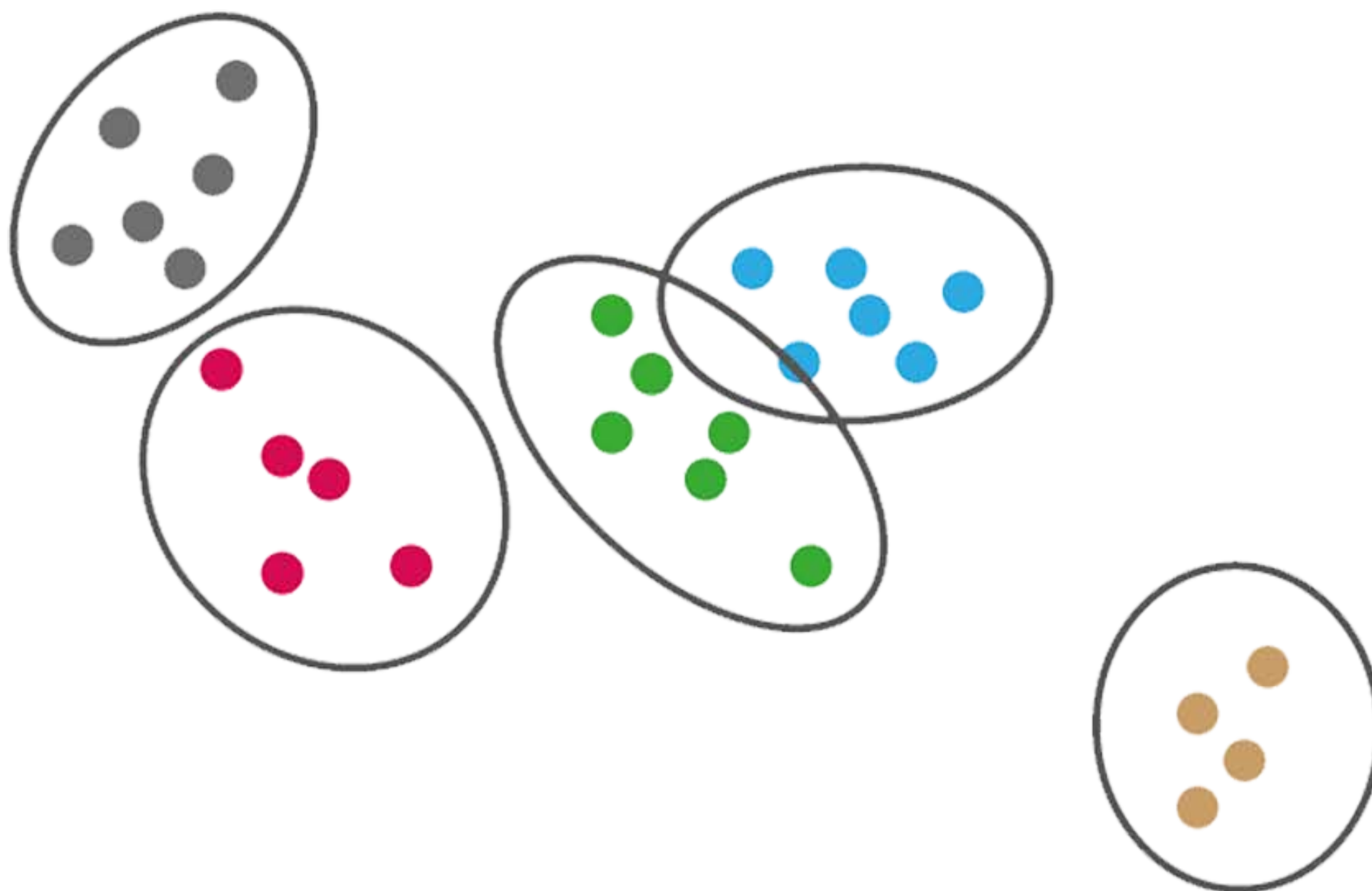


# ЗАДАЧА ОБНАРУЖЕНИЯ АНОМАЛИЙ

---

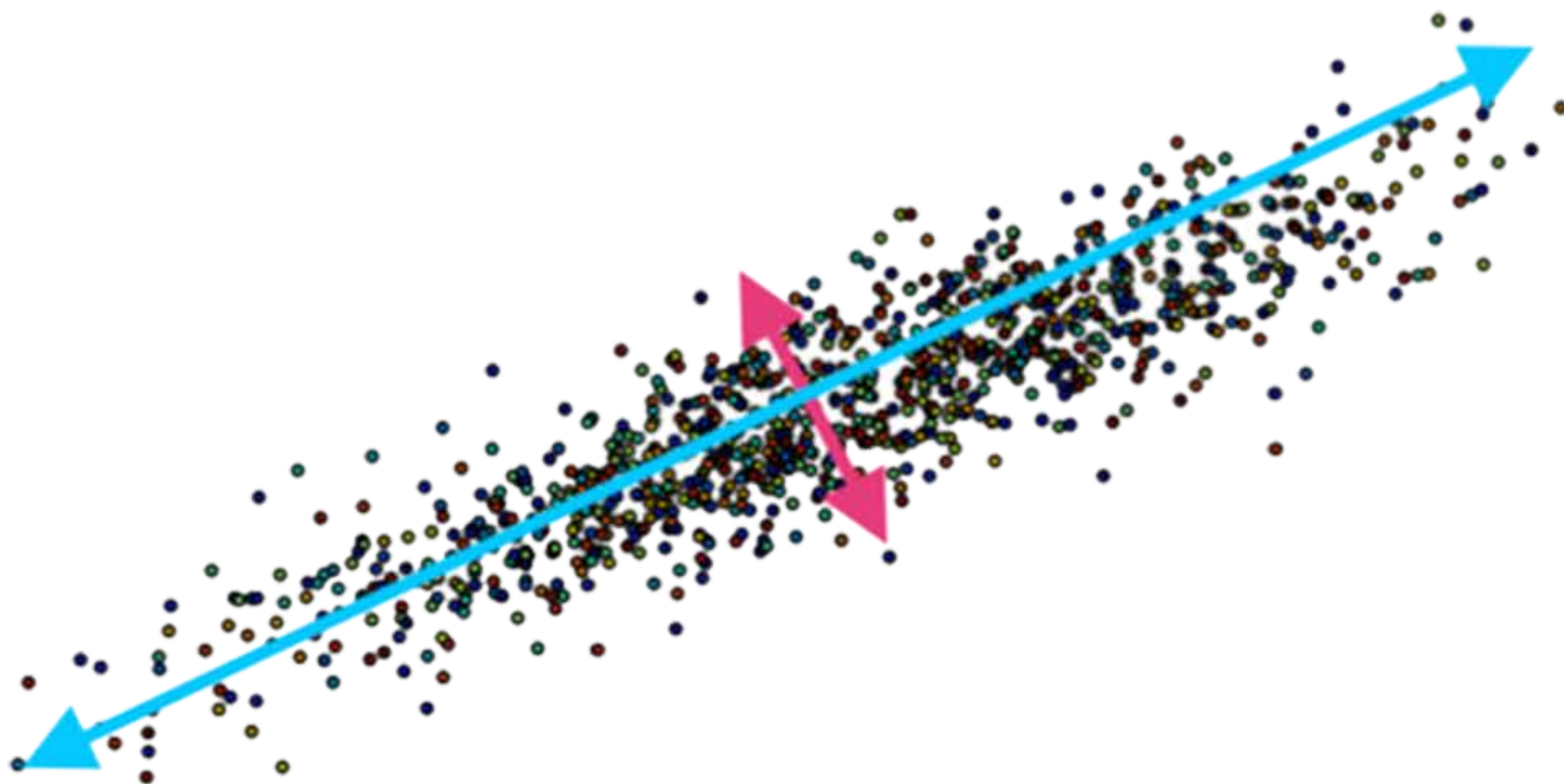
# КЛАСТЕРИЗАЦИЯ

---



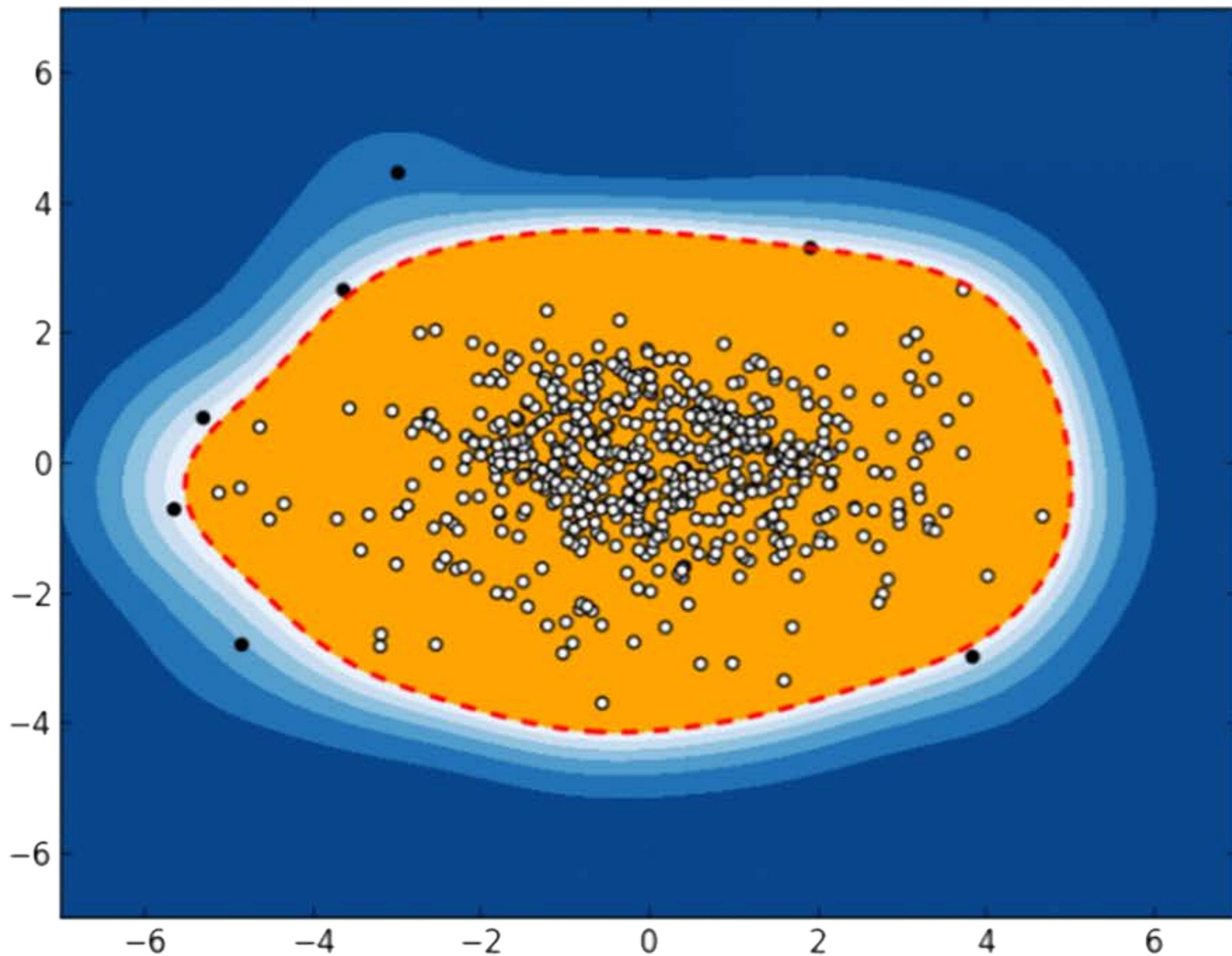
# П О Н И Ж Е Н И Е   Р А З М Е Р Н О С Т И

---



# ПОИСК АНОМАЛИЙ

---



# ПОИСК АНОМАЛИЙ

---

➤ Похож ли новый объект на остальных?



# ПОИСК АНОМАЛИЙ

---

- Объект — клиент банка в текущий момент времени
- Не выделяется ли его поведение?
- Не мошенник ли это?

# ПОИСК АНОМАЛИЙ

---

- Объект — показатели сложной компьютерной системы
- Загрузка процессоров, памяти, сети и т.д.
- Отличается ли текущее состояние системы от тех, которые мы наблюдали ранее?

# ПОИСК АНОМАЛИЙ

---

- Задача: определение тональности отзыва на банк
- Можно ли к новому отзыву применять модель, обученную на прошлых данных?
- Не изменилось ли распределение признаков?



# МЕТОДЫ ОБНАРУЖЕНИЯ АНОМАЛИЙ

---

- Методы, основанные на восстановлении плотности
- Методы, основанные на классификации

# РЕЗЮМЕ

---

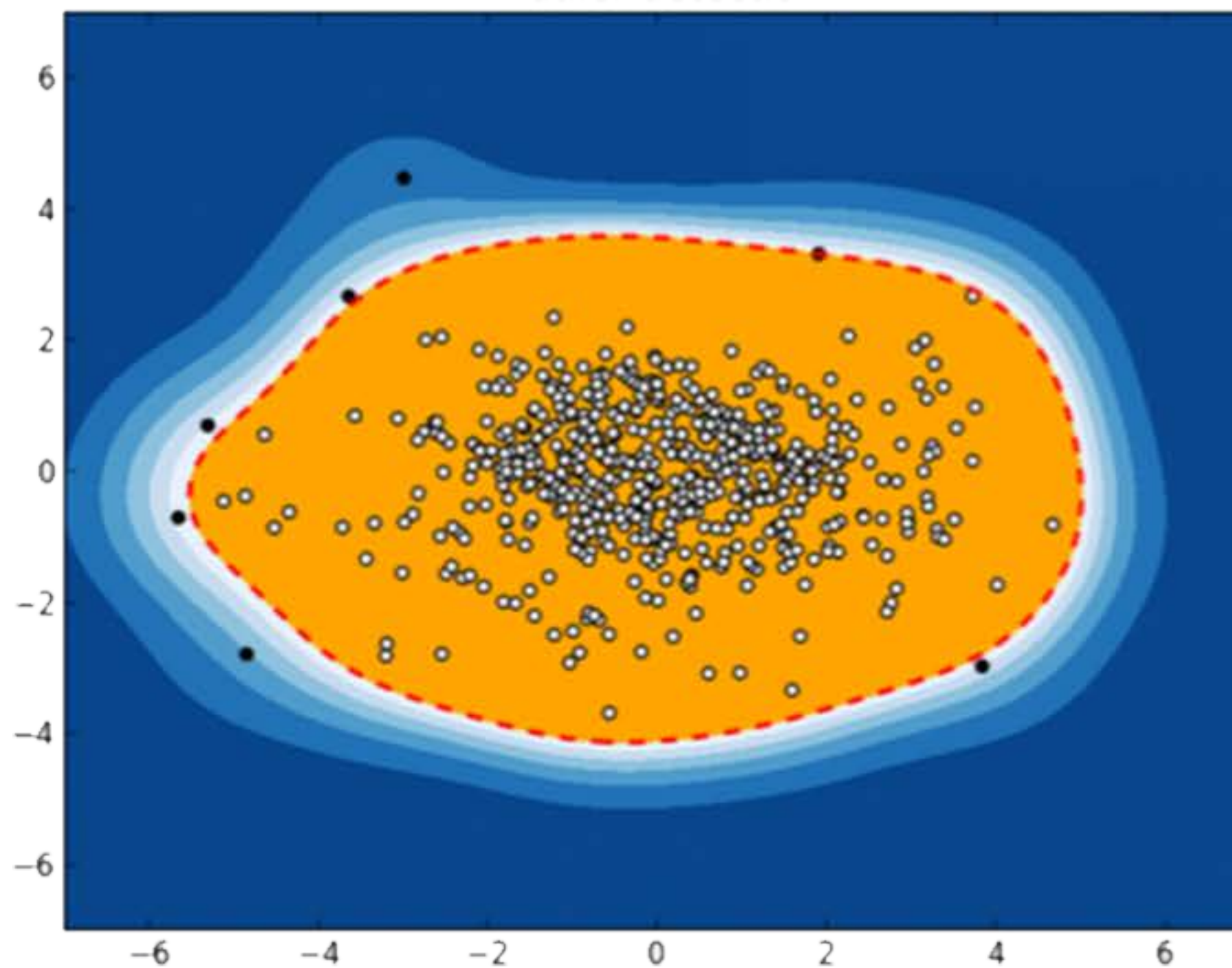
- Поиск аномалий — обнаружение объектов, которые существенно отличаются от других

# ПАРАМЕТРИЧЕСКОЕ ВОССТАНОВЛЕНИЕ ПЛОТНОСТИ

---

# ПОИСК АНОМАЛИЙ

---



# АНОМАЛИЯ

---

- Объект, который получен из другого распределения на пространстве объектов
- Как восстановить распределение?

# ВОССТАНОВЛЕНИЕ РАСПРЕДЕЛЕНИЙ

---

- › Параметрические методы
- › Непараметрические методы
- › Восстановление смесей



# ПАРАМЕТРИЧЕСКИЙ ПОДХОД

---

›  $p(x) = \phi(x|\theta)$

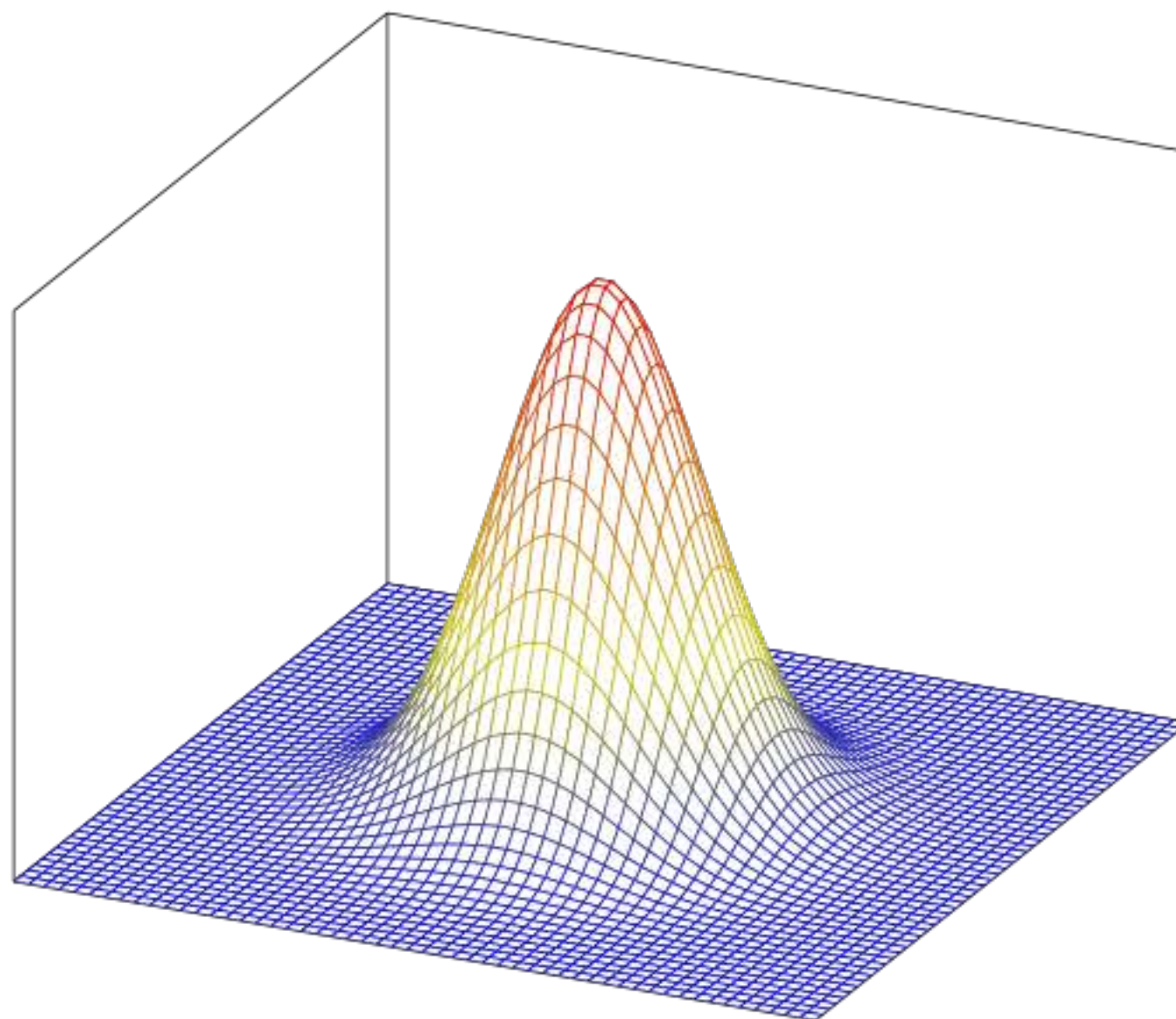
›  $\theta$  — параметры распределения

# ПРИМЕР

---

$$\phi(x|\theta) = \mathcal{N}(\mu, \Sigma)$$

- › Нормальное распределение
- › Параметры:  $\theta = (\mu, \Sigma)$



# ОБУЧЕНИЕ

---

- › Нужно подобрать параметры распределения  $\theta$
- › Вероятность у объектов из выборки должна быть как можно выше

# МЕТОД МАКСИМАЛЬНОГО ПРАВДОПОДОБИЯ

---

$$\sum_{i=1}^{\ell} \log \phi(x_i | \theta) \rightarrow \max_{\theta}$$

- Для некоторых распределений решается аналитически

# МЕТОД МАКСИМАЛЬНОГО ПРАВДОПОДОБИЯ

---

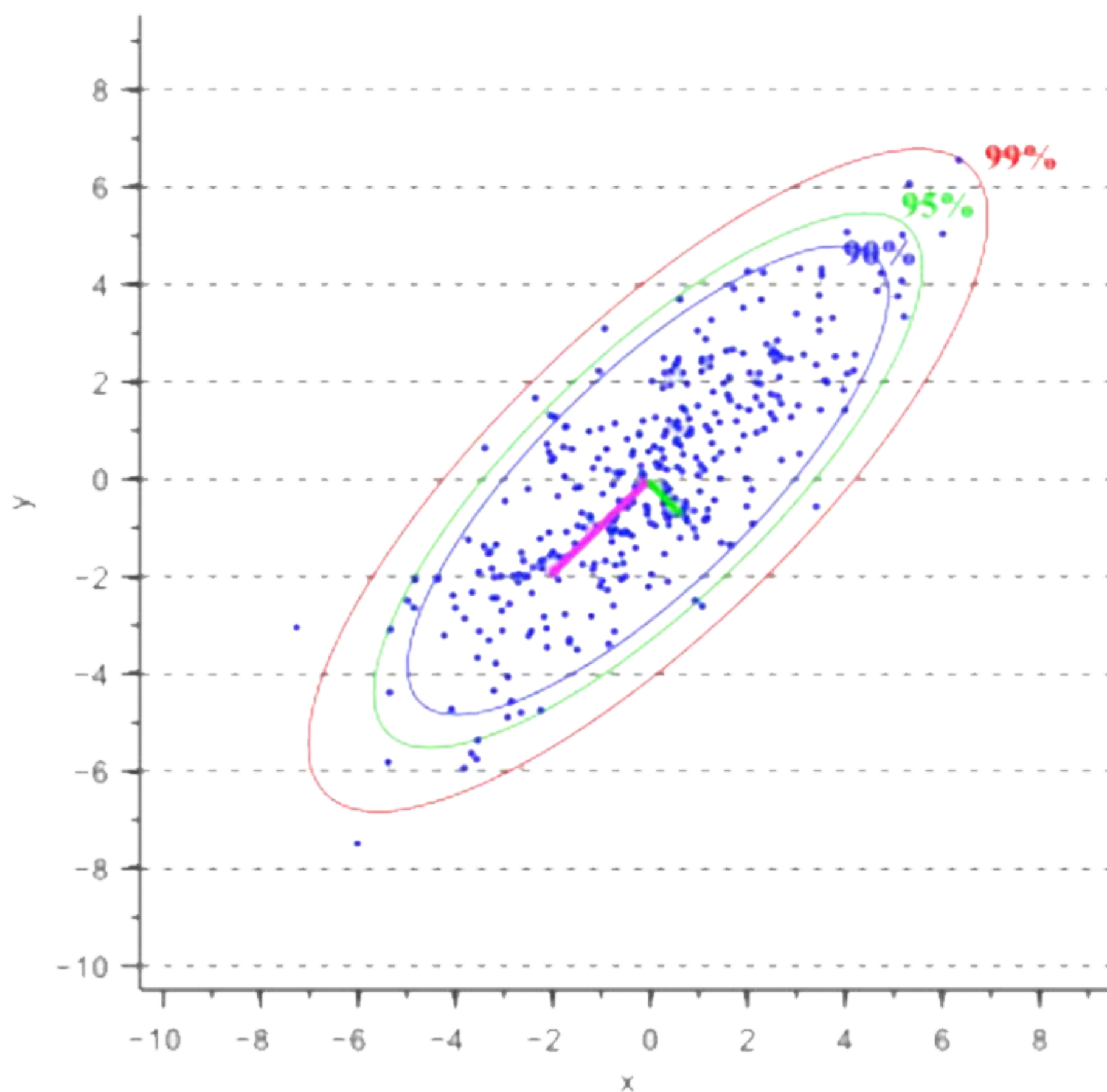
➤ Для нормального распределения:

$$\mu = \frac{1}{\ell} \sum_{i=1}^{\ell} x_i$$

$$\Sigma = \frac{1}{\ell} \sum_{i=1}^{\ell} (x_i - \mu)(x_i - \mu)^T$$

# НОРМАЛЬНОЕ РАСПРЕДЕЛЕНИЕ

---





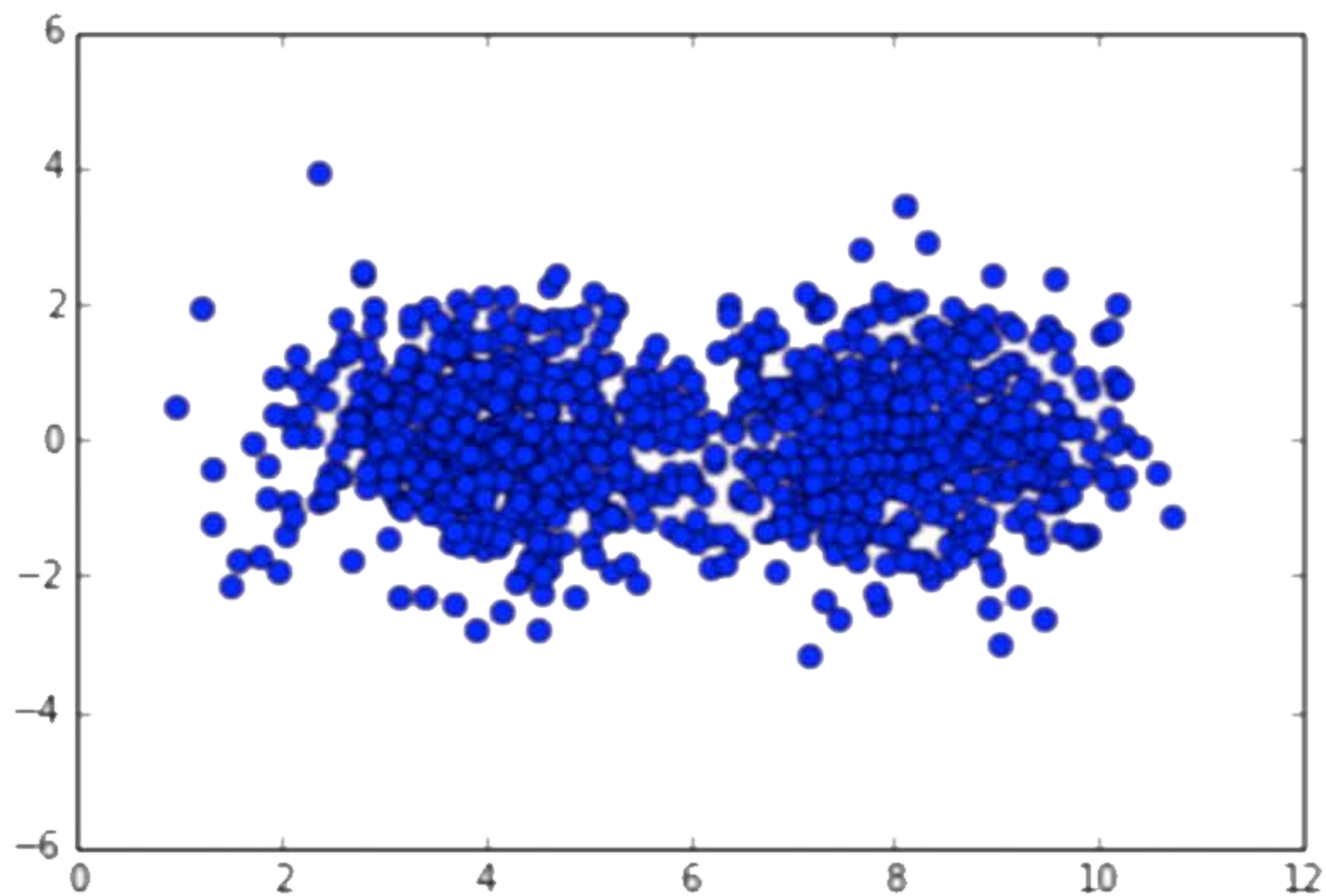
# ОБНАРУЖЕНИЕ АНОМАЛИЙ

---

- › Новый объект  $x$
- › Если  $p(x) < t$ , то это аномалия
- › Порог  $t$ :
  - ▶ Из априорных соображений
  - ▶ По известным аномалиям

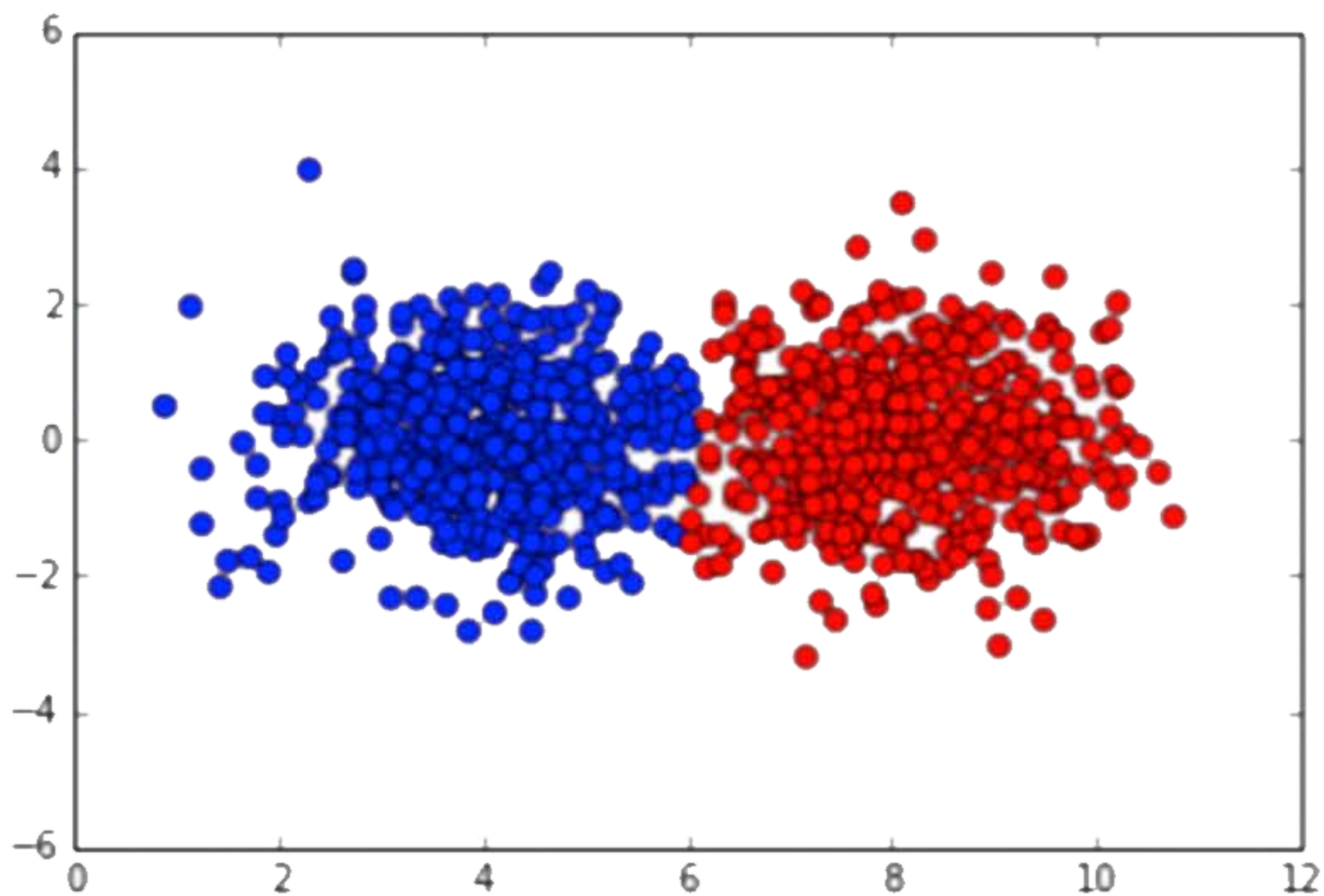
# СМЕСИ РАСПРЕДЕЛЕНИЙ

---



# СМЕСИ РАСПРЕДЕЛЕНИЙ

---



# ЕМ-АЛГОРИТМ

---

$$p(x) = \sum_{j=1}^K w_j p_j(x) \quad p_j(x) = \phi(x|\theta_j)$$

► Е-шаг:

$$g_{ji} = p(j|x_i) = \frac{w_j p_j(x_i)}{p(x_i)}$$

► М-шаг:

$$w_j = \frac{1}{N} \sum_{i=1}^N g_{ji}$$

$$\theta_j = \operatorname{argmax}_{\theta} \sum_{i=1}^N g_{ji} \ln \varphi(\theta; x_i)$$

# РЕЗЮМЕ

---

- Аномалия — объект из другого распределения
- Параметрические методы восстановления распределений

# НЕПАРАМЕТРИЧЕСКОЕ ВОССТАНОВЛЕНИЕ ПЛОТНОСТИ

---

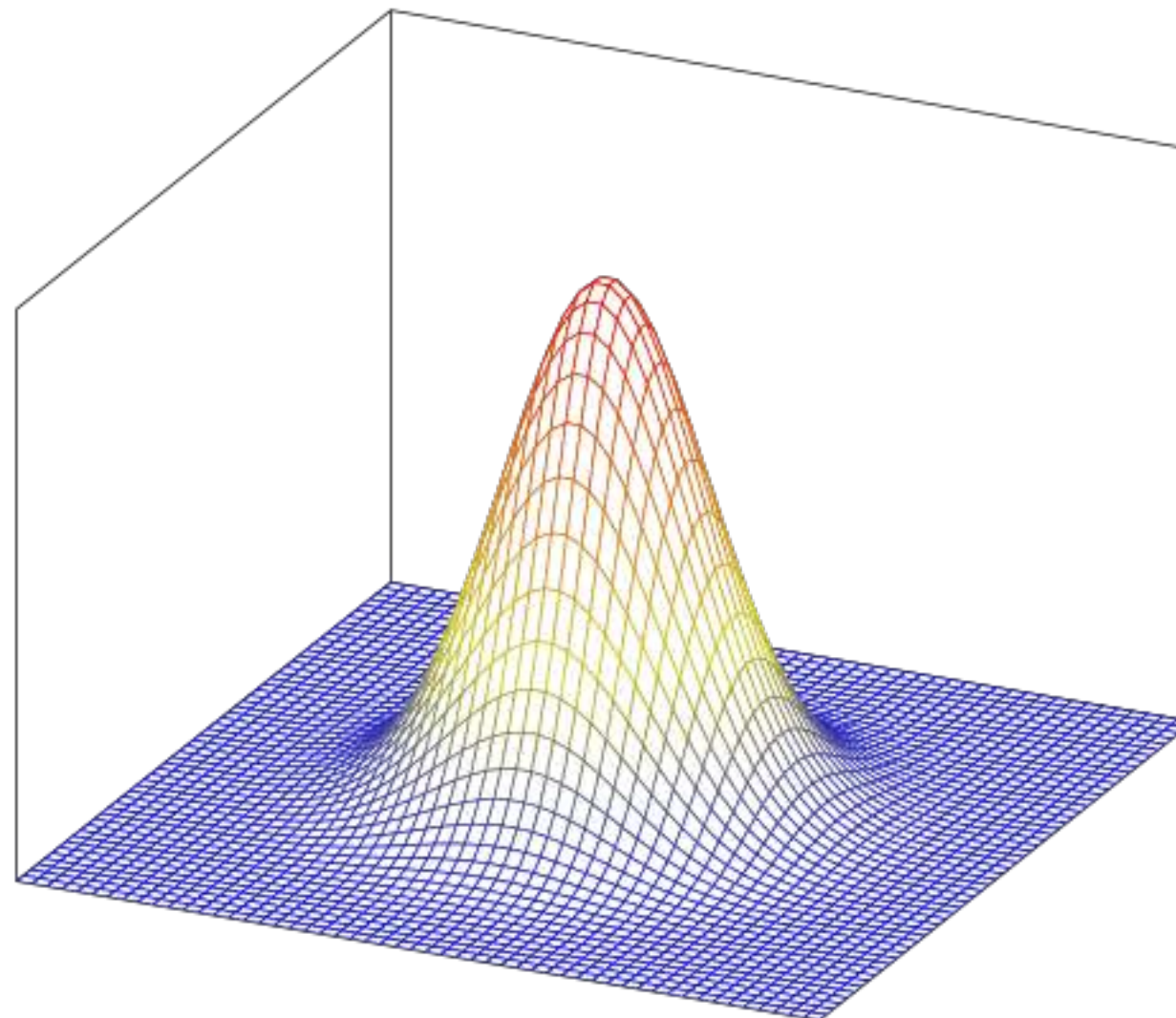


# ПАРАМЕТРИЧЕСКОЕ ВОССТАНОВЛЕНИЕ

---

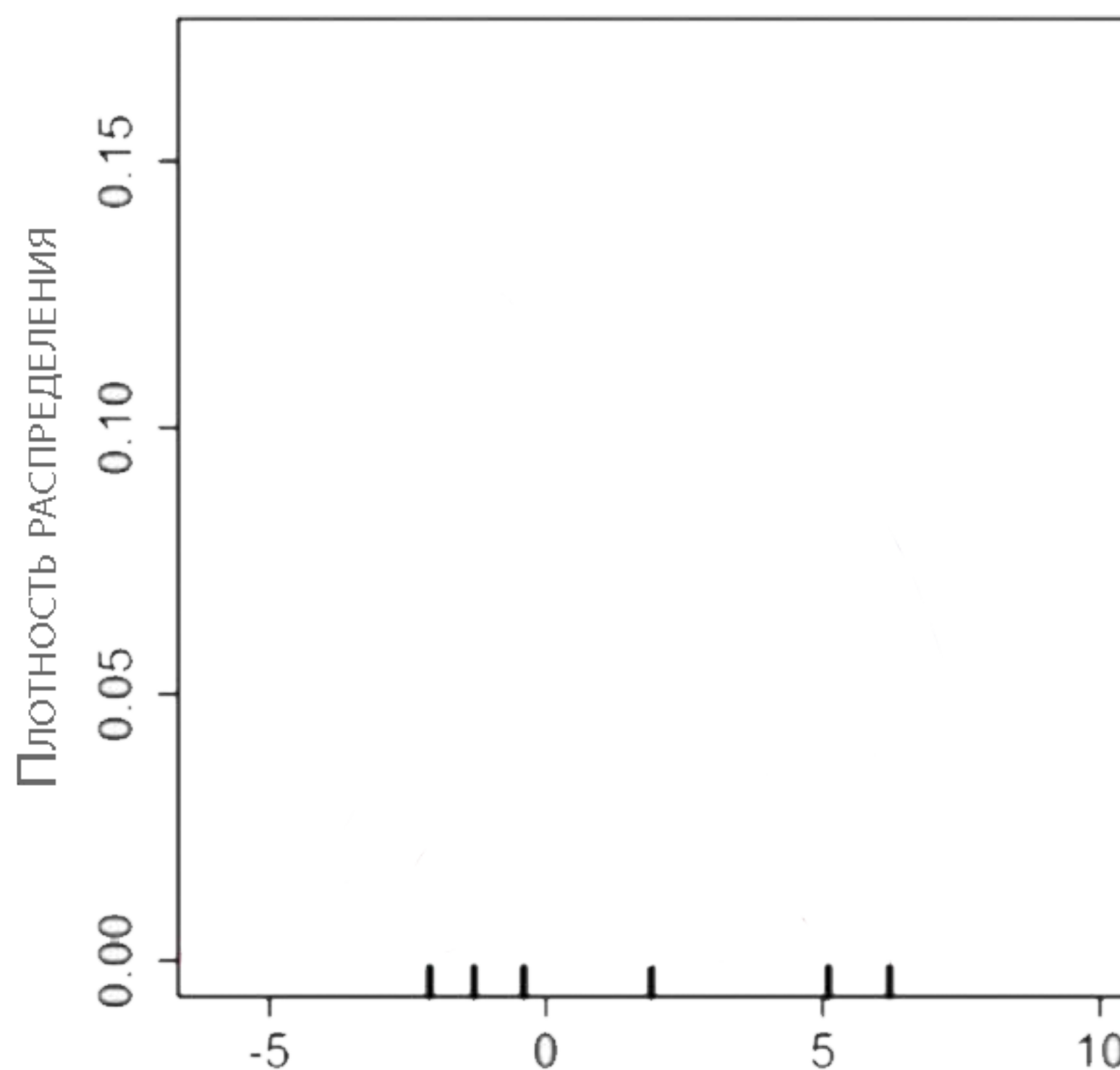
$$\phi(x|\theta) = \mathcal{N}(\mu, \Sigma)$$

- › Нормальное распределение
- › Параметры:  $\theta = (\mu, \Sigma)$



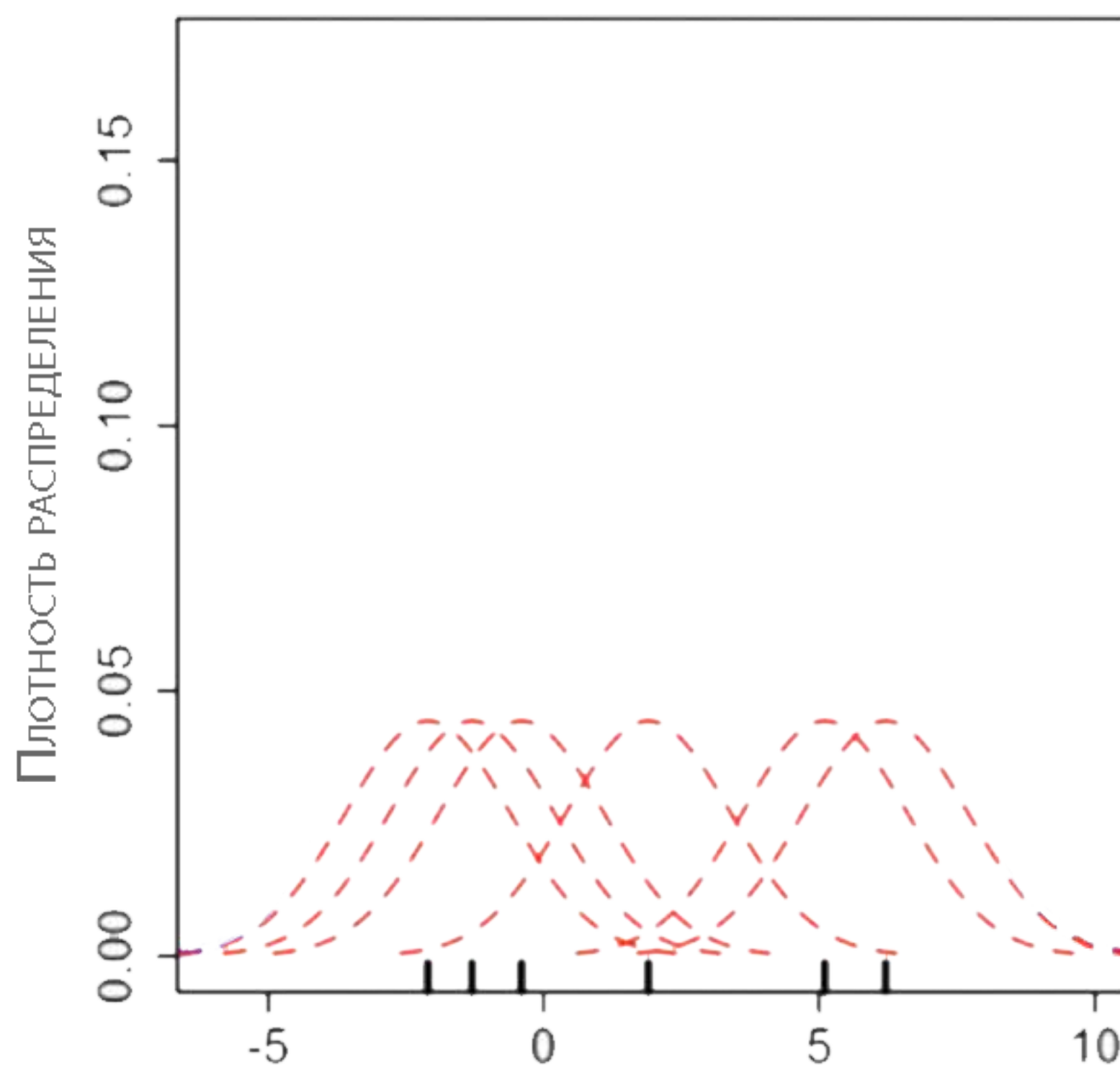
# НЕПАРАМЕТРИЧЕСКОЕ ВОССТАНОВЛЕНИЕ

---



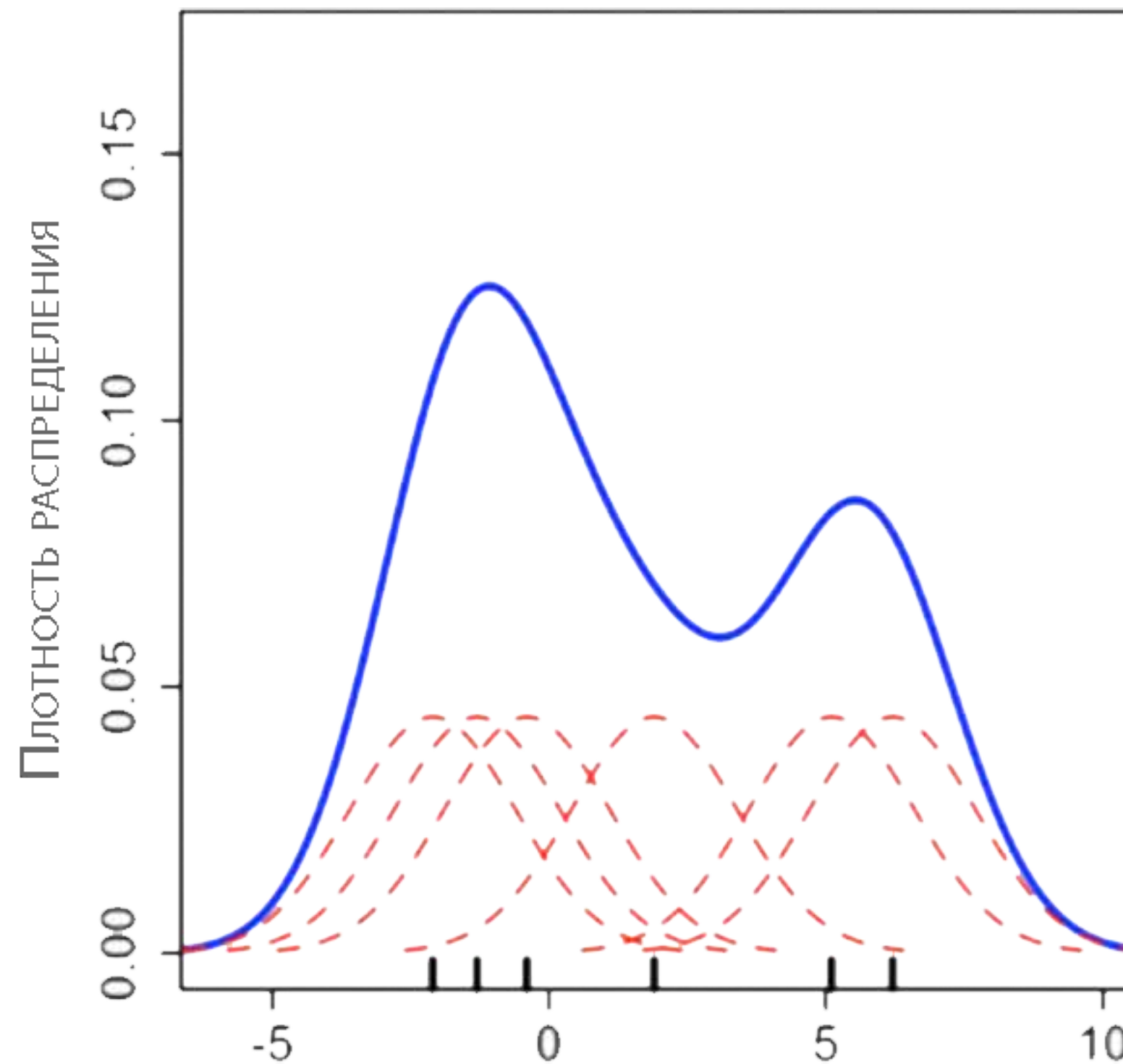
# НЕПАРАМЕТРИЧЕСКОЕ ВОССТАНОВЛЕНИЕ

---



# НЕПАРАМЕТРИЧЕСКОЕ ВОССТАНОВЛЕНИЕ

---



# ОЦЕНКА ПАРЗЕНА-РОЗЕНБЛАТТА

---

$$p_h(x) = \frac{1}{\ell h} \sum_{i=1}^{\ell} K\left(\frac{x - x_i}{h}\right)$$

# ОЦЕНКА ПАРЗЕНА-РОЗЕНБЛАТТА

---

$$p_h(x) = \frac{1}{\ell h} \sum_{i=1}^{\ell} K\left(\frac{x - x_i}{h}\right)$$

Ядро



# ОЦЕНКА ПАРЗЕНА-РОЗЕНБЛАТТА

---

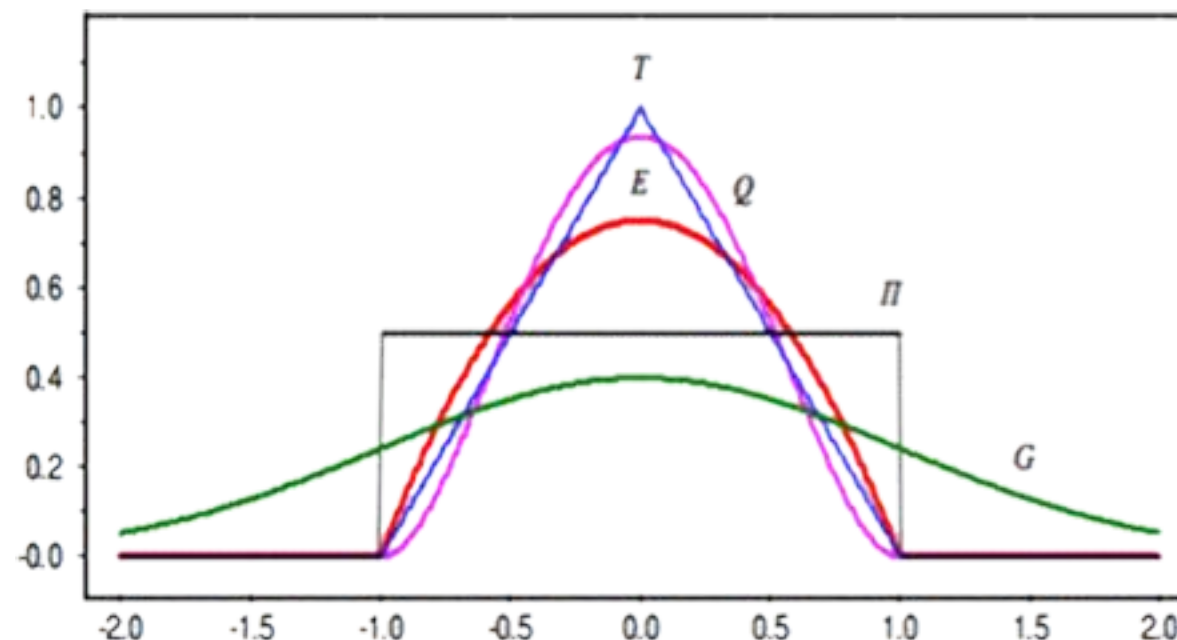
$$p_h(x) = \frac{1}{\ell h} \sum_{i=1}^{\ell} K\left(\frac{x - x_i}{h}\right)$$

- ›  $K(r)$  — ядро
- › Чётная функция
- ›  $\int K(r) dr = 1$

# ЯДРО

---

- ›  $E(r) = \frac{3}{4}(1 - r^2)[|r| \leq 1]$  — оптимальное
- ›  $Q(r) = \frac{15}{16}(1 - r^2)^2[|r| \leq 1]$  — кватрическое
- ›  $T(r) = (1 - |r|)[|r| \leq 1]$  — треугольное
- ›  $G(r) = (2\pi)^{-\frac{1}{2}} \exp(-\frac{1}{2}r^2)$  — гауссовское
- ›  $\Pi(r) = \frac{1}{2}[|r| \leq 1]$  — прямое



# ОЦЕНКА ПАРЗЕНА-РОЗЕНБЛАТТА

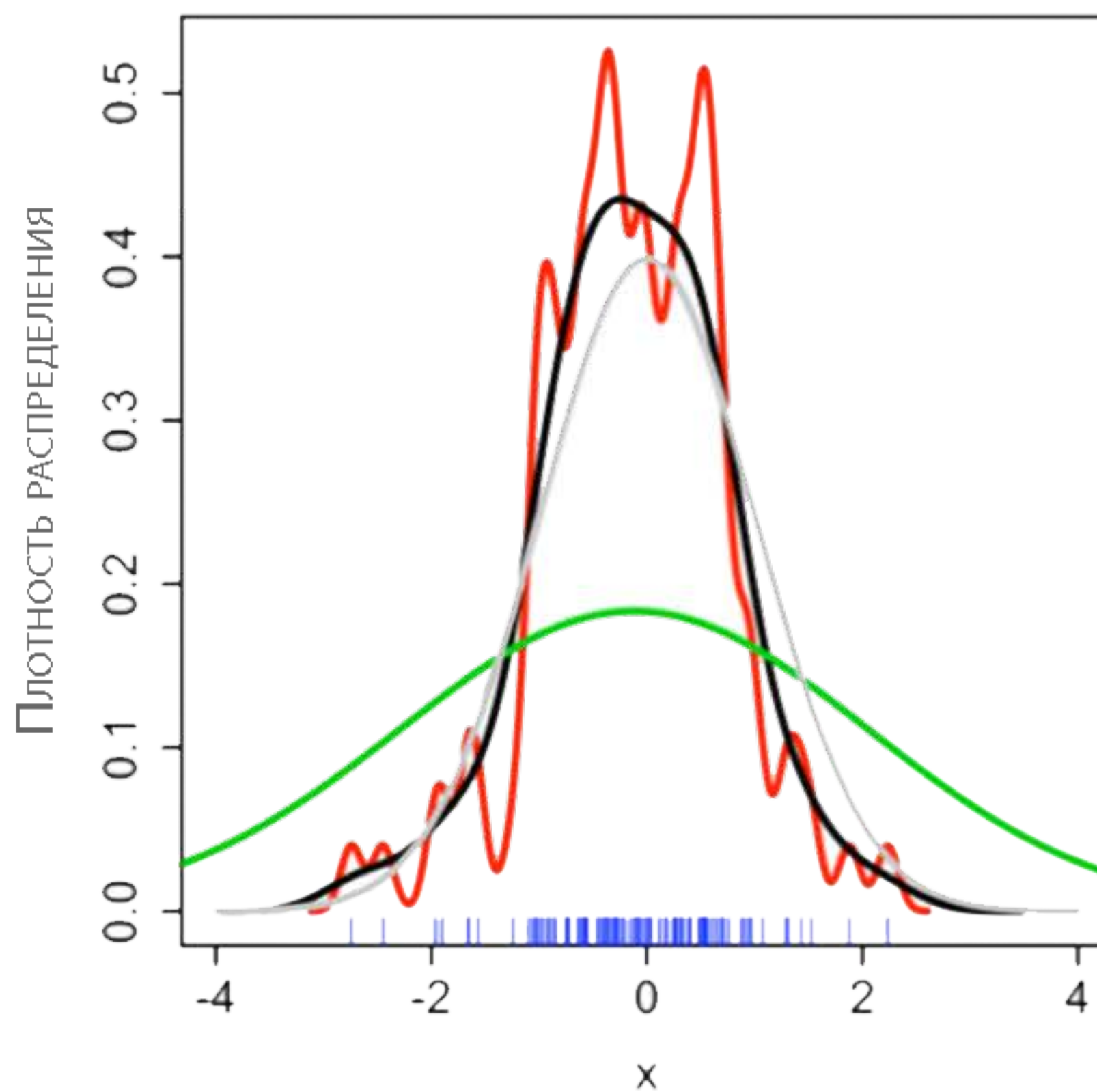
---

$$p_h(x) = \frac{1}{\ell h} \sum_{i=1}^{\ell} K\left(\frac{x - x_i}{h}\right)$$

»  $h$  — ширина окна

# ШИРИНА ОКНА

---



# МНОГОМЕРНЫЙ СЛУЧАЙ

---

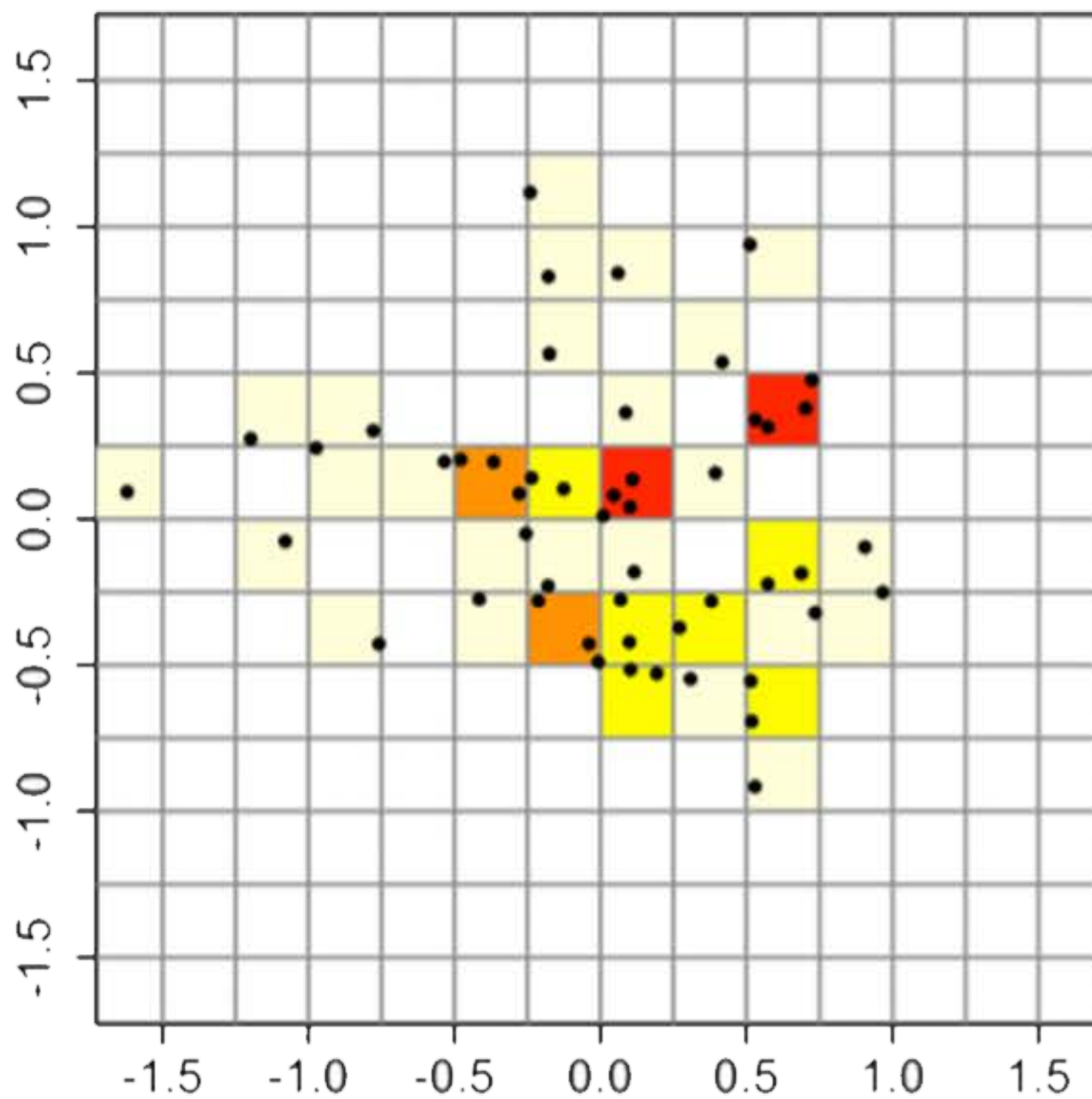
$$p_h(x) = \frac{1}{\ell V(h)} \sum_{i=1}^{\ell} K\left(\frac{\rho(x, x_i)}{h}\right)$$

$$\blacktriangleright V(h) = \int K\left(\frac{\rho(x, x_i)}{h}\right) dx \quad \text{—}$$

нормировочная константа

# МНОГОМЕРНЫЙ СЛУЧАЙ

- » Чем выше размерность, тем больше объектов нужно





# МНОГОМЕРНЫЙ СЛУЧАЙ

---

- Число объектов, необходимых для качественного оценивания, растёт экспоненциально с ростом размерности

# ОБНАРУЖЕНИЕ АНОМАЛИЙ

---

- › Новый объект  $x$
- › Если  $p(x) < t$ , то это аномалия
- › Порог  $t$ :
  - ▶ Из априорных соображений
  - ▶ По известным аномалиям

# РЕЗЮМЕ

---

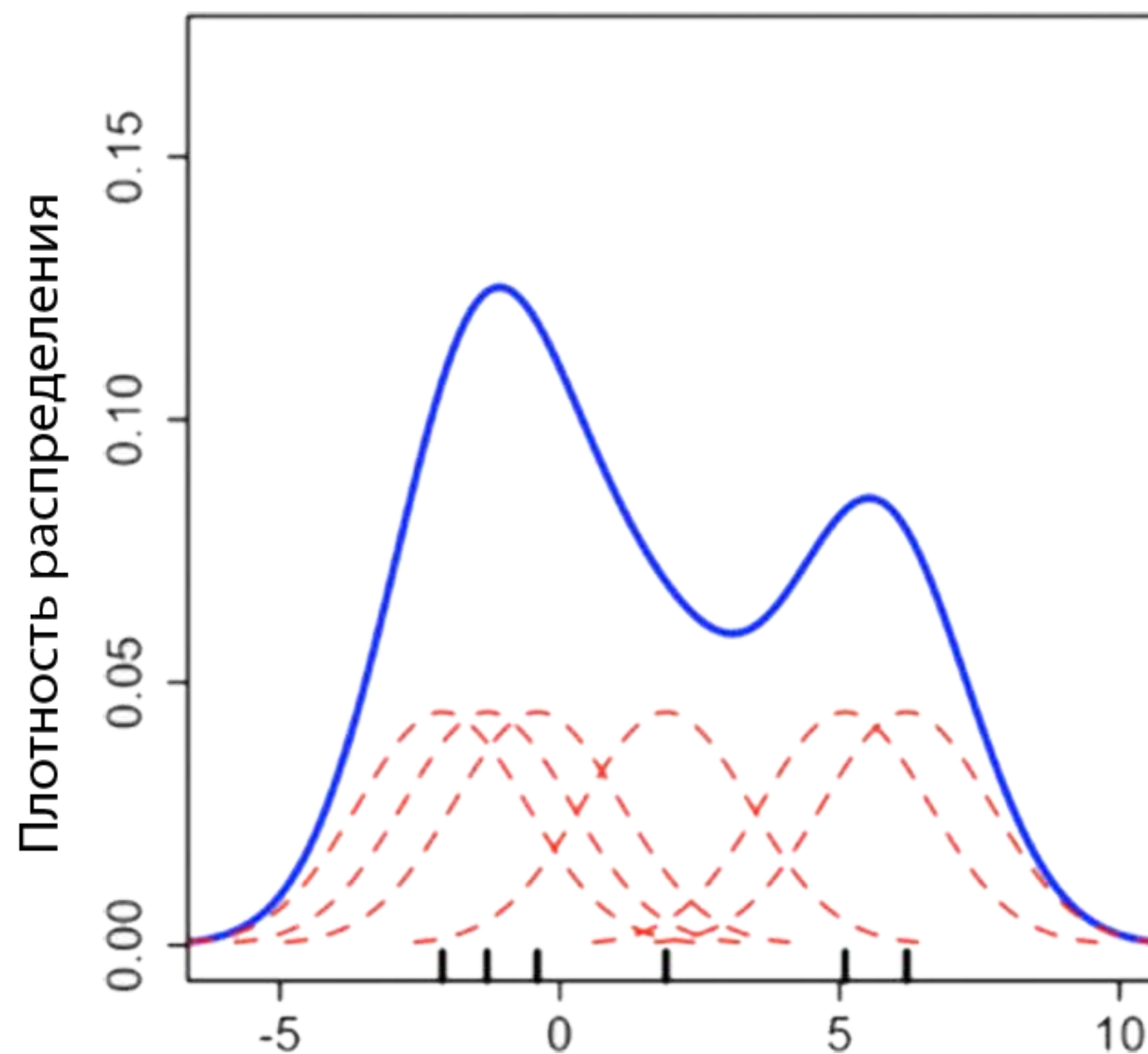
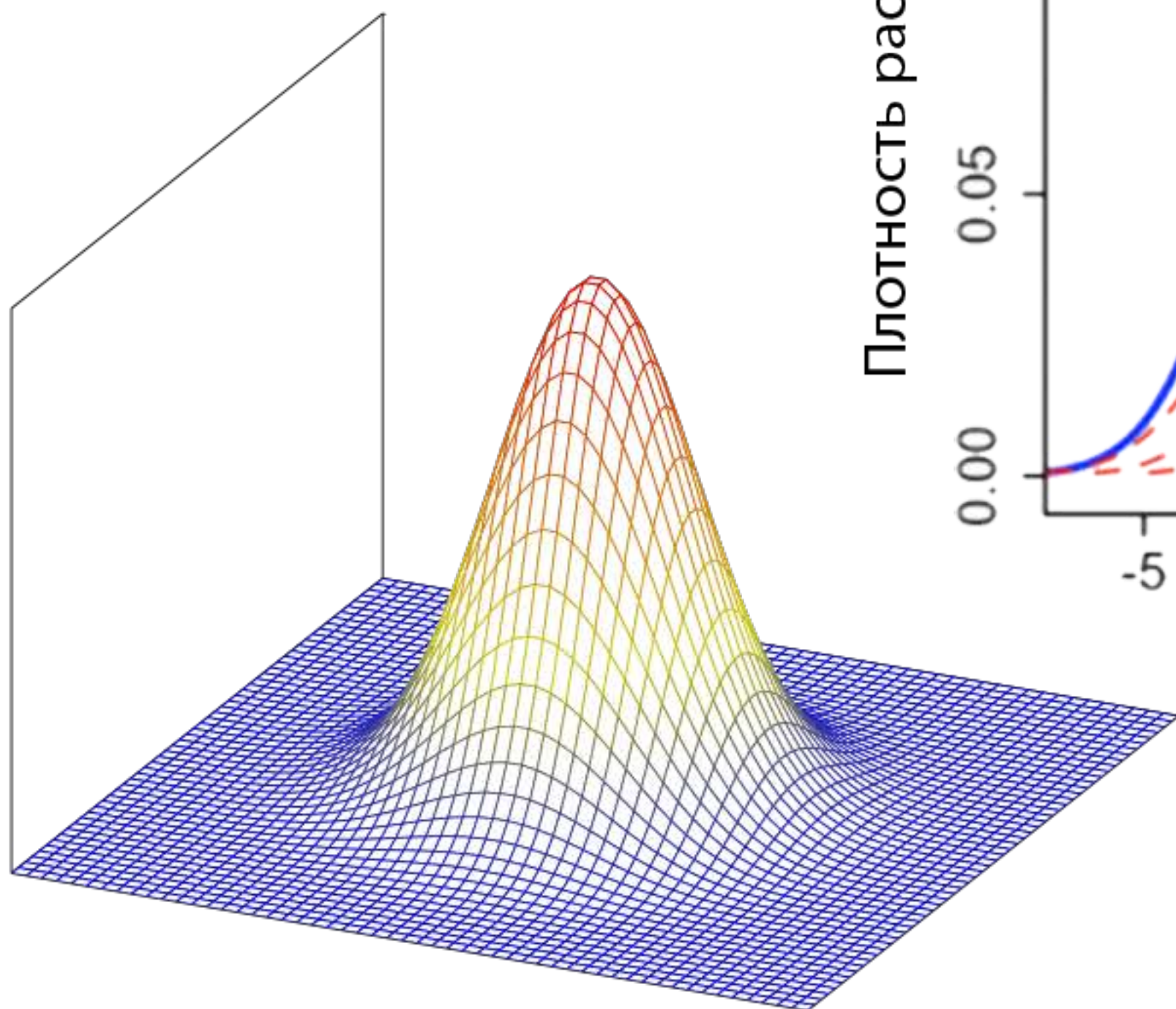
- › Непараметрический подход — для восстановления сложных плотностей
- › Параметры: ядро и ширина окна
- › В многомерно случае требуется большая выборка

# ОДНОКЛАССОВЫЙ SVM

---

# ВОССТАНОВЛЕНИЕ ПЛОТНОСТИ

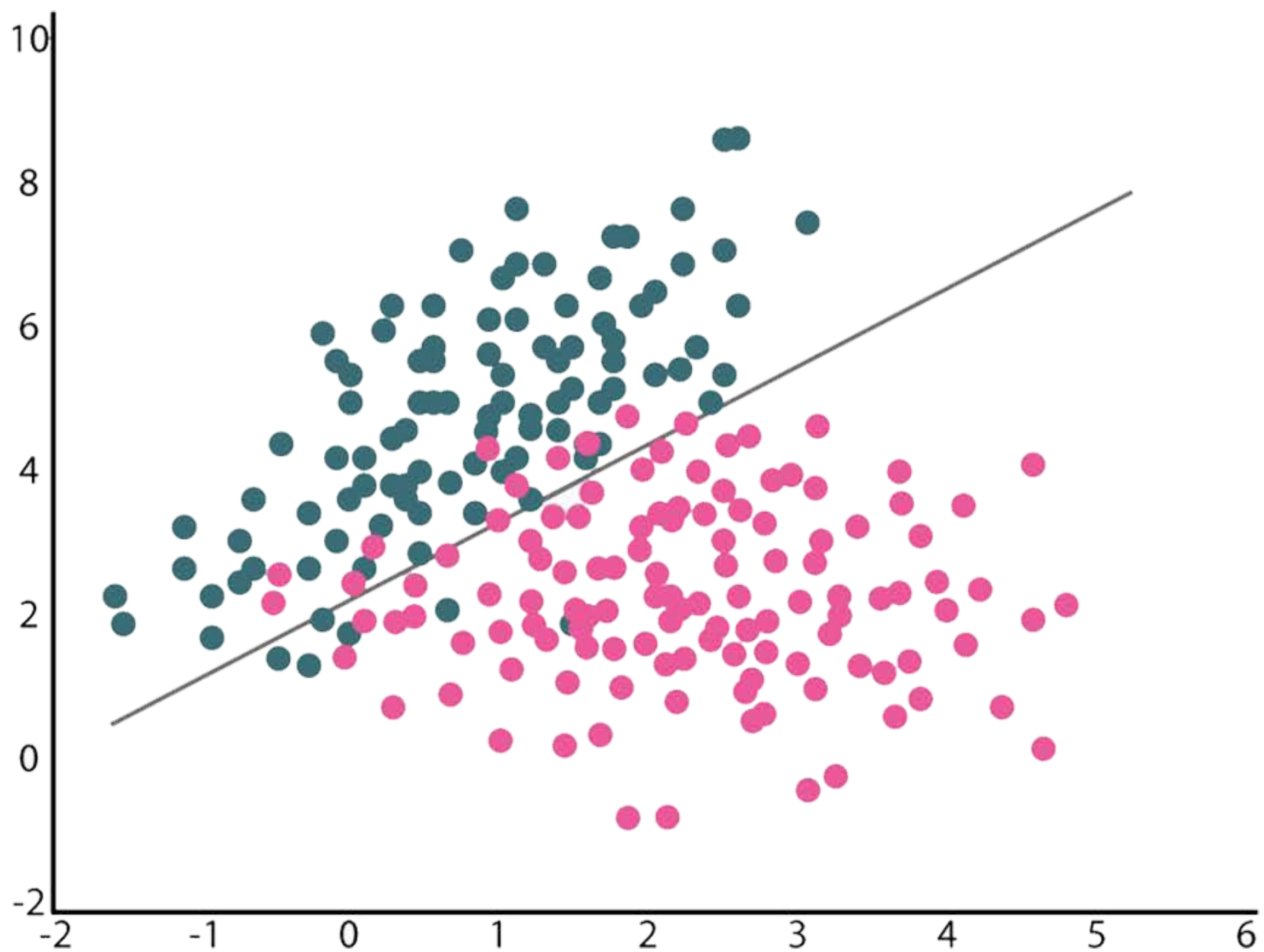
---





# КЛАССИФИКАЦИЯ

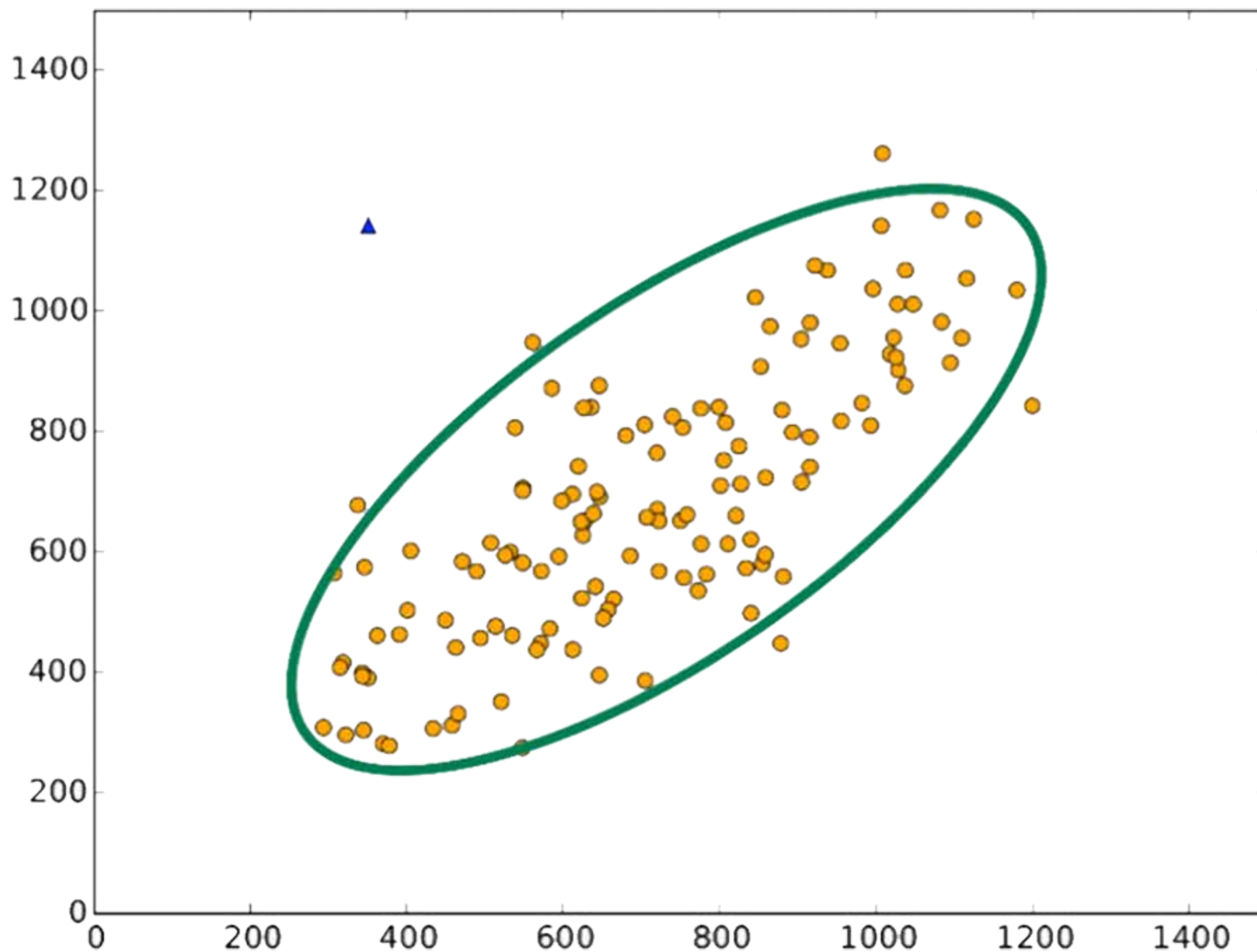
---





# ОБНАРУЖЕНИЕ АНОМАЛИЙ

---



# ОБНАРУЖЕНИЕ АНОМАЛИЙ

---

- Обучающая выборка — нормальные объекты
- Начало координат — аномалия
- Задача: отделить выборку гиперплоскостью от нуля
- Максимизация отступа

# ОДНОКЛАССОВЫЙ SVM

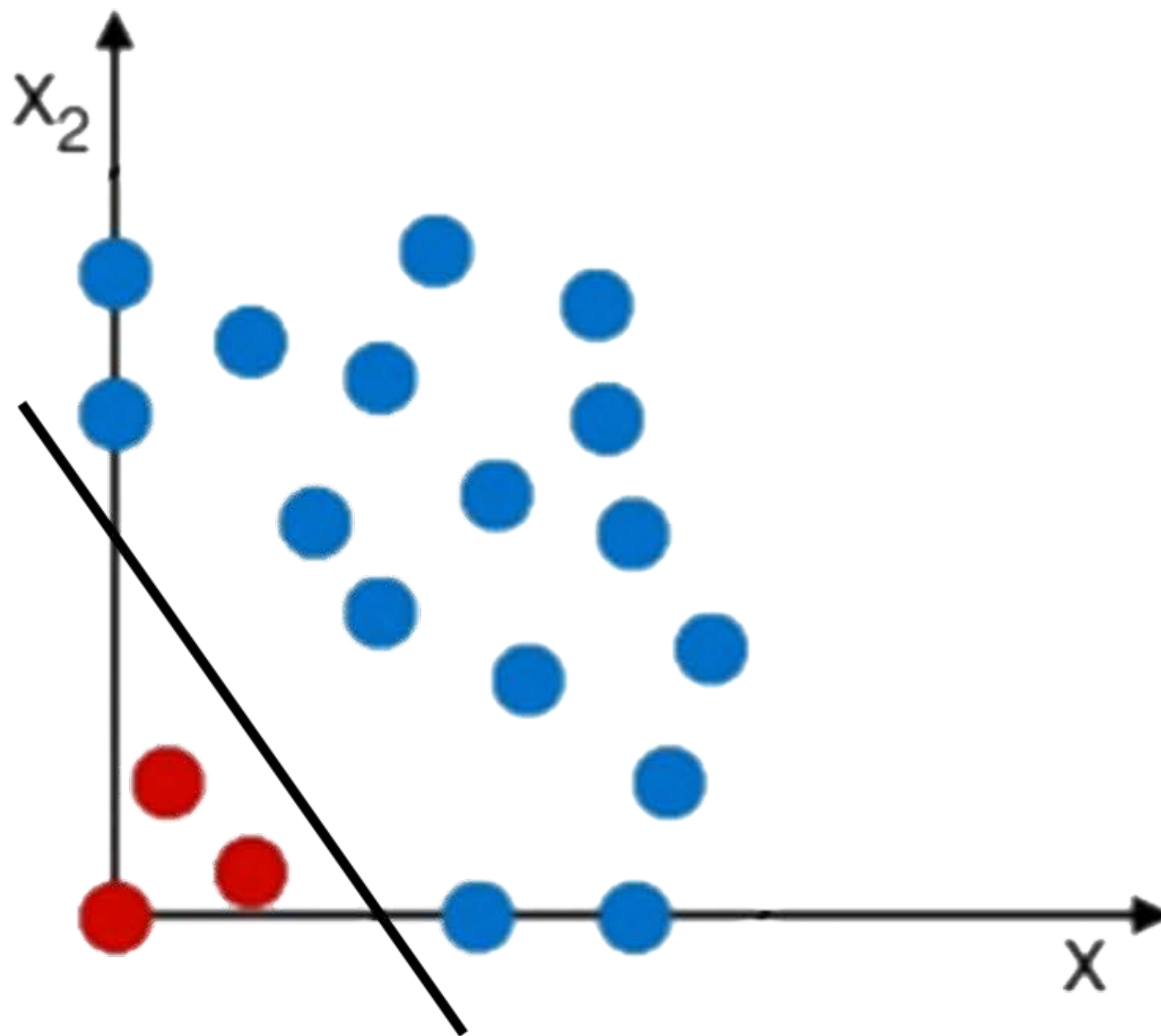
---

$$\begin{cases} \frac{1}{2} \| \mathbf{w} \|^2 + \frac{1}{v\ell} \sum_{i=1}^{\ell} \xi_i - \rho \rightarrow \min_{\mathbf{w}, \xi, \rho} \\ \langle \mathbf{w}, x_i \rangle \geq \rho - \xi_i, \quad \xi_i \geq 0 \end{cases}$$

»  $v$  — верхняя оценка на долю аномалий на выборке

# ОДНОКЛАССОВЫЙ SVM

---



# ЯДРОВОЙ ПЕРЕХОД

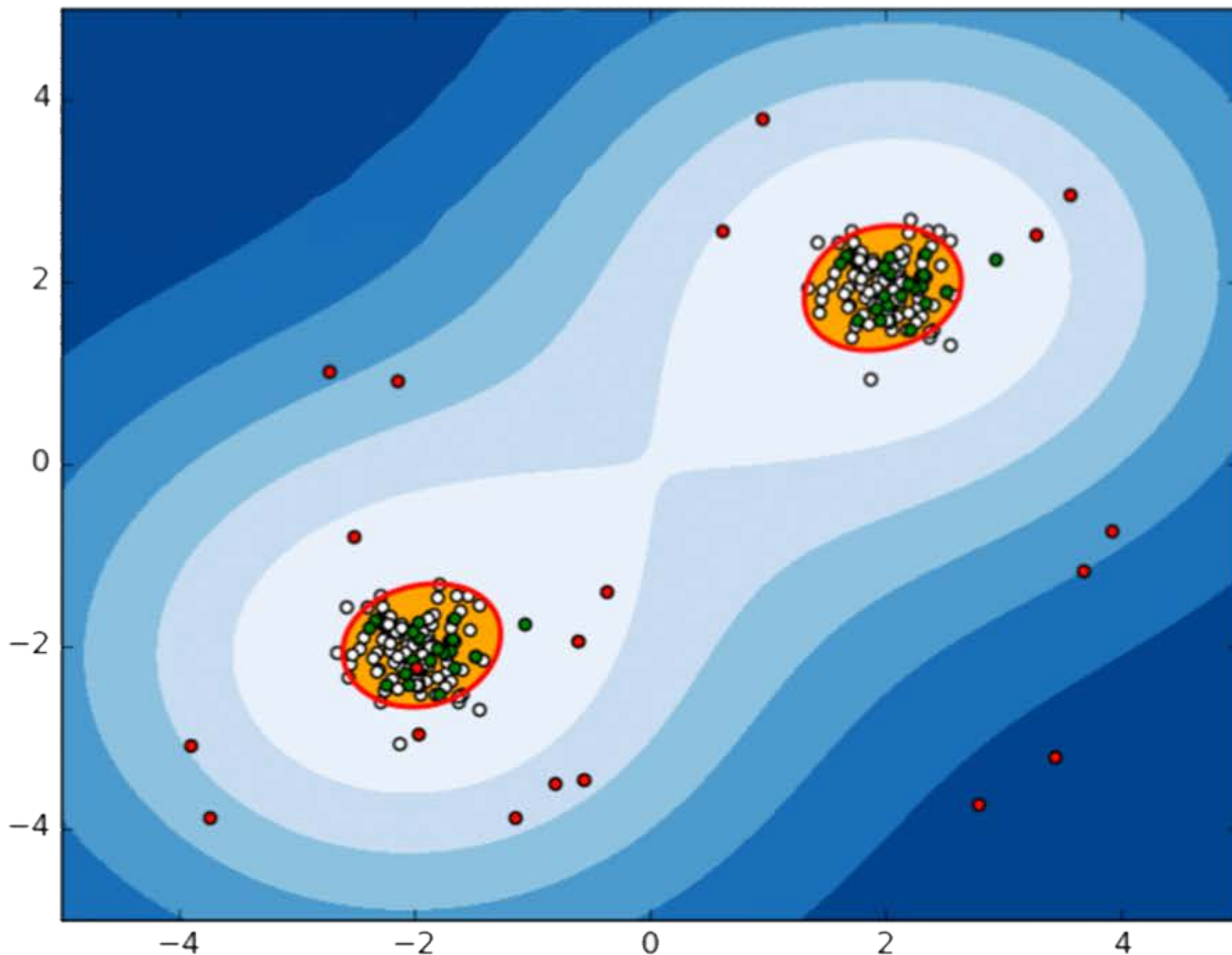
---

- › SVM позволяет строить нелинейные разделяющие поверхности
- › Ядровой переход (kernel trick)
- › Популярный выбор: RBF-ядро

$$K(x, z) = \exp\left(\frac{\|x - z\|^2}{\sigma^2}\right)$$



# ЯДРОВОЙ ПЕРЕХОД



# РЕЗЮМЕ

---

- Одноклассовый SVM отделяет выборку от начала координат
- Имеет смысл при использовании ядер