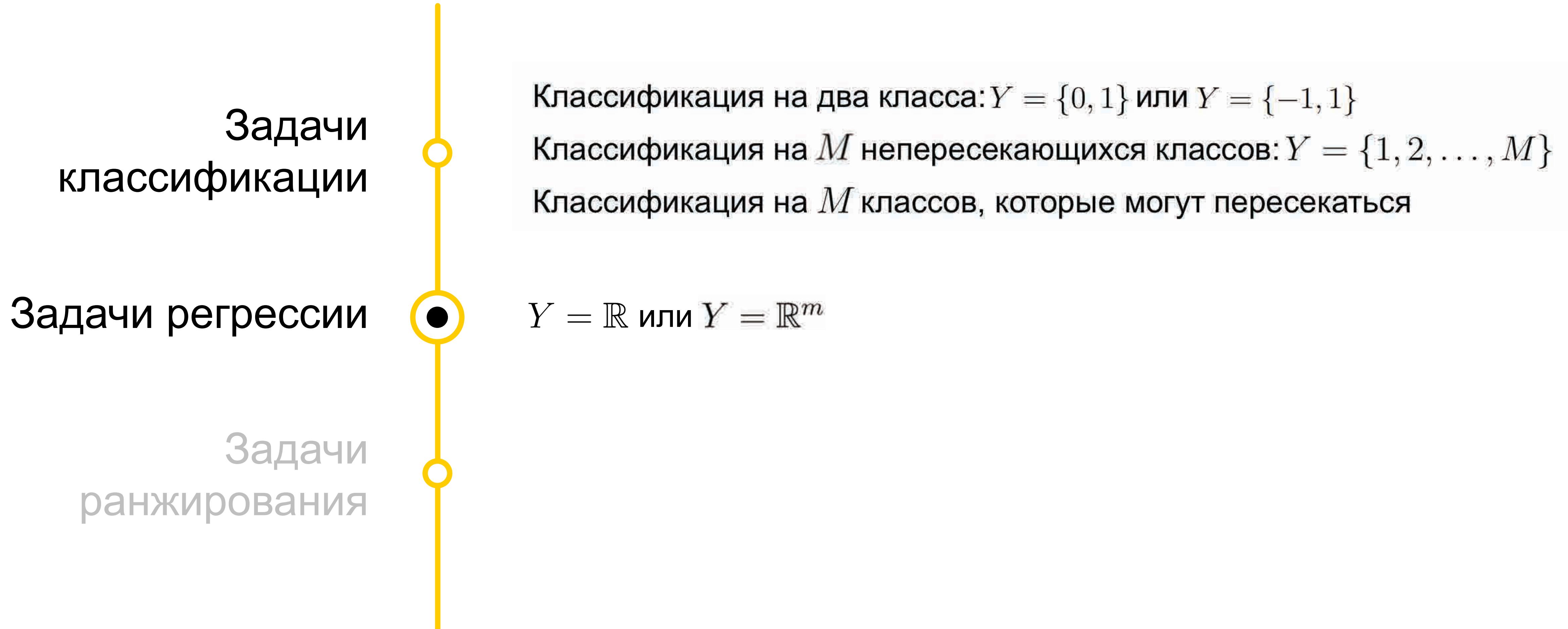
Основные понятия машинного обучения

Найти: $a : X \rightarrow Y$ – алгоритм, решающую функцию, приближающую неизвестную функцию y на всём множестве X

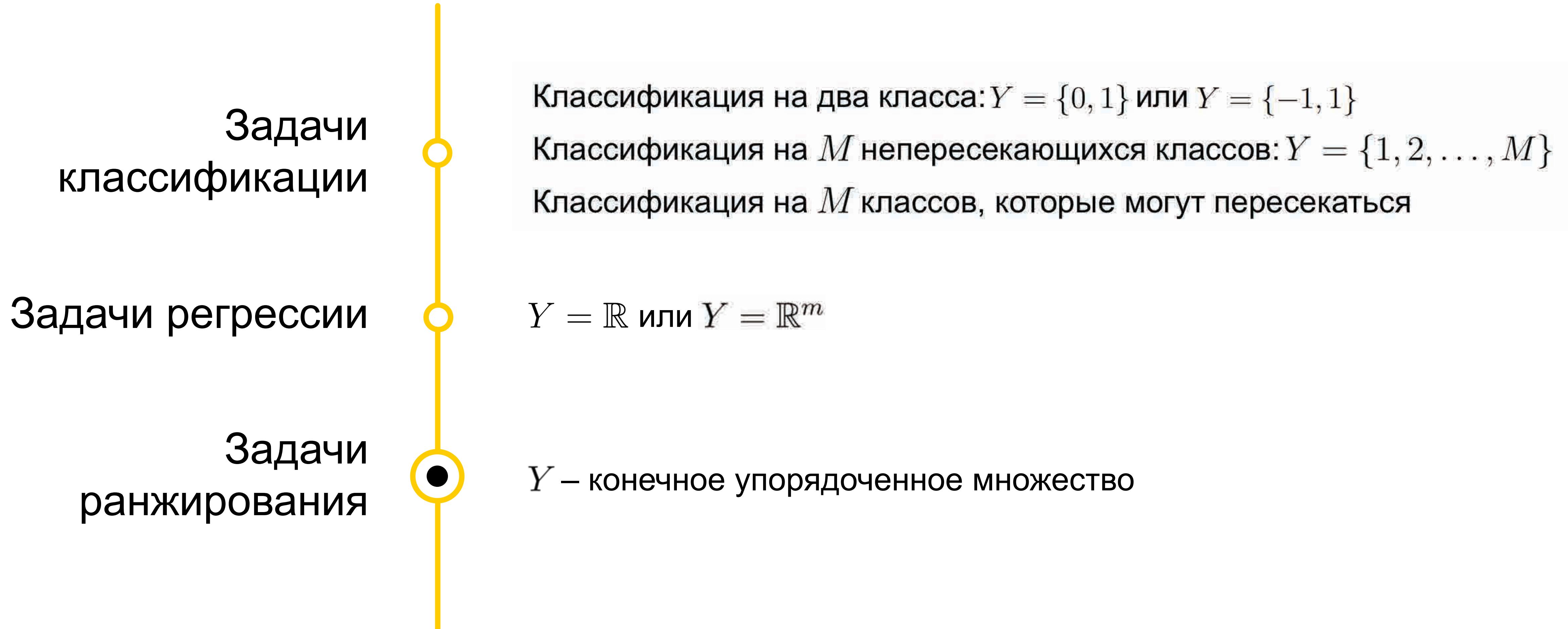
Типы задач машинного обучения



Типы задач машинного обучения



Типы задач машинного обучения



Модель машинного обучения

Модель
(предсказательная модель)
это параметрическое
семейство функций

$$A = \{g(x, \theta) \mid \theta \in \Theta\}$$

Где: $g : X \times \Theta \rightarrow Y$ – фиксированная функция
 Θ – множество допустимых значений параметра θ

Пример: линейная модель

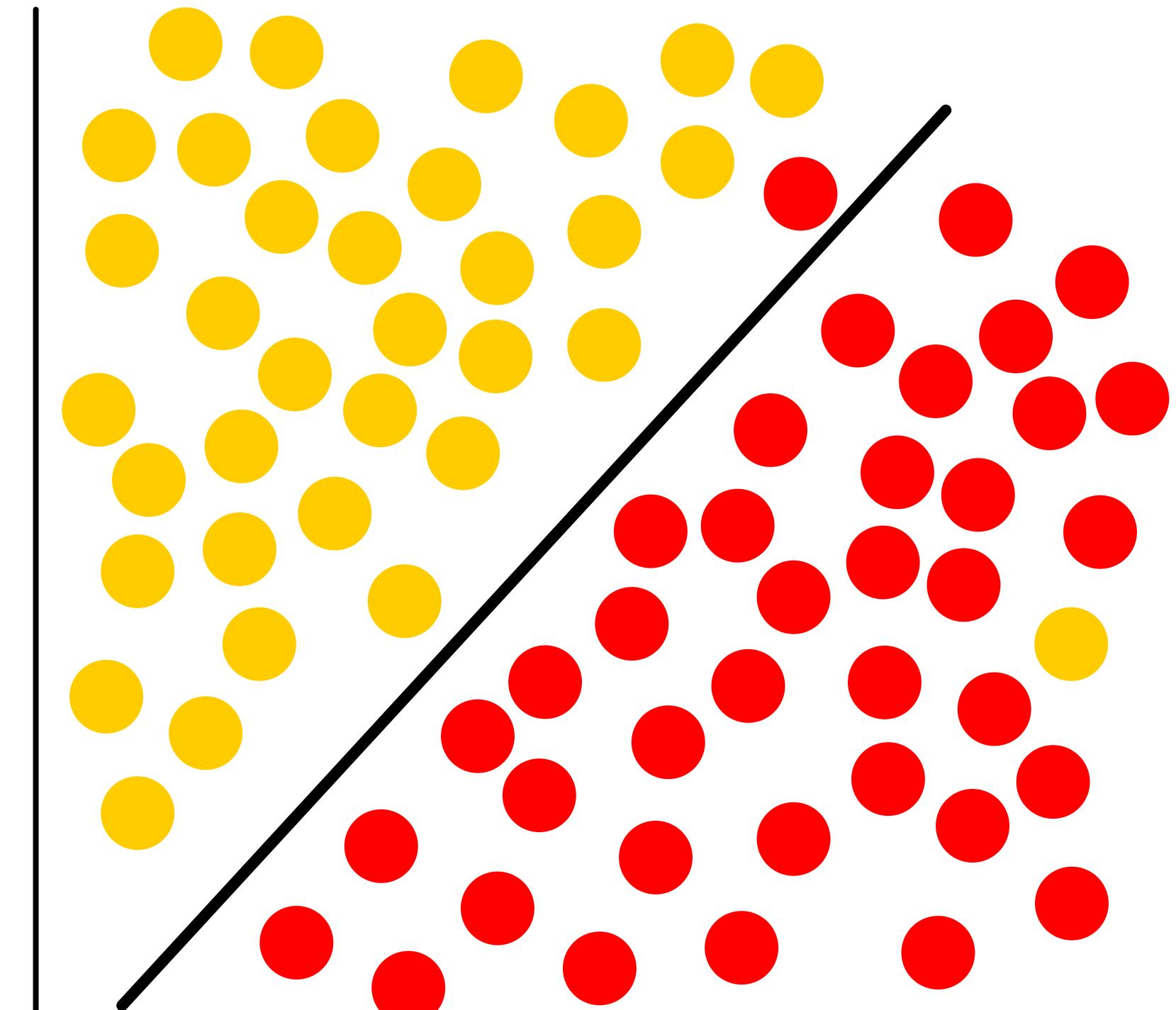
Пусть вектор параметров – это $\theta = (\theta_1, \dots, \theta_n)$, $\Theta = \mathbb{R}^n$.

Тогда для задач регрессии и ранжирования можно положить

$$g(x, \theta) = \sum_{j=1}^n \theta_j f_j(x), \quad Y = \mathbb{R}$$

Для задачи классификации обычно полагают

$$g(x, \theta) = \operatorname{sgn} \sum_{j=1}^n \theta_j f_j(x), \quad Y = \{-1, 1\}$$



Этапы обучения

Метод обучения

это отображение вида
 $\mu : (X \times Y)^l \rightarrow A$,

которое произвольной
выборке

$\{(x_i, y_i) \mid i = 1, \dots, l\}$

ставит в соответствие
алгоритм $a \in A$

В задачах обучения всегда есть два этапа:

1. Этап обучения

метод μ по выборке $\bar{x} \in X^l$
строит алгоритм обучения $a = \mu(\bar{x})$

2. Этап применения

алгоритм a для новых объектов $x \in X$
выдаёт ответы $a(x)$

ФУНКЦИЯ ПОТЕРЬ

Функция потерь $\mathcal{L}(a, x)$ – это величина ошибки алгоритма
 $a \in A$ на объекте $x \in X$

Функция потерь для задач классификации:

$\mathcal{L}(a, x) = [a(x) \neq y(x)]$ – индикатор ошибки

Функции потерь для задач регрессии:

$\mathcal{L}(a, x) = |a(x) - y(x)|$ – абсолютное значение ошибки

или

$\mathcal{L}(a, x) = (a(x) - y(x))^2$ – квадратичная ошибка

Минимизация эмпирического риска

Функционал качества – это эмпирический риск алгоритма a на X^l :

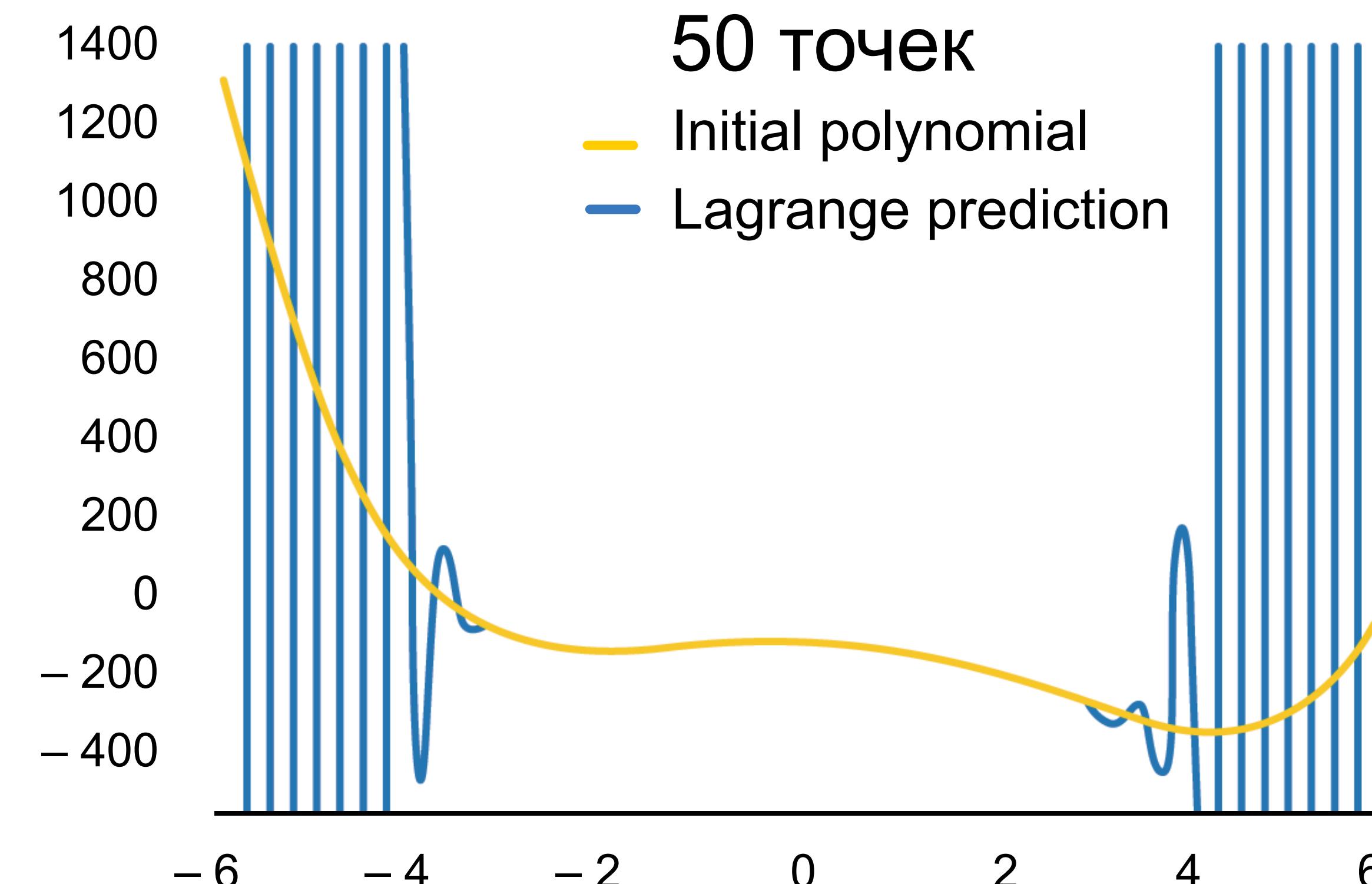
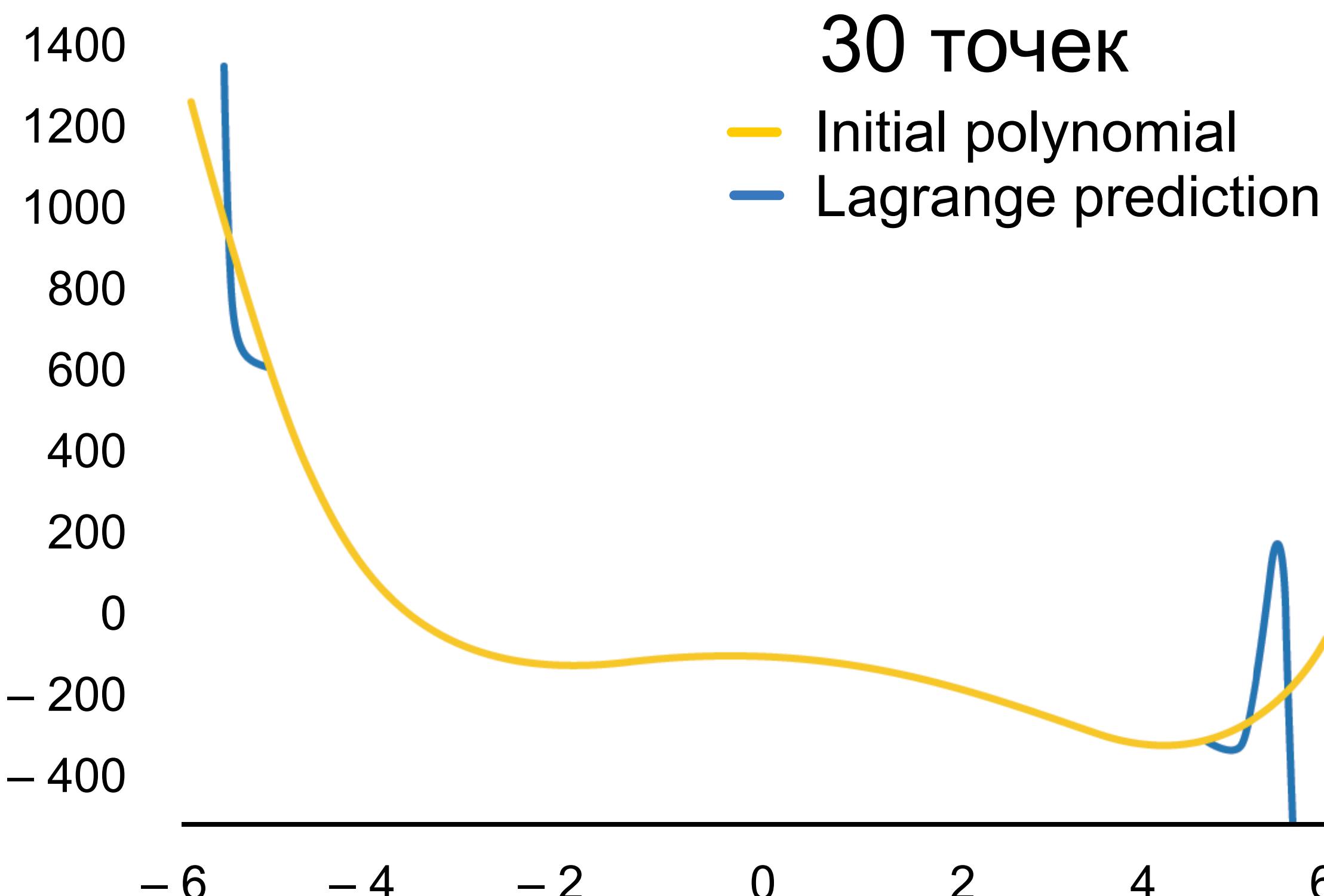
$$Q(a, X^l) = \frac{1}{l} \sum_{j=1}^l \mathcal{L}(a, x_j)$$

Вот каким образом мы сводим задачу обучения к обычной задаче оптимизации:

$$\mu(X^l) = \arg \min_{a \in A} Q(a, X^l)$$

Проблема переобучения

Рассмотрим полином $x^4 - 3x^3 - 17x^2 + 4$ на отрезке $[-6; 6]$. Случайный маленький шум
На новой функции возьмём n точек и просто проведём через них интерполяционный многочлен
Лагранжа



Проверка качества обучения

Как проверить, сильно ли вы переобучились?

разделить выборку на две части, обучаться на одной части, проверять на другой:

Кросс-проверка

$$Q(\mu(X^l), X^k) \rightarrow \min$$

Скользящий контроль (leave-one-out)

$$\frac{1}{L} \sum_{i=1}^L \mathcal{L}(\mu(X^L \setminus \{x_i\}), x_i) \rightarrow \min$$

$$X^L = X_n^l \cup X_n^k, \quad n = 1, \dots, N$$

$$\frac{1}{N} \sum_{n=1}^N Q(\mu(X_n^l), X_n^k) \rightarrow \min$$

Цикл решения задачи





Как применяется
Машинное обучение
в жизни

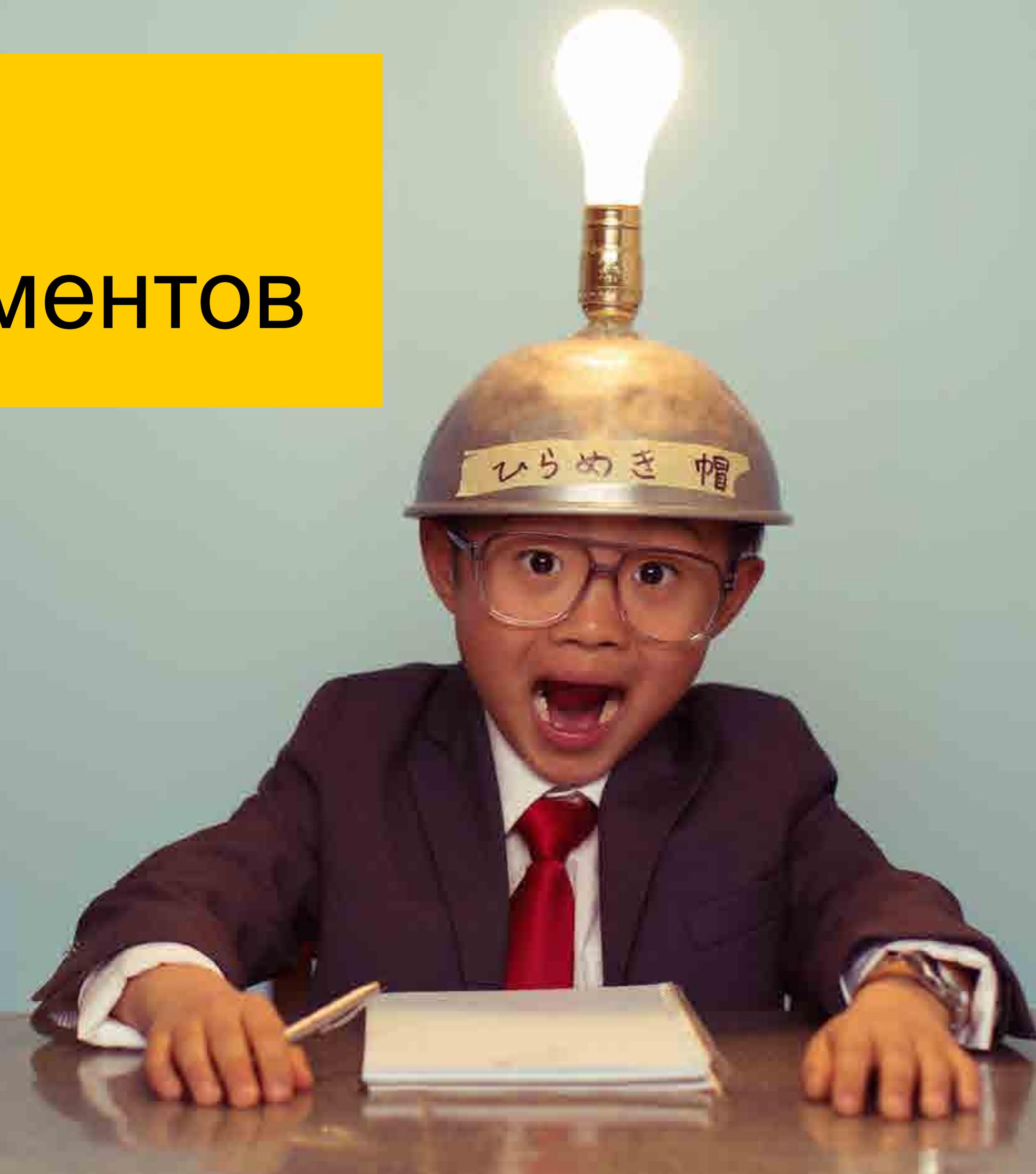
Медицина



Банковская сфера



Категоризация текстовых документов



Прогнозирование стоимости недвижимости



Ранжирование поисковой выдачи



Поиск изображений



Рекомендательная система



Метрические методы

Гипотезы компактности и непрерывности

Гипотеза компактности (для задач классификации):

Близкие объекты, как правило, лежат в одном классе

Гипотеза непрерывности (для задач регрессии):

Близким объектам соответствуют близкие ответы

Формализация понятия близости

Задана функция расстояния $\rho : X \times X \rightarrow [0, \infty)$

Например, любая метрика:

$$\rho(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \text{ — евклидова}$$

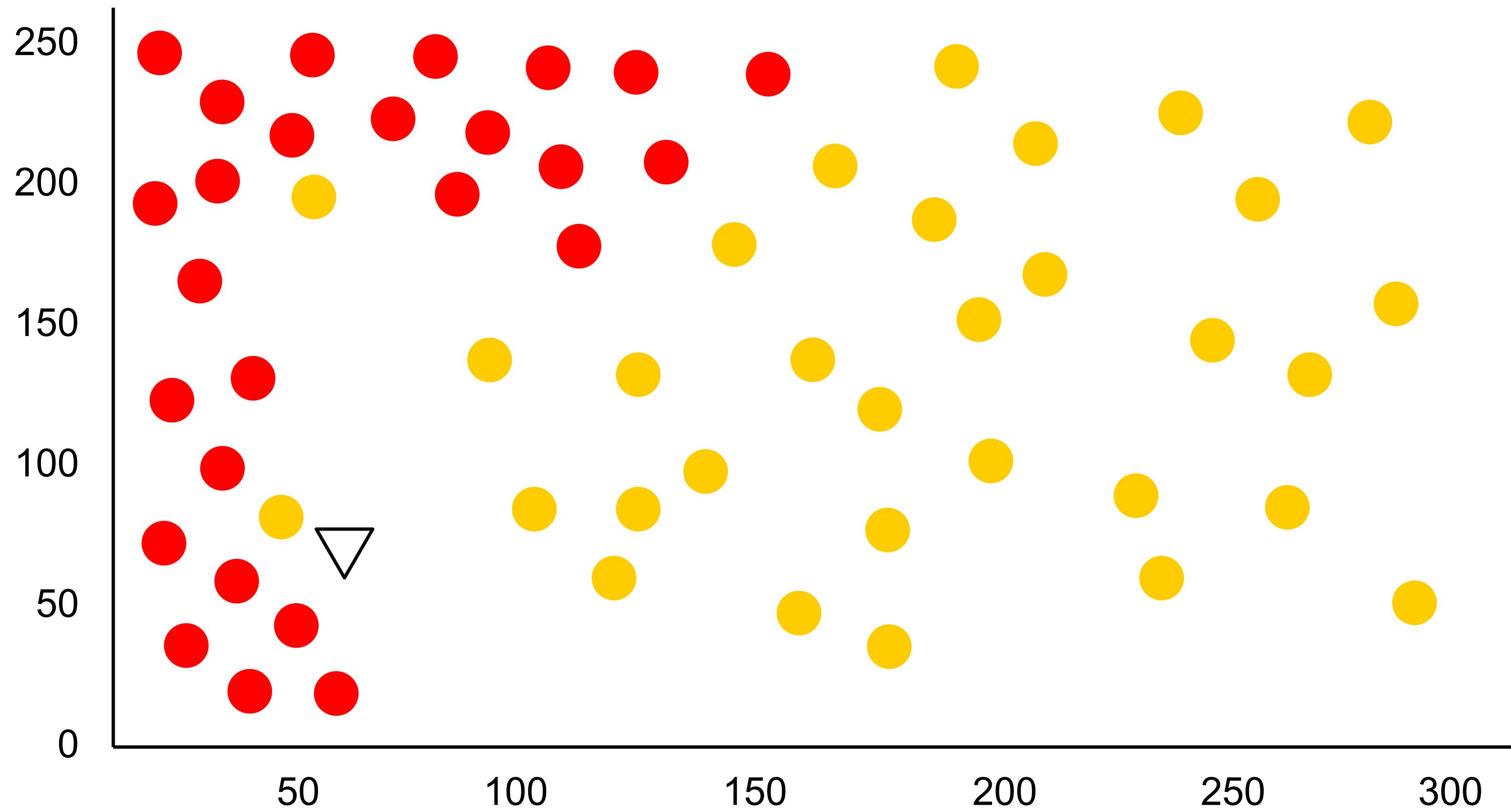
$$\rho(x, y) = \sum_{i=1}^n |x_i - y_i| \text{ — метрика } L_1$$

$$\rho(x, y) = \max_{i=1, \dots, n} |x_i - y_i| \text{ — метрика } L_\infty$$

Ещё можно вводить веса у признаков

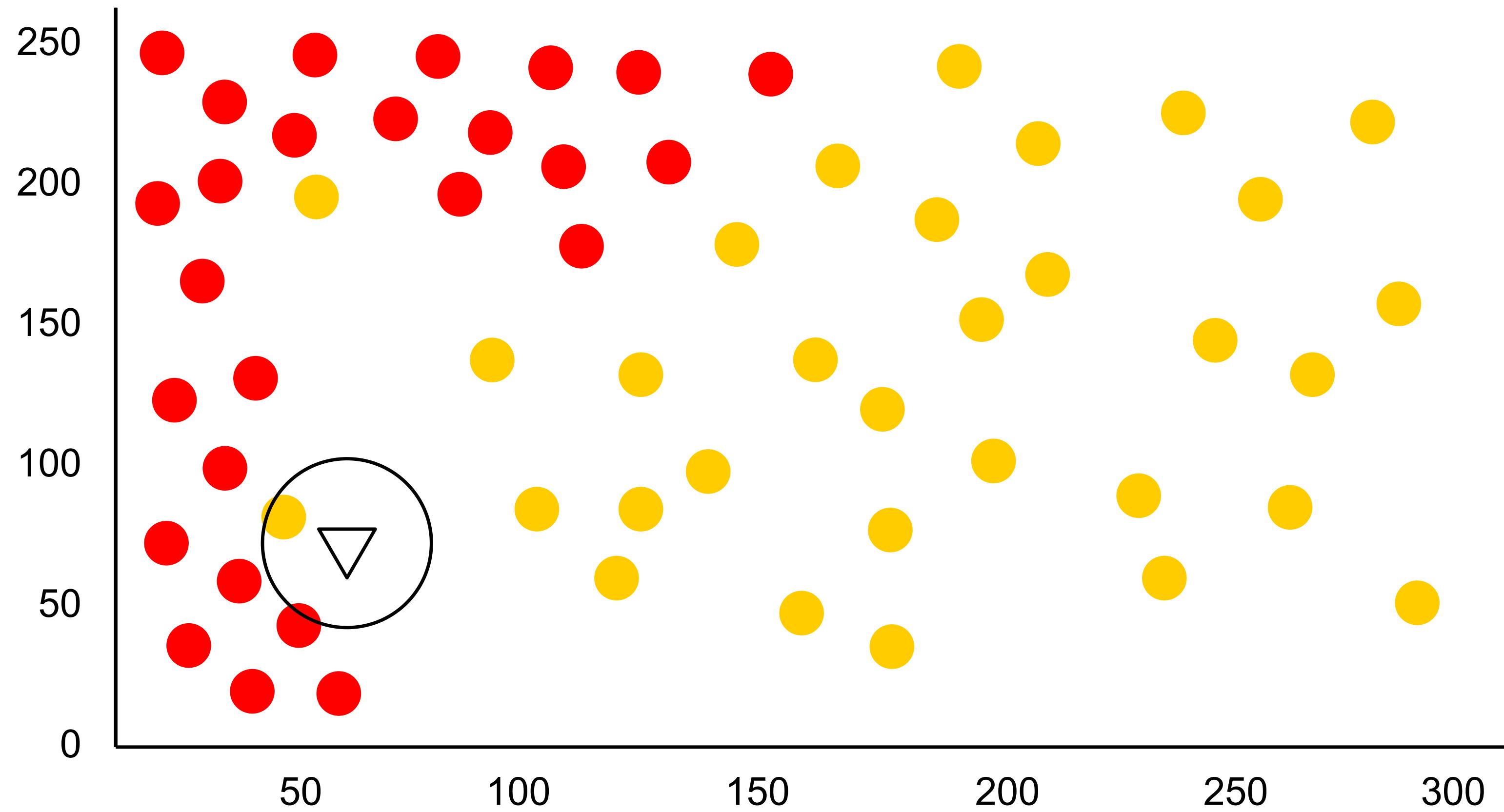
Метод одного ближайшего соседа

Красный или жёлтый новый объект?

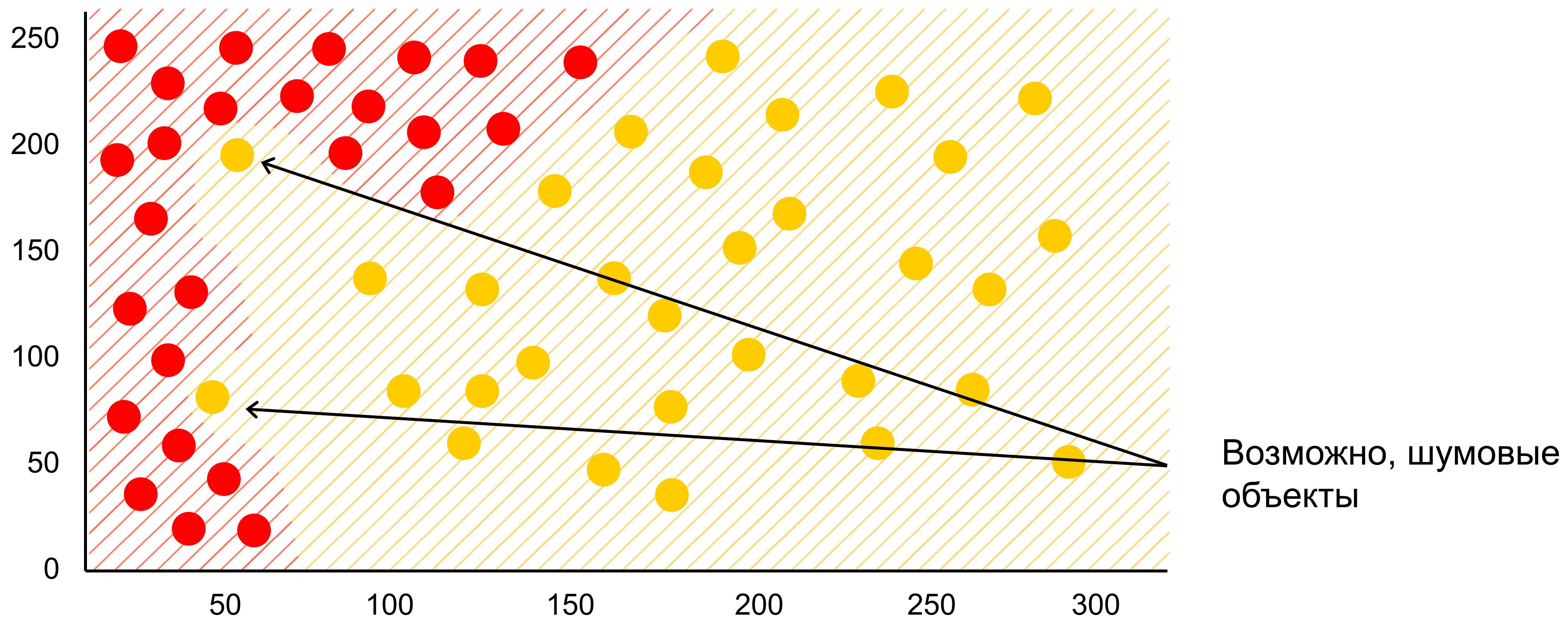


Метод одного ближайшего соседа

Пусть новый объект принадлежит к тому же классу,
что и его ближайший сосед

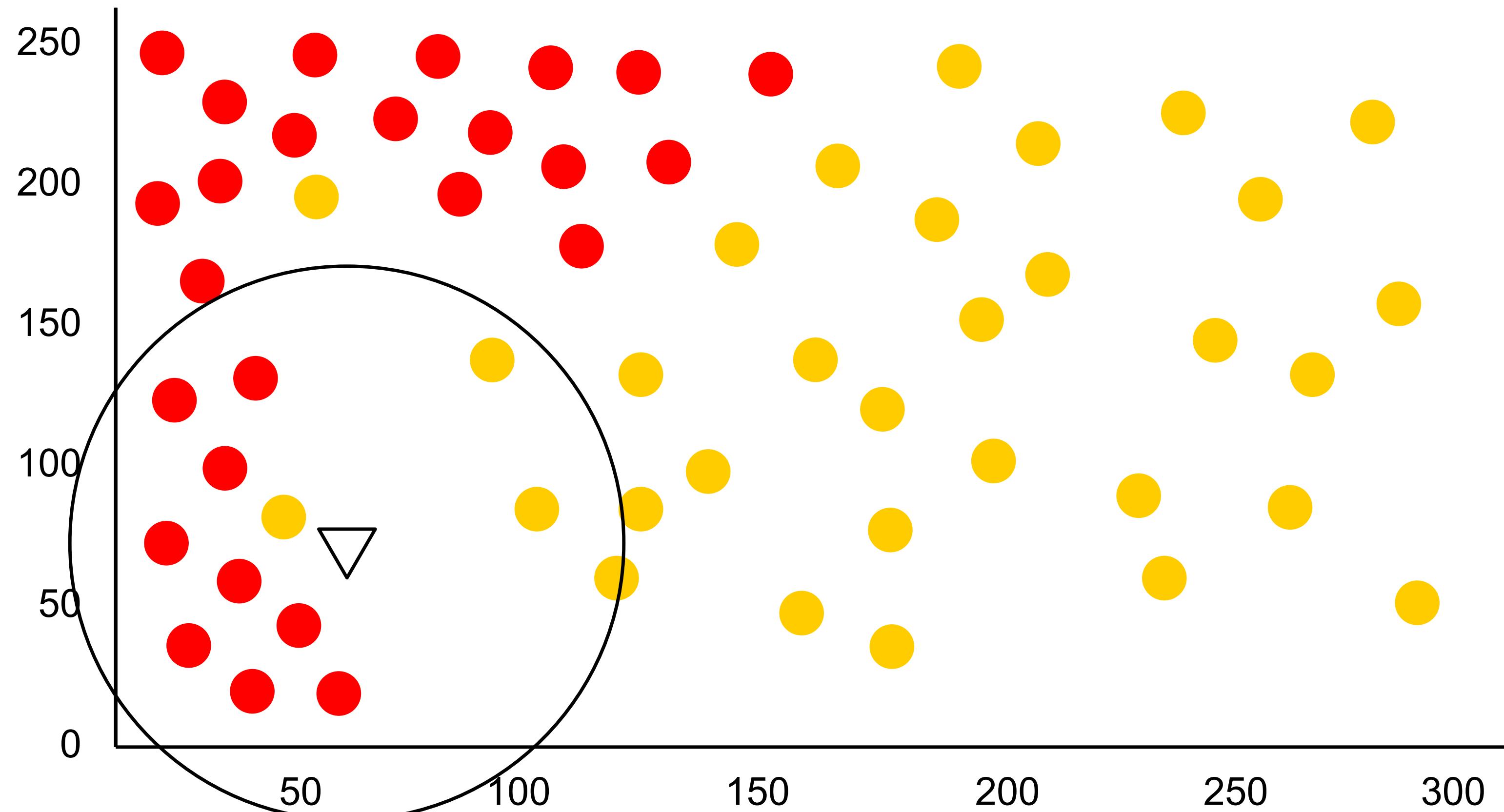


Граница разделения классов



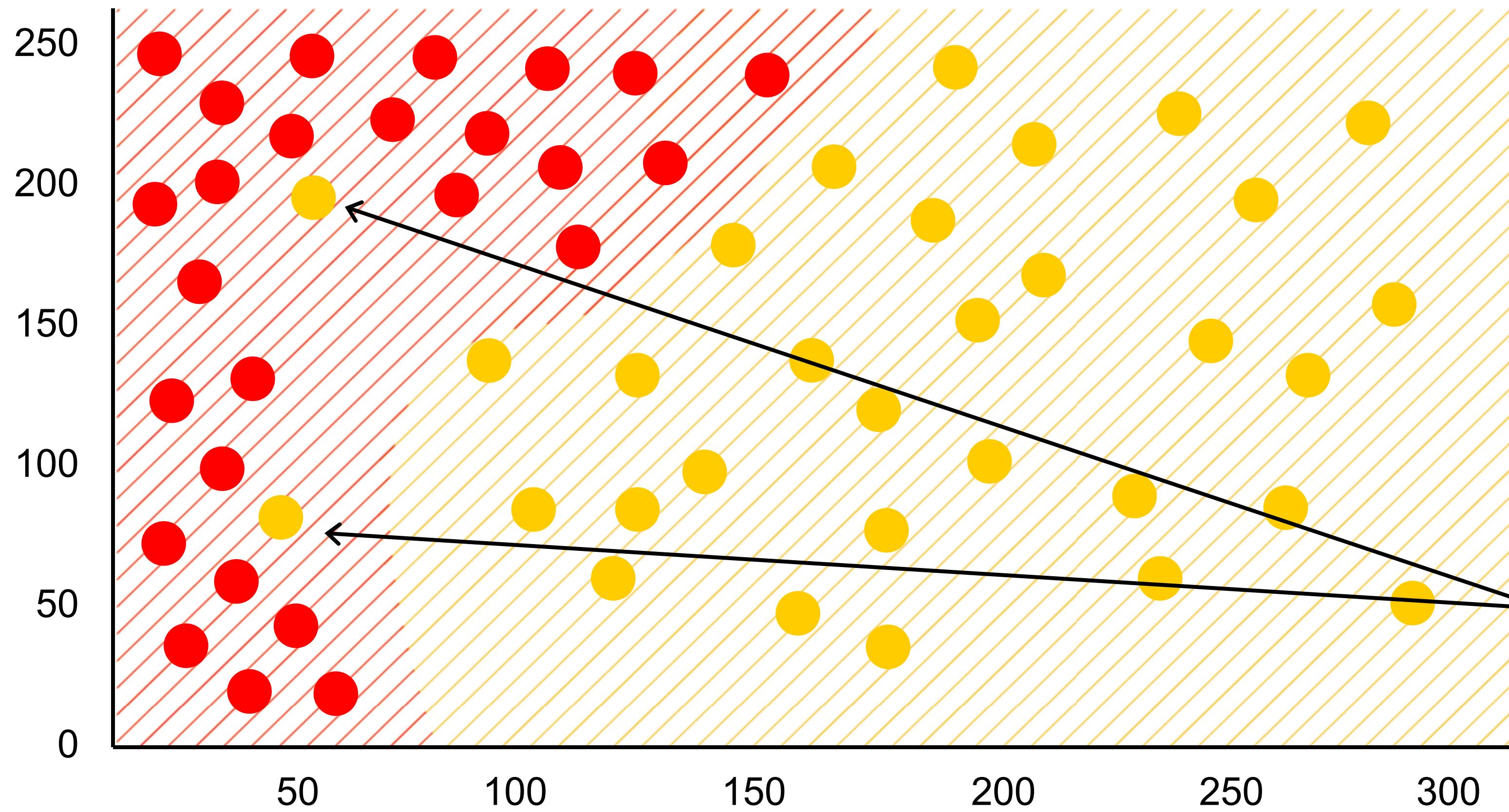
Метод k ближайших соседей

Новый объект принадлежит тому же классу,
что и большинство из k его соседей



Метод k ближайших соседей

Граница разделения классов $k=5$



Оказывается, алгоритм
дает ошибку на обучающей
выборке! А это и не плохо

A person is sitting cross-legged on a large stack of books, reading a book. The books are arranged in several rows, creating a tiered effect. The person is wearing a light-colored shirt and dark pants. The background is a plain, light-colored wall.

Логические методы

Логические закономерности

Логическая закономерность – это предикат $R : X \rightarrow \{0, 1\}$ удовлетворяющий двум требованиям:

1. Интерпретируемость:

R записывается на естественном языке

R зависит от небольшого числа признаков (1–7)

2. Информативность относительно
одного из классов: $c \in Y$

$$\#\{x_i \mid R(x_i) = 1 \text{ и } y_i = c\} \rightarrow \max$$

$$\#\{x_i \mid R(x_i) = 1 \text{ и } y_i \neq c\} \rightarrow \min$$

Часто используемые виды закономерностей

Пороговое условие (решающий пень):

$$R(x) = [f_j(x) \leq a_j] \text{ или } [a_j \leq f_j(x) \leq b_j]$$

Конъюнкция пороговых условий:

$$R(x) = \bigwedge_{j \in J} [a_j \leq f_j(x) \leq b_j]$$

Синдром – выполнение не менее d условий из J :

$$R(x) = \left[\sum_{j \in J} [a_j \leq f_j(x) \leq b_j] \geq d \right]$$

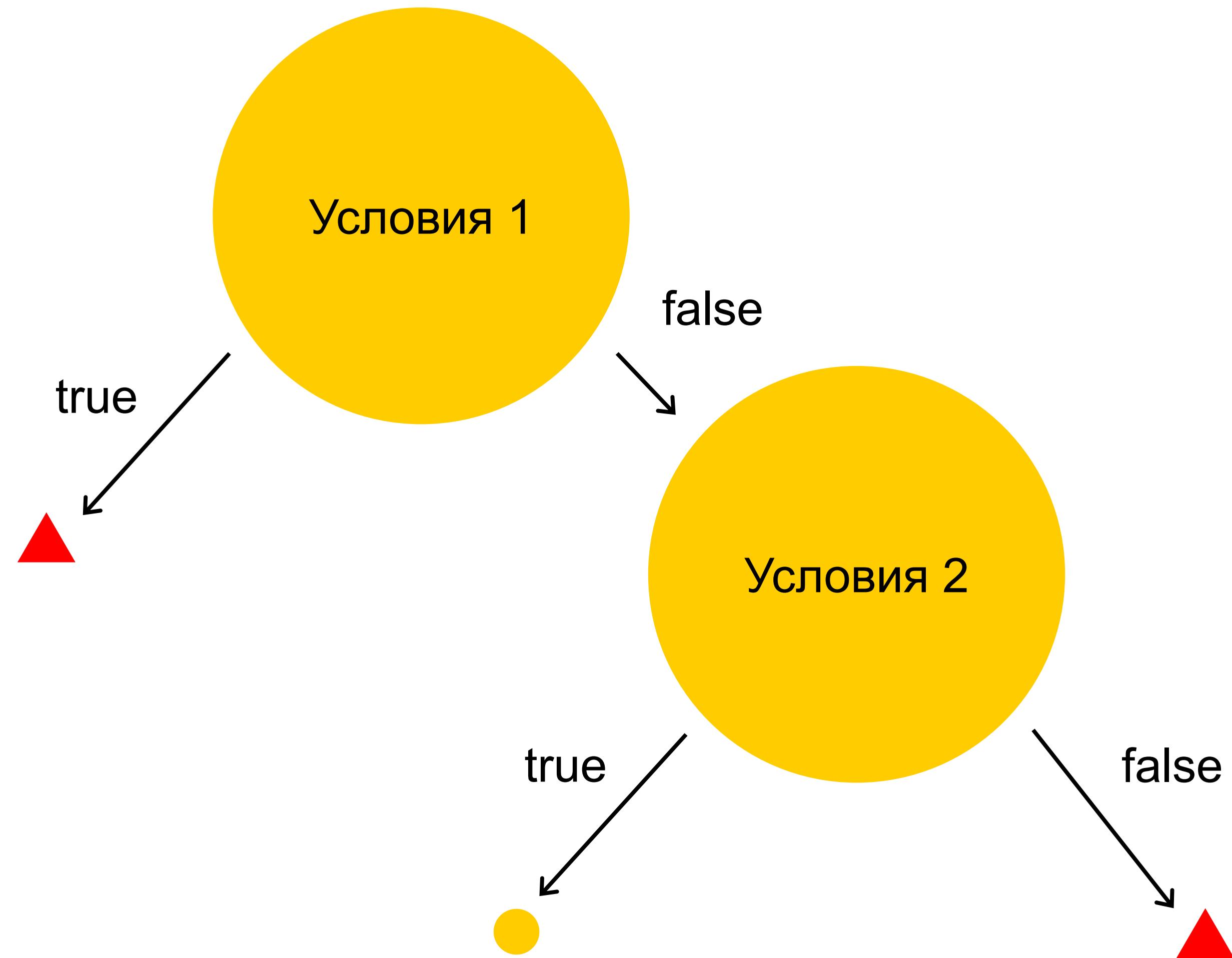
Полуплоскость – линейная пороговая функция:

$$R(x) = \left[\sum_{j \in J} w_j f_j(x) \geq w_0 \right]$$

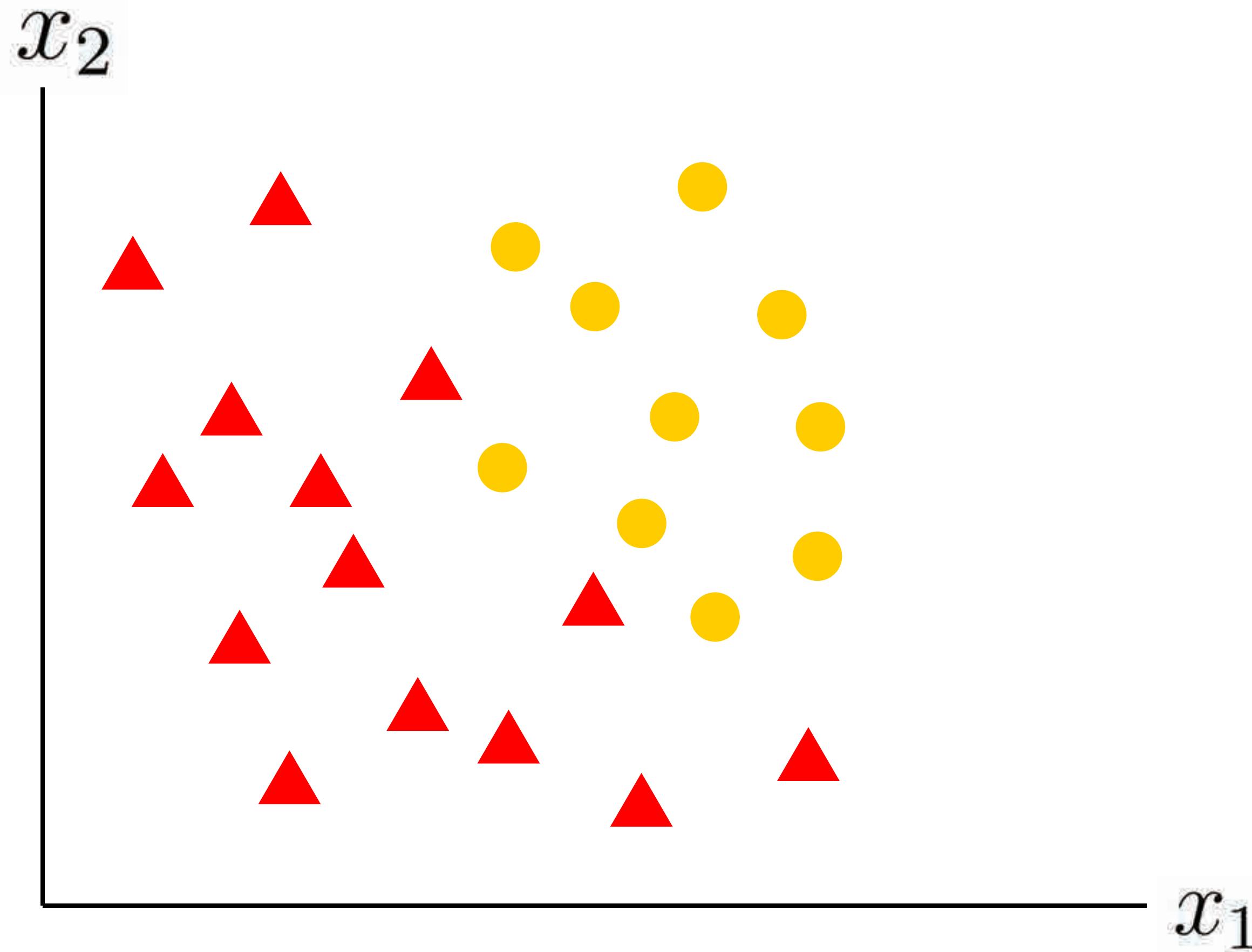
Шар – пороговая функция близости:

$$R(x) = [\rho(x, x_0) \leq w_0]$$

Решающие деревья



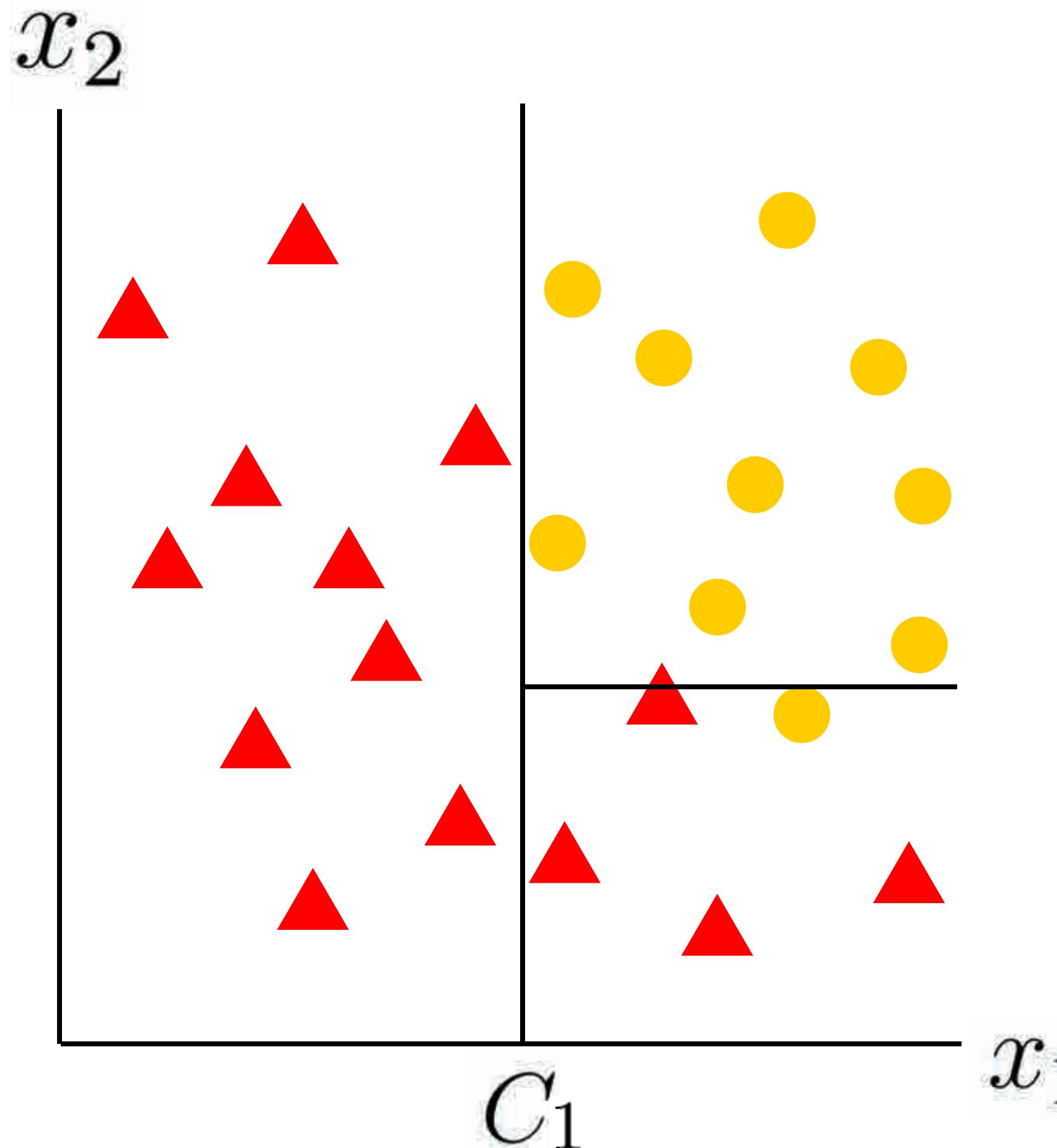
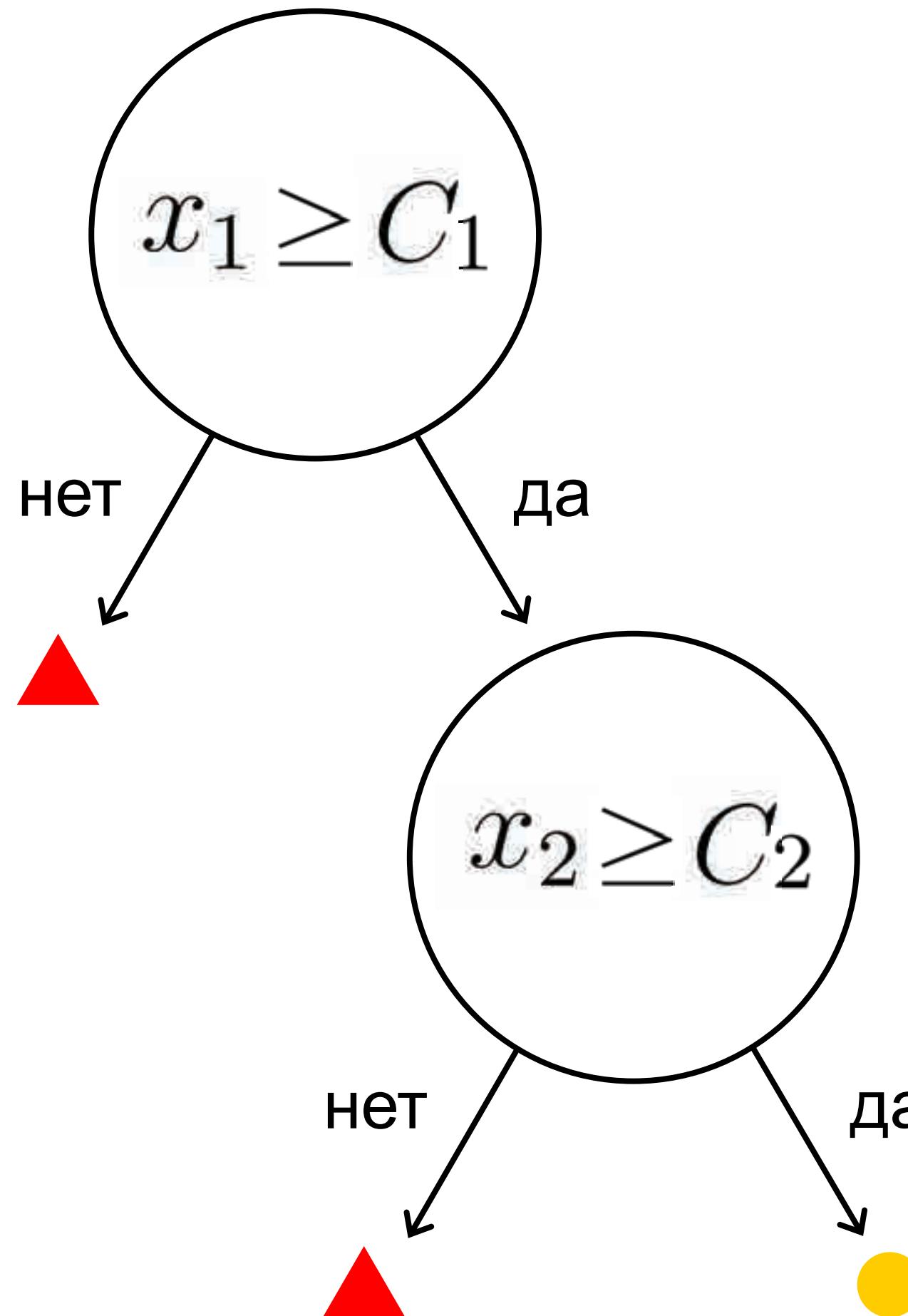
Решающие деревья



Какие условия будут
в дереве?

Попробуем использовать
пороговые условия
перехода в виде пороговых
правил: $x > C$

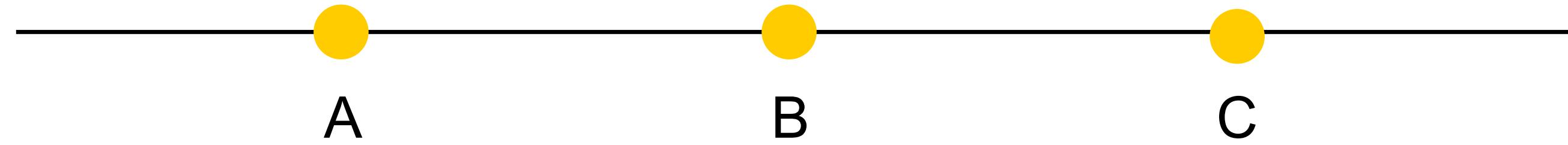
Решающие деревья



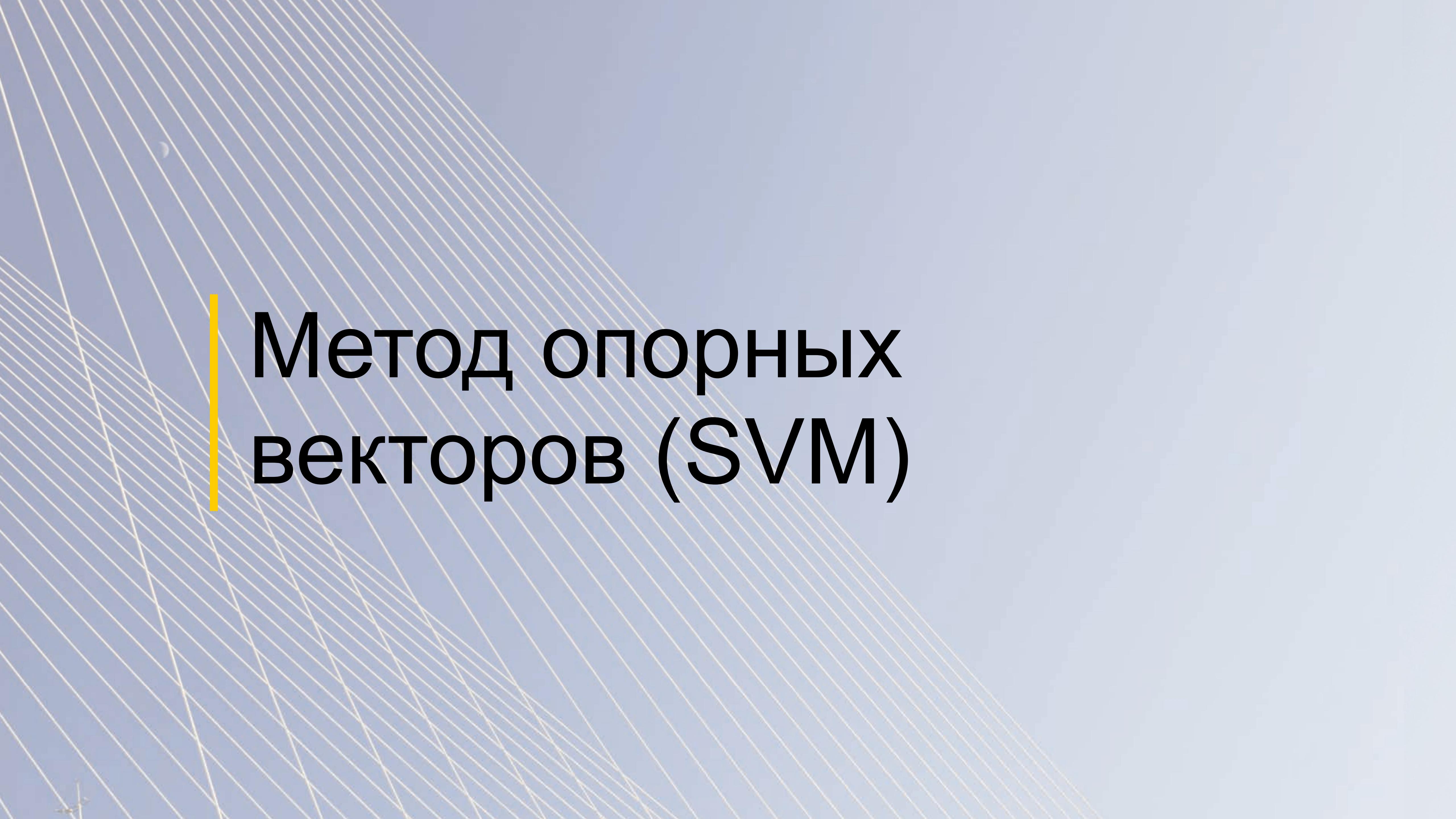
Каждый раз берём наиболее «информационное» разделение текущей области

Преимущества логических алгоритмов перед метрическими

1. Придумать правильную меру сходства – значит почти решить задачу, это сложно. А решающие деревья не используют метрики
2. Единственное, что используют деревья – точка В ближе к А, чем С по данному признаку



3. Устойчивы к монотонным преобразованиям признаков



Метод опорных векторов (SVM)

Линейные методы классификации

Дано:

Обучающая выборка $X^l = \{(x_i, y_i) \mid i = 1, \dots, l\}$

x_i – объекты, векторы из множества $X = \mathbb{R}^n$

y_i – ответы из множества $Y = \{-1; +1\}$

Найти:



Параметры $w \in \mathbb{R}^n$, $w_0 \in \mathbb{R}$ линейной модели классификации



$a(x; w, w_0) = \text{sign} (\langle x, w \rangle - w_0)$

Понятие отступа

Отступ
положителен

объект правильно
расклассифицирован

Отступ
отрицателен

объект неправильно
расклассифицирован

$$M_i(w, w_0) = (\langle x_i, w \rangle - w_0)y_i$$

– отступ объекта x_i

Задача SVM

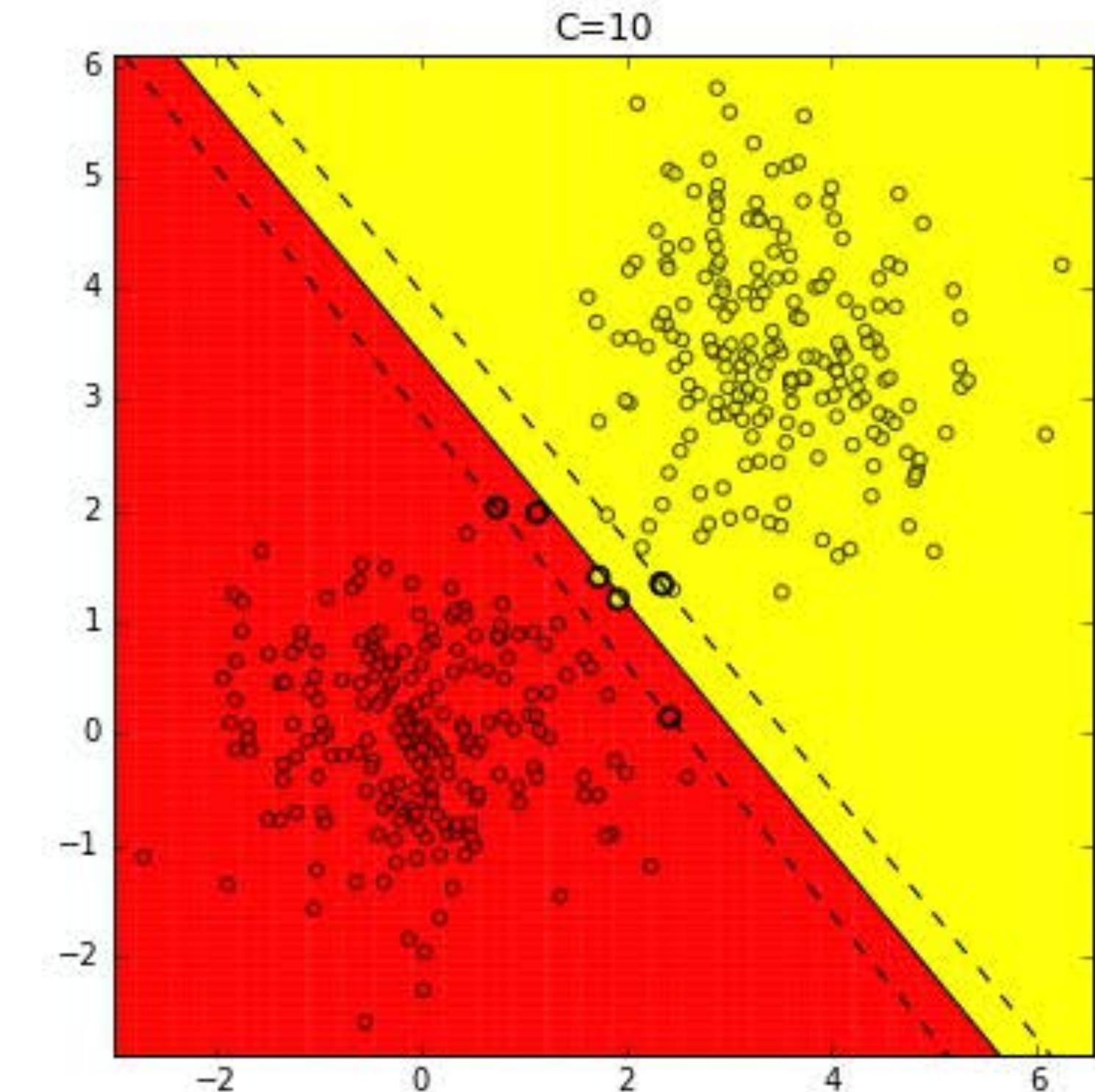
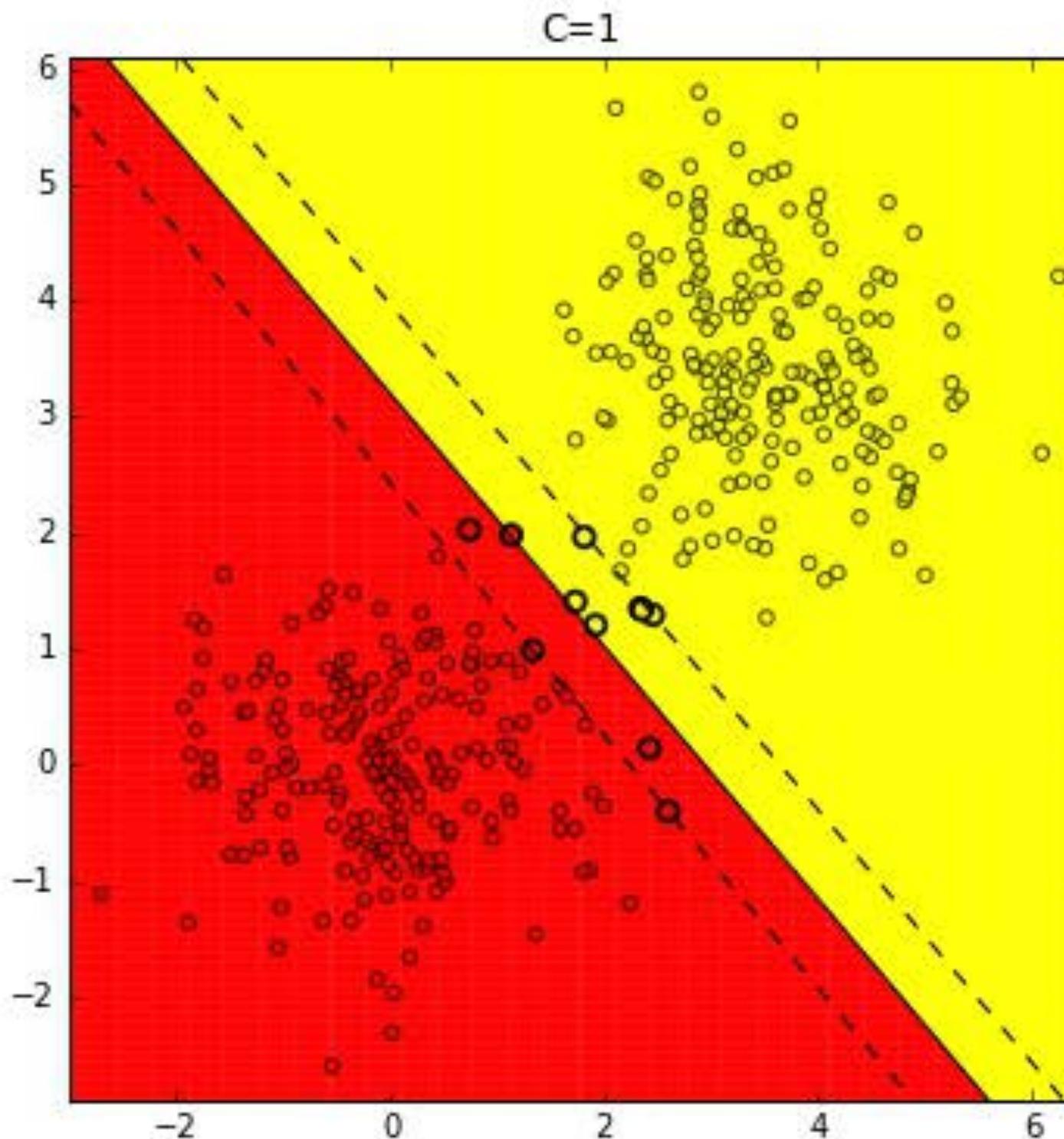
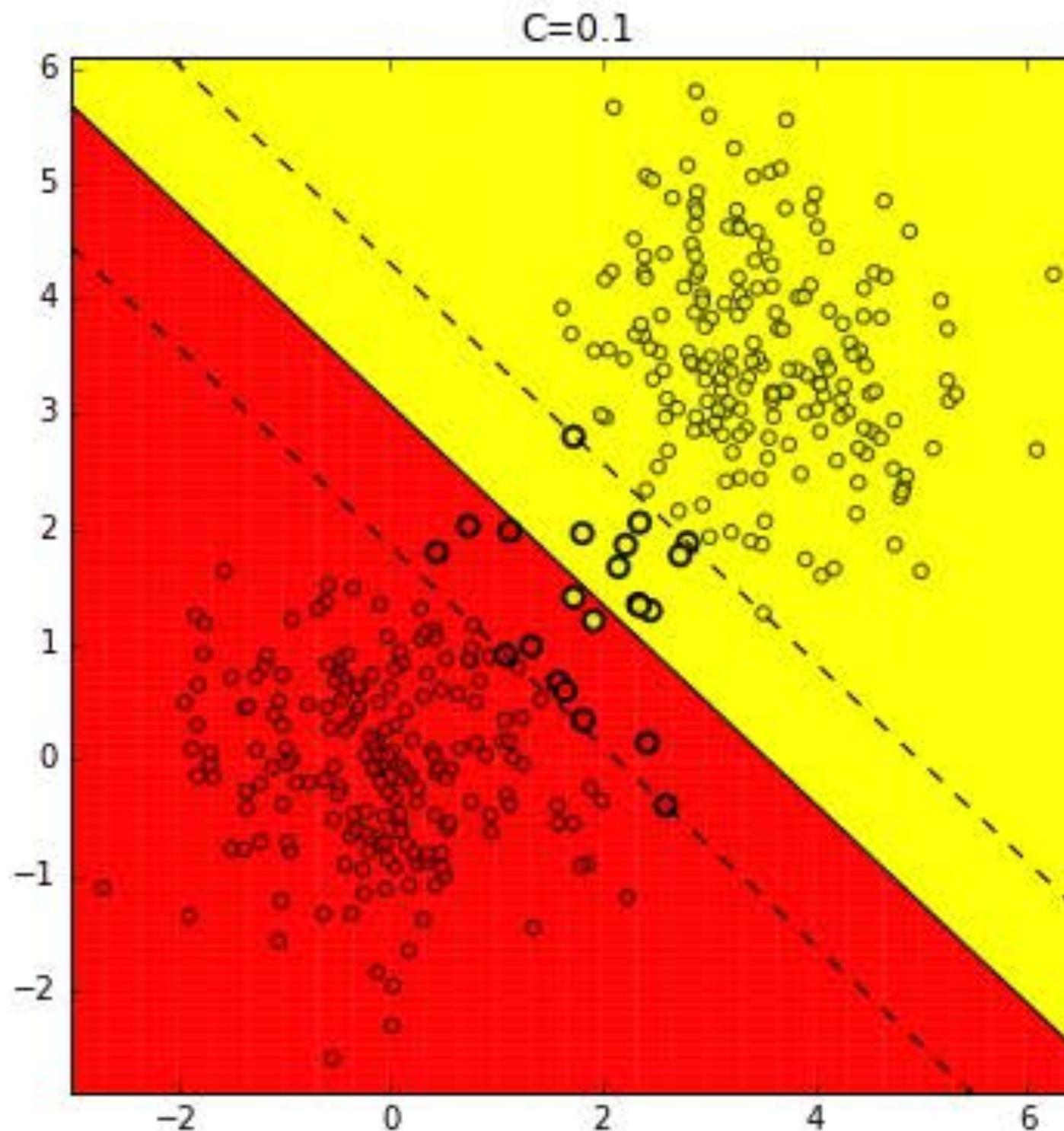
Минимизация эмпирического риска

$$\sum_{i=1}^l [a(x_i, w, w_0) \neq y_i] = \sum_{i=1}^l [M_i(w, w_0) < 0] \rightarrow \min_{w, w_0}$$

Заменим получившуюся кусочно-постоянную функцию оценкой сверху, непрерывной по параметрам:

$$\sum_{i=1}^l [M_i(w, w_0) < 0] \leq \sum_{i=1}^l [1 - M_i(w, w_0)]_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}$$

Применение SVM с разными константами



$$\sum_{i=1}^l [1 - M_i(w, w_0)]_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}$$

Двойственная задача

Эквивалентная задача (после применения теоремы Кароша-Куна-Таккера):

$$\begin{cases} -\sum_{i=1}^l \lambda_i + \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \min_{\lambda}, \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, l, \\ \sum_{i=1}^l \lambda_i y_i = 0. \end{cases}$$

Нелинейный SVM

В двойственной задаче элементы x_i участвуют только в скалярных произведениях.

И постановка задачи, и сам классификатор зависят только от способности вычислять скалярные произведения объектов.

Новая идея: заменить скалярное произведение функцией от двух переменных, если её можно считать скалярным произведением в некотором (пусть другом) пространстве.

Понятие ядра

Переход к спрямляющему пространству
более высокой размерности

Функция $K : (X, X) \rightarrow \mathbb{R}$ называется ядром,
если $\forall x, x' \in X \quad K(x, x') = \langle \psi(x), \psi(x') \rangle$

для некоторой функции $\psi : X \rightarrow \mathcal{H}$,
где \mathcal{H} – пространство со скалярным
произведением

Популярные ядра

$$K(x, x') = \langle x, x' \rangle^2$$

– квадратичное ядро

$$K(x, x') = \langle x, x' \rangle^d$$

– полиномиальное ядро
с мономами степени d

$$K(x, x') = (\langle x, x' \rangle + 1)^d$$

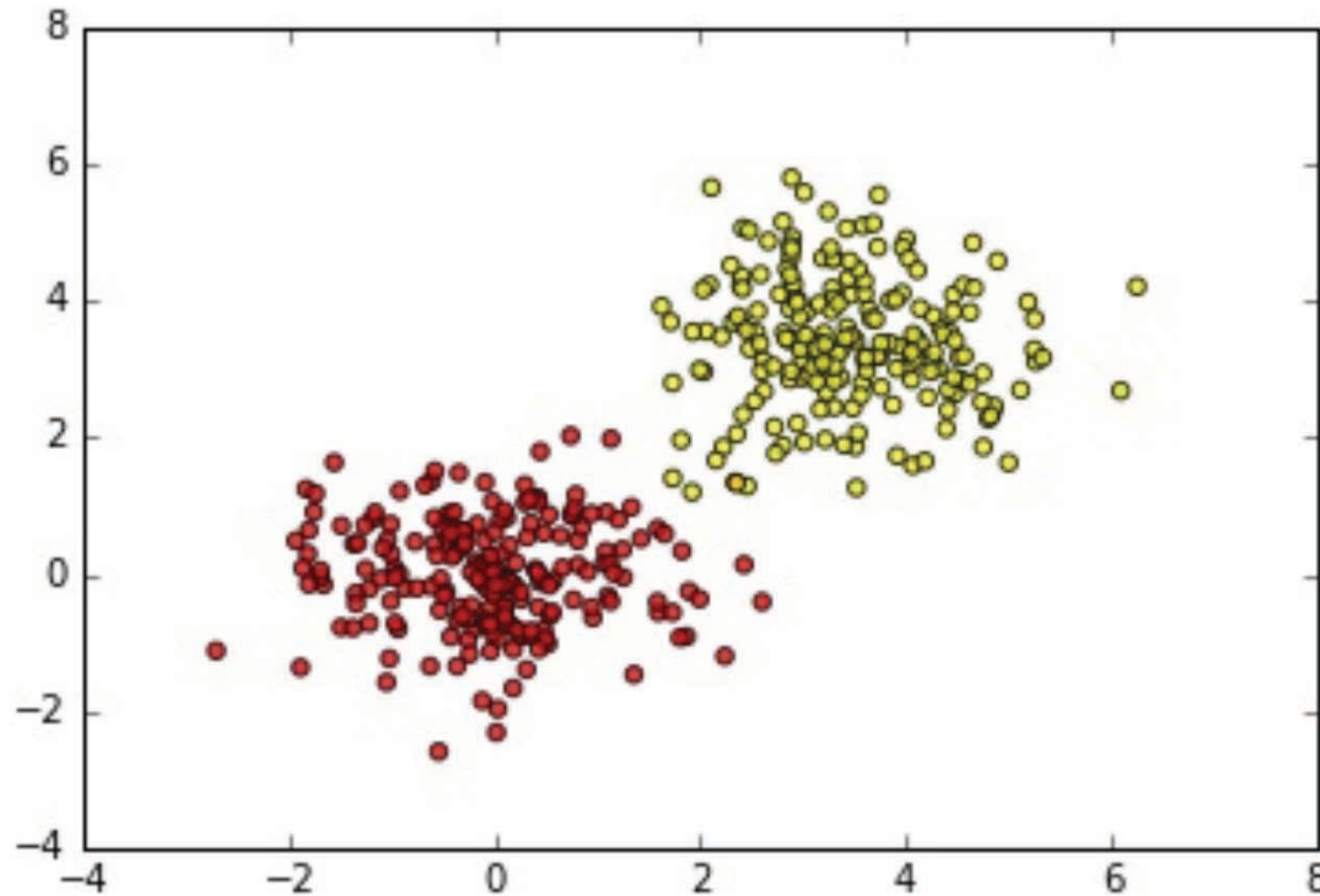
– полиномиальное ядро
с мономами степени $\leq d$

$$K(x, x') = \exp(-\gamma \|x - x'\|^2)$$

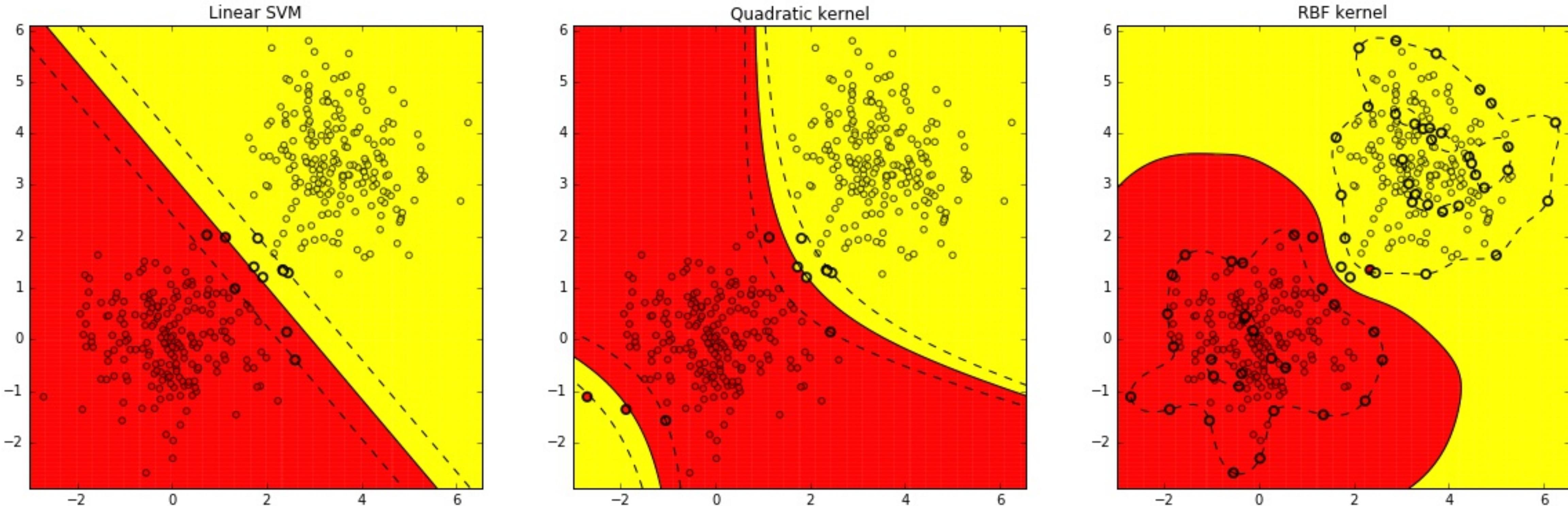
– сеть радиальных базисных
функций (RBF-ядро)

Гиперплоскость в пространстве \mathcal{H} соответствует нелинейной разделяющей поверхности в X

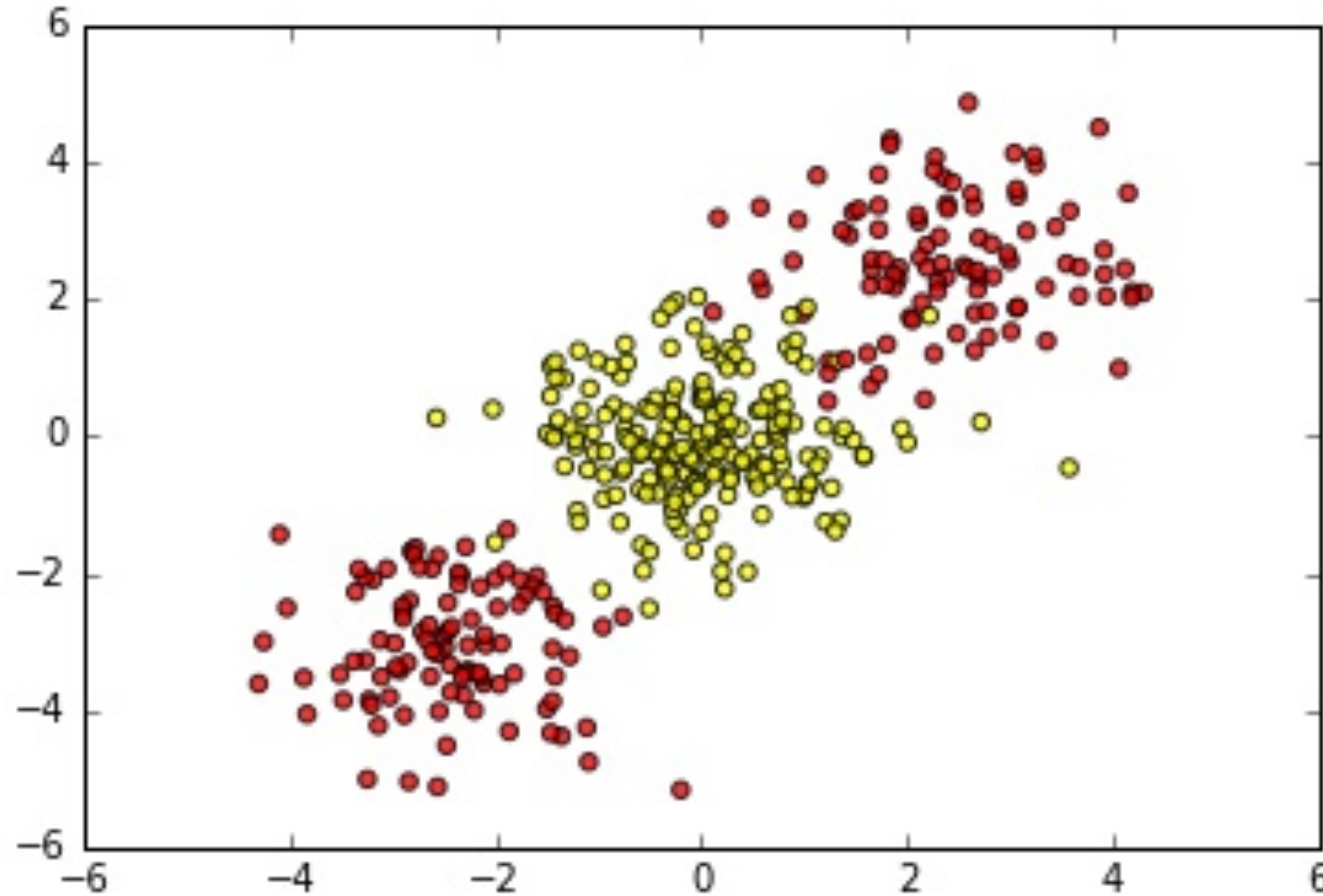
Линейно разделимая выборка



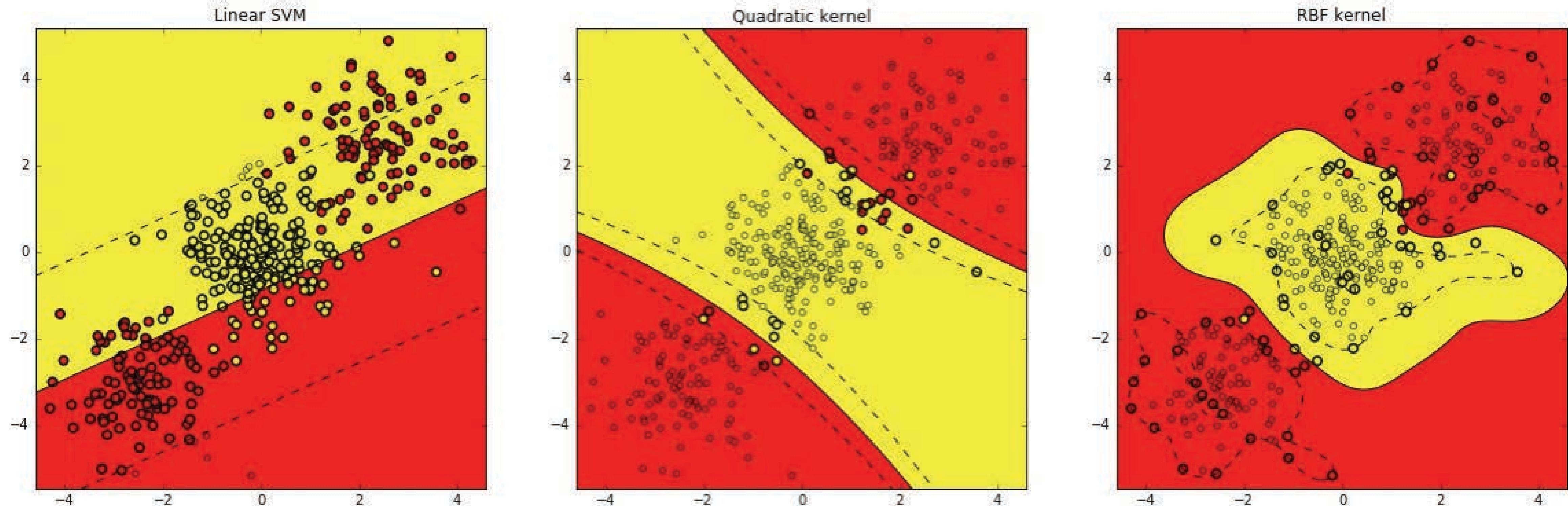
Применение SVM с разными ядрами



Линейно не разделимая выборка



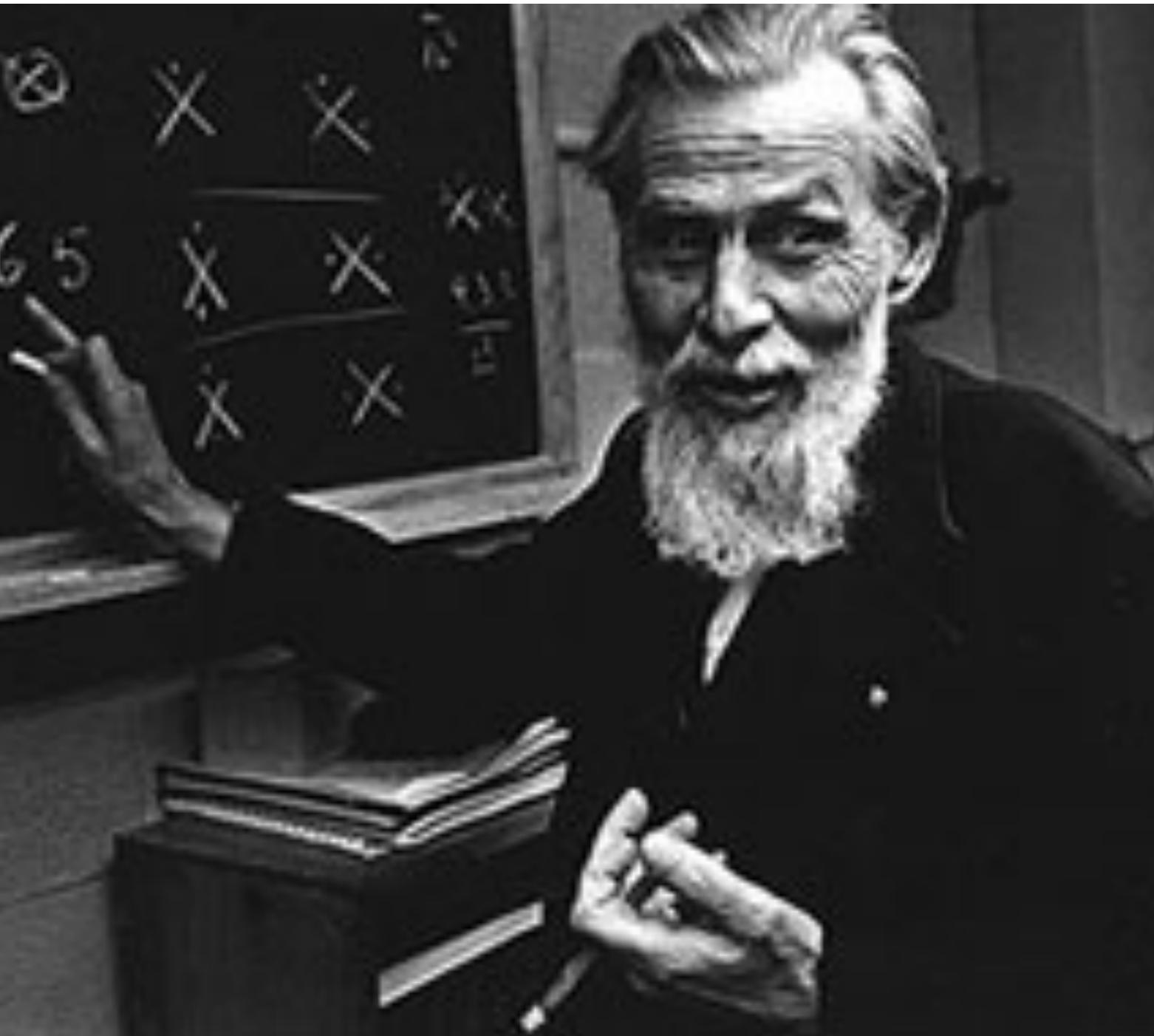
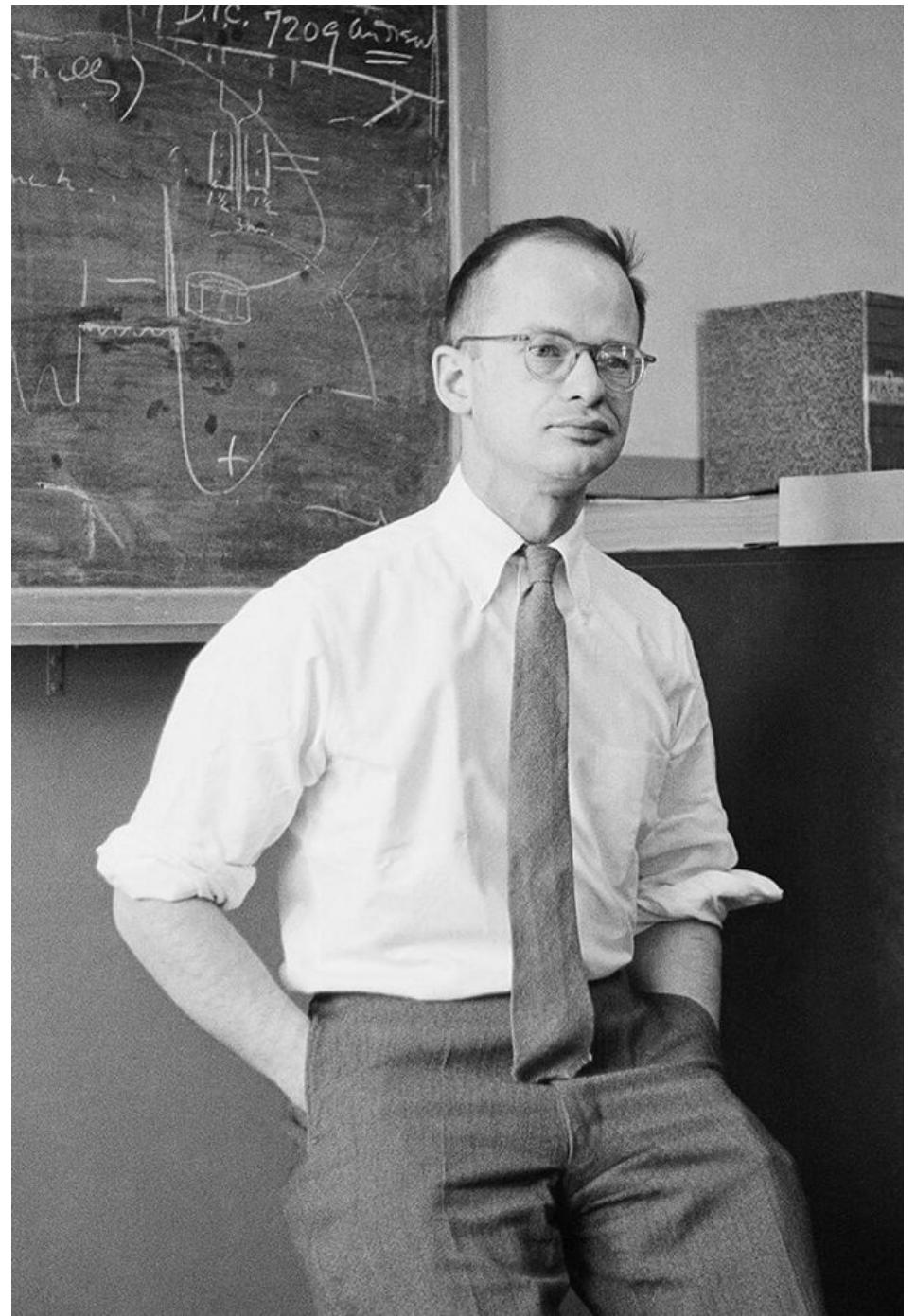
Применение SVM с разными ядрами



A server rack filled with many orange network cables. The cables are bundled together and run across the front of the server units. The server units have various ports and labels visible, such as 'LA1', 'LA2', 'LA3', etc. The background is dark, and the cables are brightly lit.

Нейронные сети

Немного истории



Уоррен Маккалок, Уолтер Питтс
Первый формальный нейрон – 1943

Фрэнк Розенблатт
Первая искусственная
нейронная сеть – 1957

Линейная модель нейрона Маккалока-Питтса

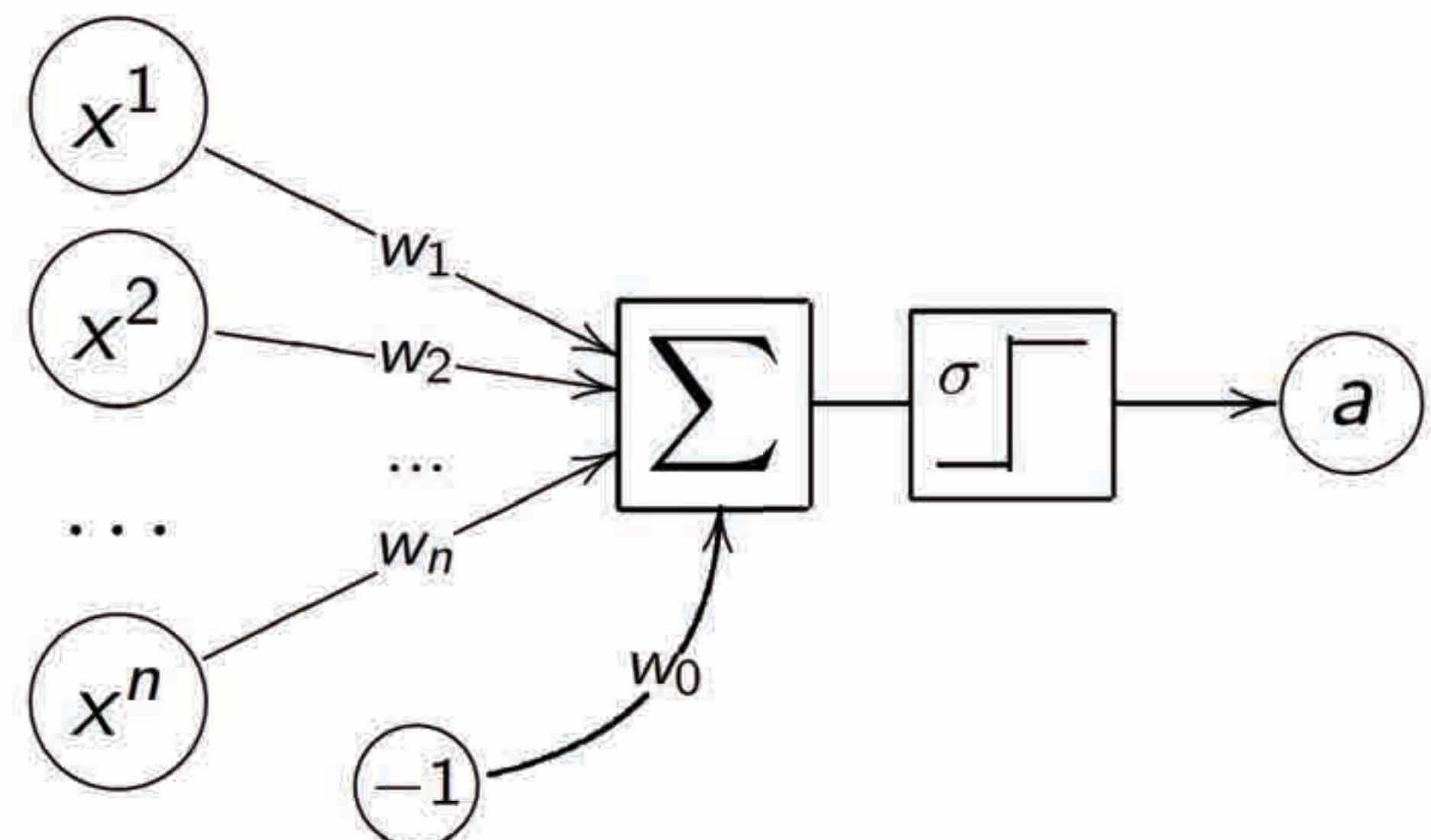
$$a(x, w) = \sigma(\langle x, w \rangle) = \sigma \left(\sum_{j=1}^n w_j f_j(x) - w_0 \right)$$

где $\sigma(z)$ – функция активации
(например, sign)

w_j – весовые коэффициенты сигналов
 w_0 – порог активации

$w, x \in \mathbb{R}^{n+1}$, если ввести константный признак

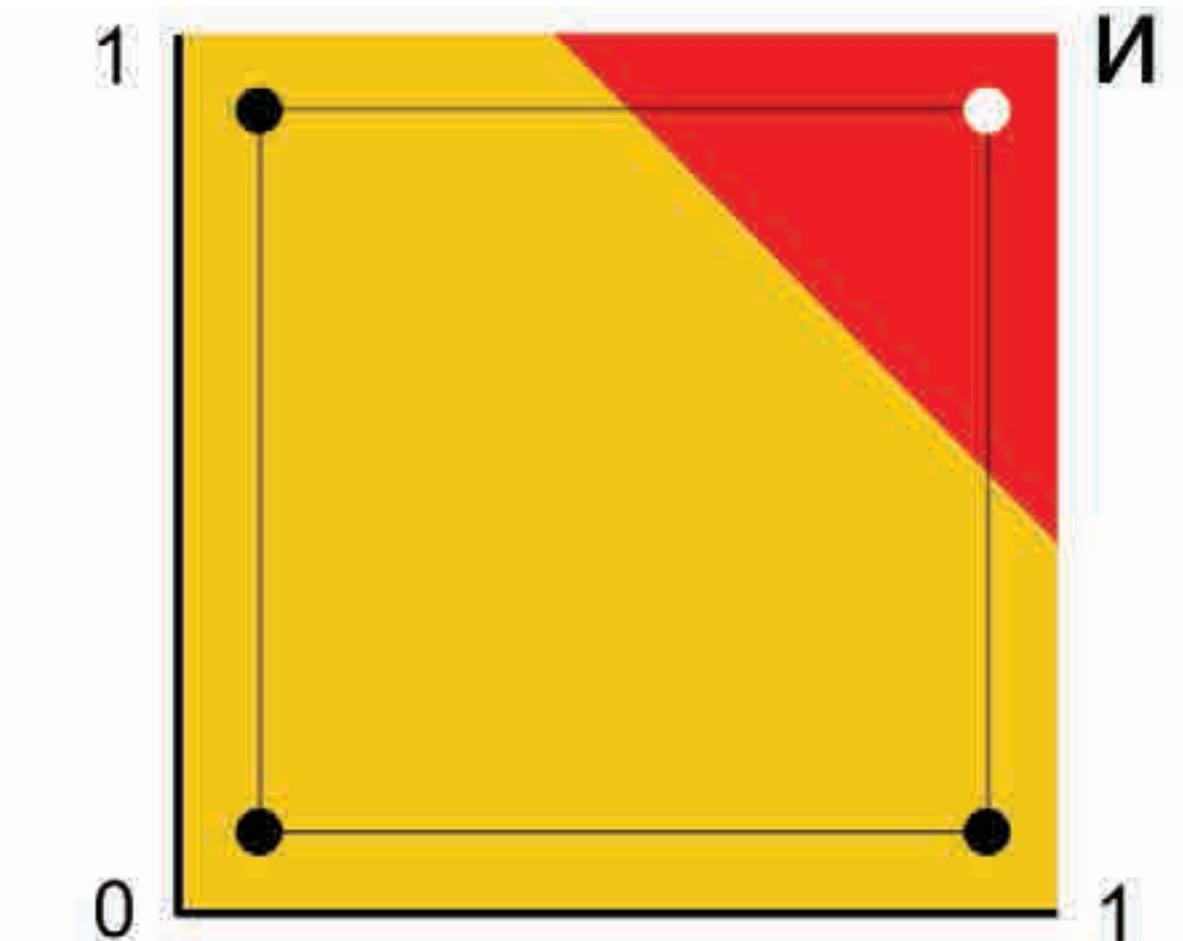
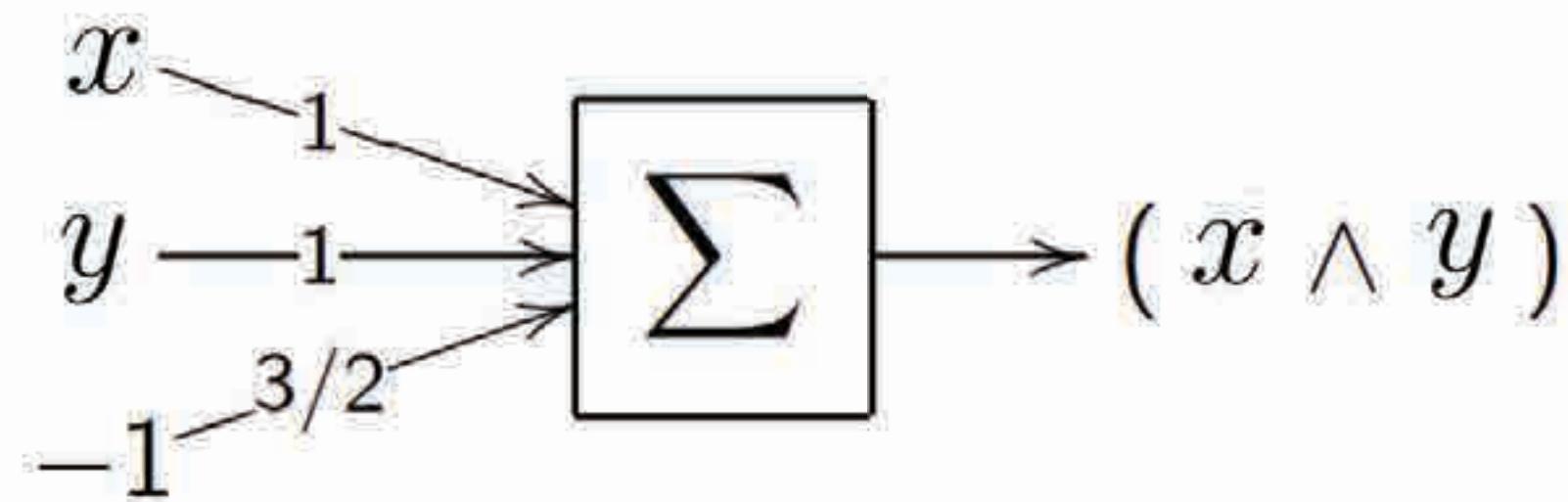
$f_0(x) \equiv 1$



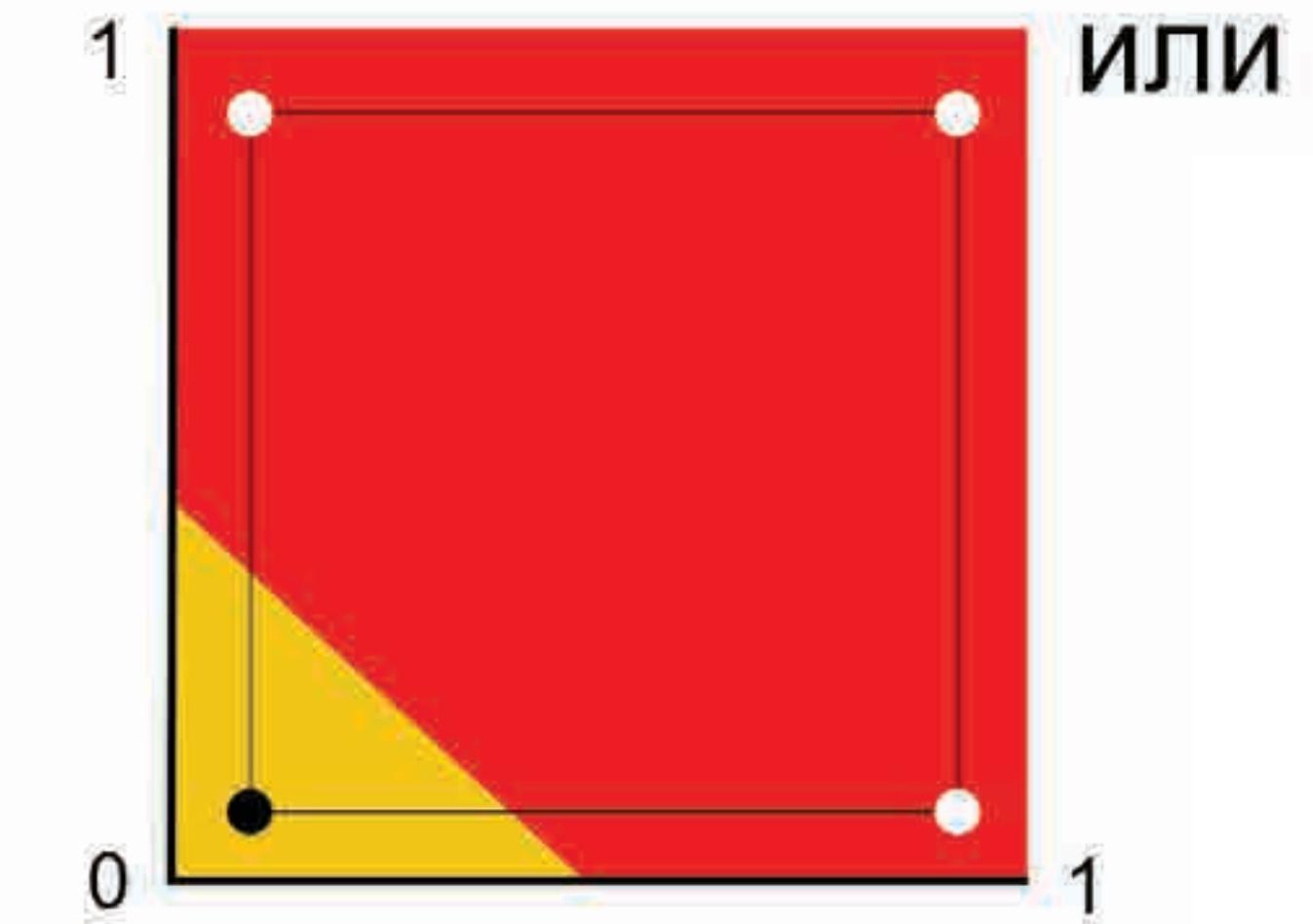
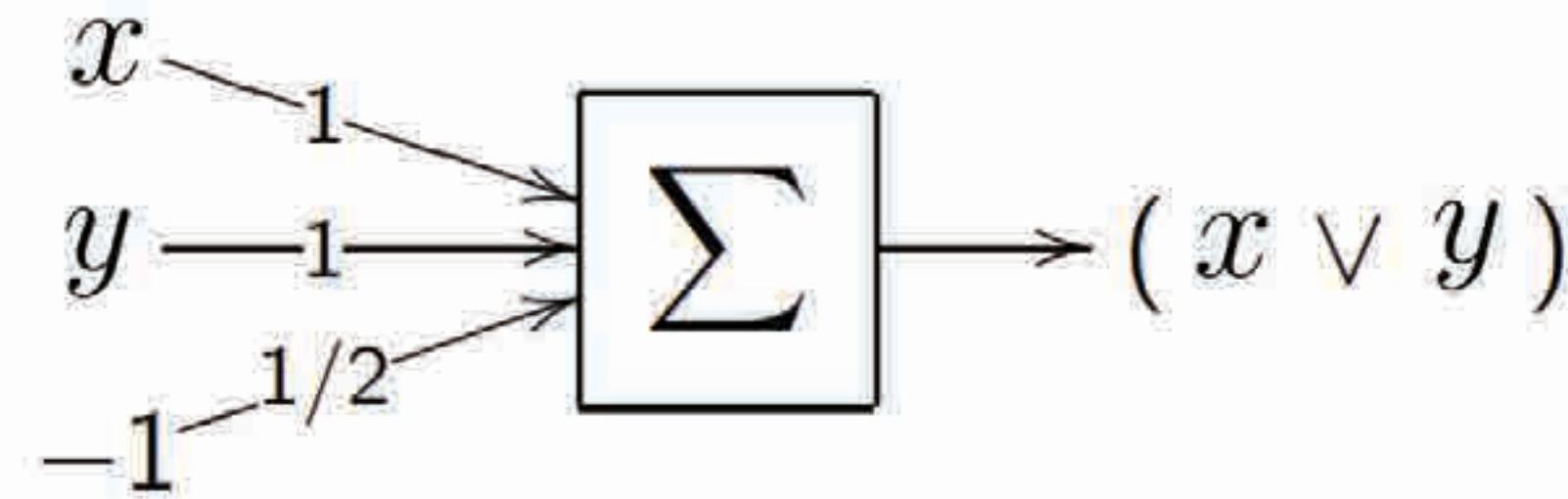
Нейронная реализация логических функций

Функции И, ИЛИ от бинарных переменных x и y

$$x \wedge y = \left[x + y - \frac{3}{2} > 0 \right]$$



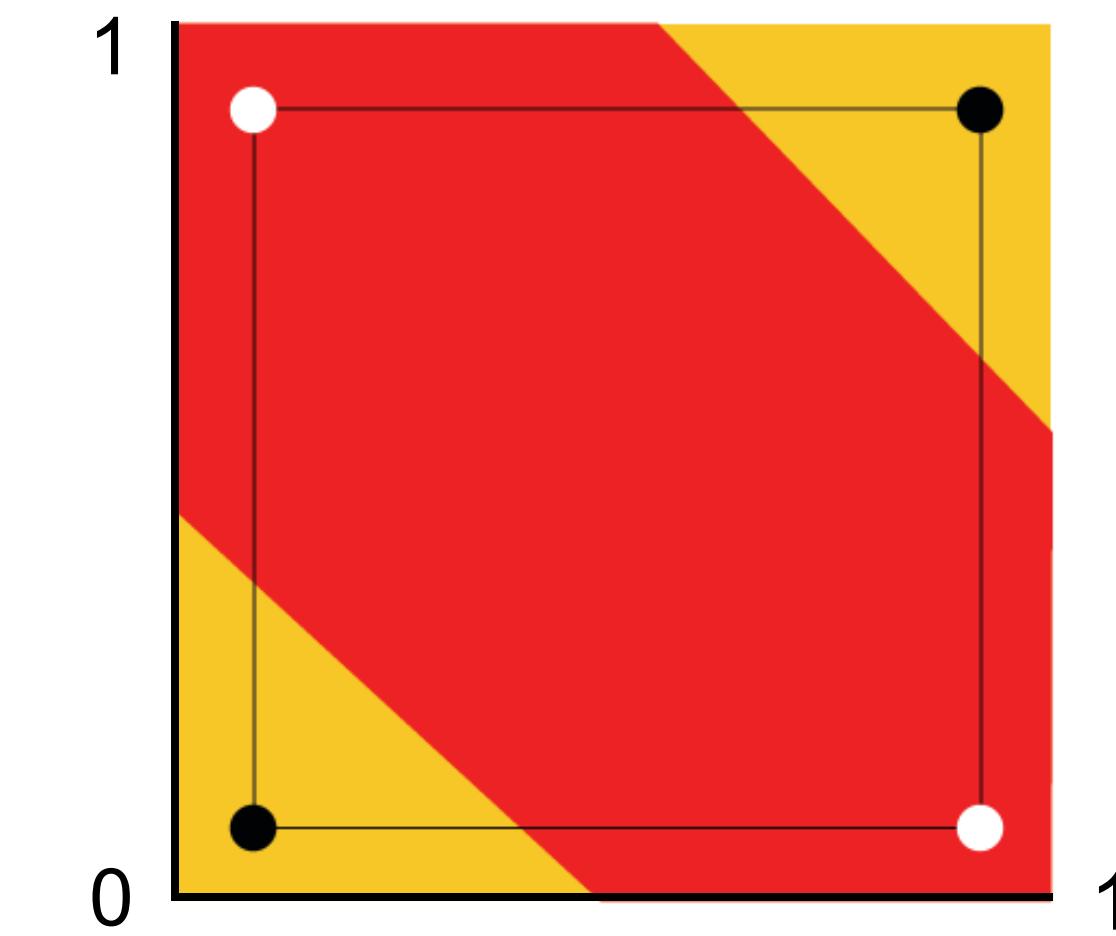
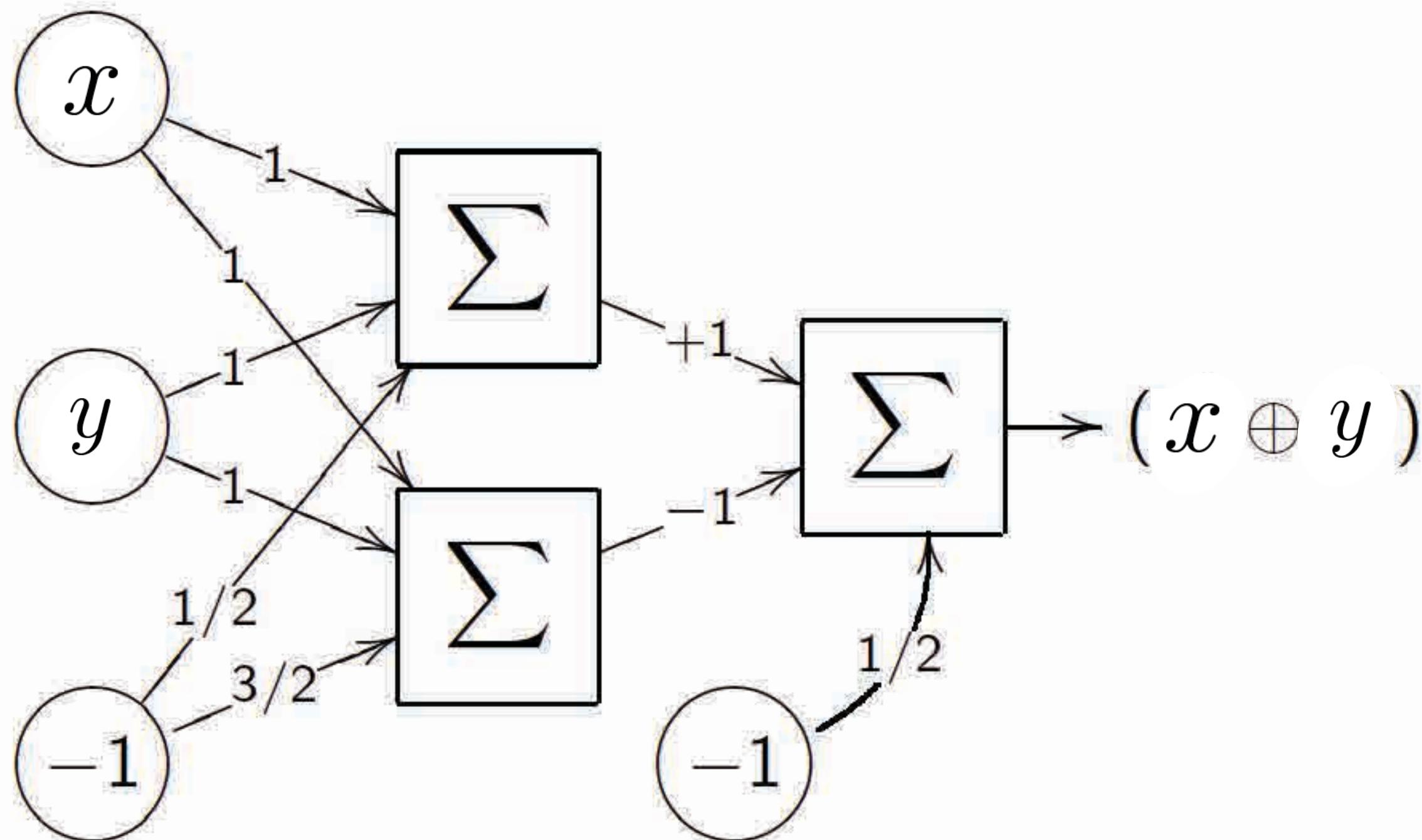
$$x \vee y = \left[x + y - \frac{1}{2} > 0 \right]$$



Реализация исключающего ИЛИ

Функция $x \oplus y = [x \neq y]$ не реализуема ни одним нейроном. Сеть (двухслойная суперпозиция) функций И, ИЛИ, НЕ

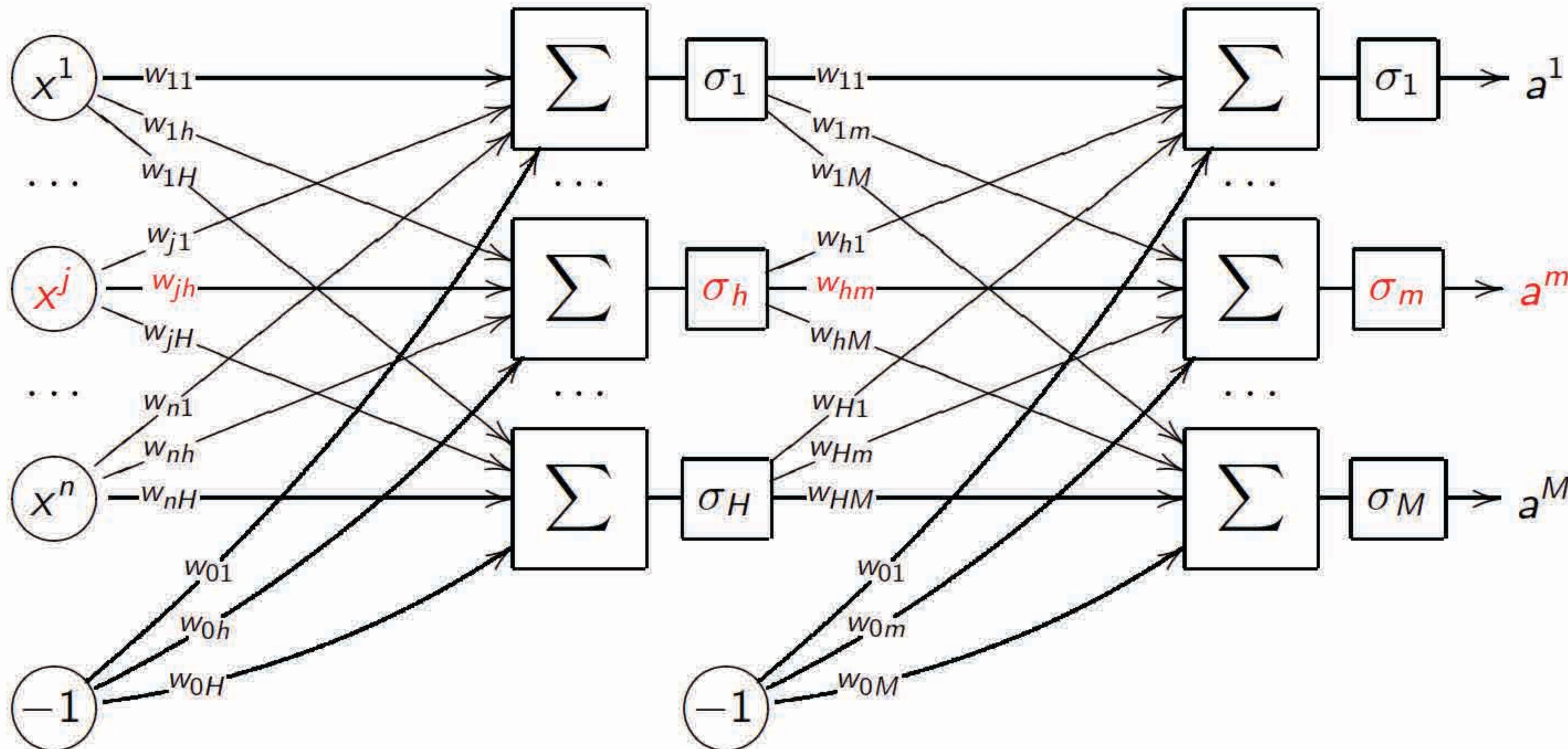
$$x \oplus y = \left[(x \vee y) - (x \wedge y) - \frac{1}{2} > 0 \right]$$



Представление функций нейросетью

- Двухслойная сеть в $\{0, 1\}^n$ позволяет реализовать произвольную булеву функцию
- Двухслойная сеть в \mathbb{R}^n позволяет отделить произвольный многогранник
- Трёхслойная сеть в \mathbb{R}^n позволяет отделить произвольную многогранную область, не обязательно выпуклую и связную
- С помощью линейных операций и одной нелинейной функции активации можно приблизить любую непрерывную функцию с любой заданной точностью

Многослойная нейронная сеть



Персепtron Розенблатта

Задача классификации с двумя классами $y_i \in \{0, 1\}$
и бинарными признаками $f_j(x) \in \{0, 1\}$

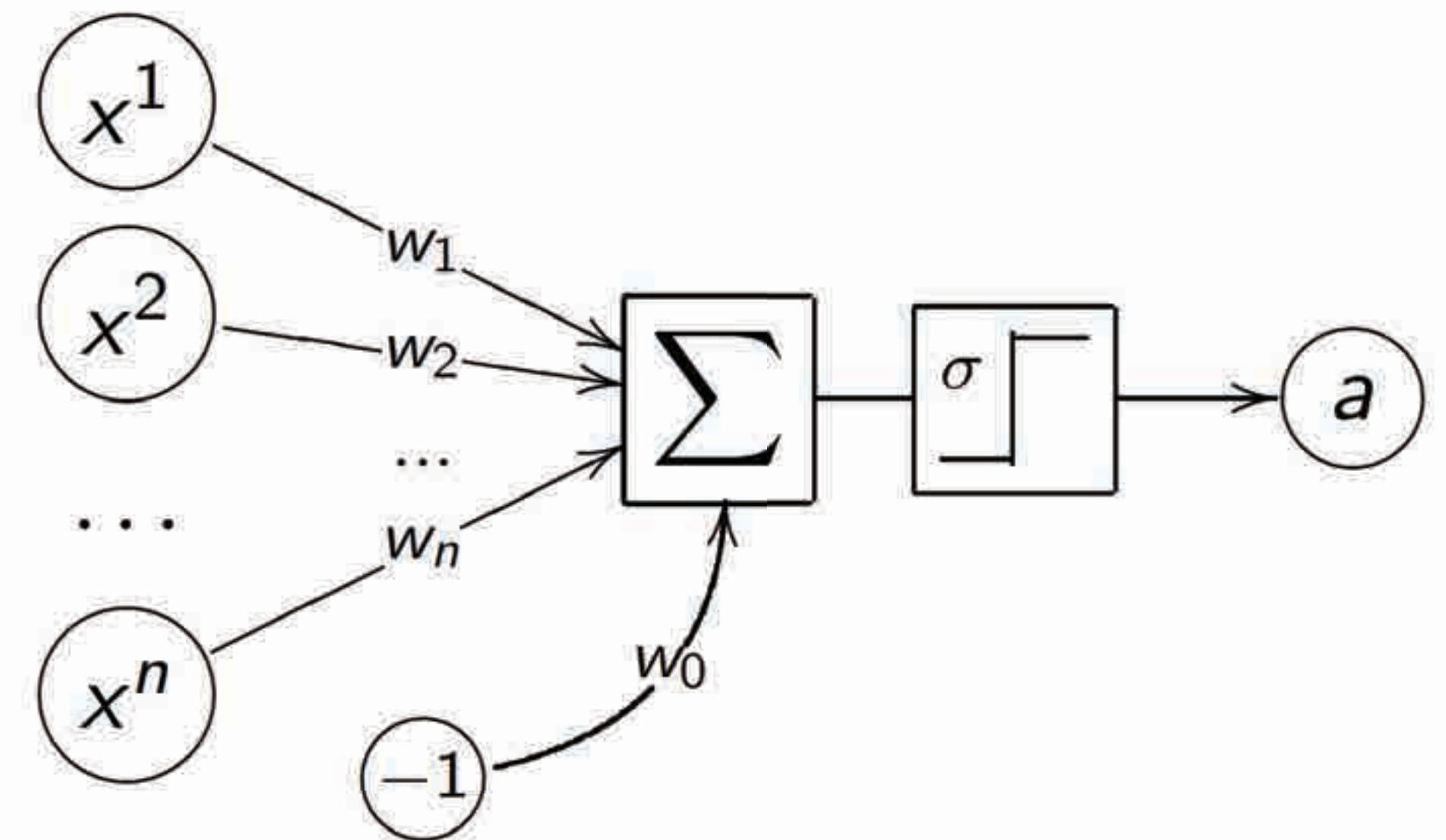
$$a(x) = [\langle x, w \rangle > 0]$$

Исправление вектора весов в случае ошибки:

Если $a(x_i, w) = y_i$, то w менять не нужно

Если $a(x_i, w) = 0$, а $y_i = 1$, то $w_j := w_j + h \cdot f_j(x_i)$

Если $a(x_i, w) = 1$, а $y_i = 0$, то $w_j := w_j - h \cdot f_j(x_i)$



Математика

$$2 \times 2 = 4$$



Математические дисциплины

Математический анализ

Теория алгоритмов

Теория информации

Линейная алгебра

Функциональный анализ

Теория вероятностей

Математическая статистика

Оптимизация