

ЗАДАЧА КЛАСТЕРИЗАЦИИ

РАНЕЕ: ОБУЧЕНИЕ НА РАЗМЕЧЕННЫХ ДАННЫХ (SUPERVISED LEARNING)

» Обучающая выборка:

x_1, \dots, x_ℓ — объекты

y_1, \dots, y_ℓ — ответы

РАНЕЕ: ОБУЧЕНИЕ НА РАЗМЕЧЕННЫХ ДАННЫХ (SUPERVISED LEARNING)

» Обучающая выборка:

x_1, \dots, x_ℓ — объекты

y_1, \dots, y_ℓ — ответы

» Тестовая выборка:

$x_{\ell+1}, \dots, x_{\ell+u}$

РАНЕЕ: ОБУЧЕНИЕ НА РАЗМЕЧЕННЫХ ДАННЫХ (SUPERVISED LEARNING)

» Обучающая выборка:

x_1, \dots, x_ℓ — объекты

y_1, \dots, y_ℓ — ответы

» Тестовая выборка:

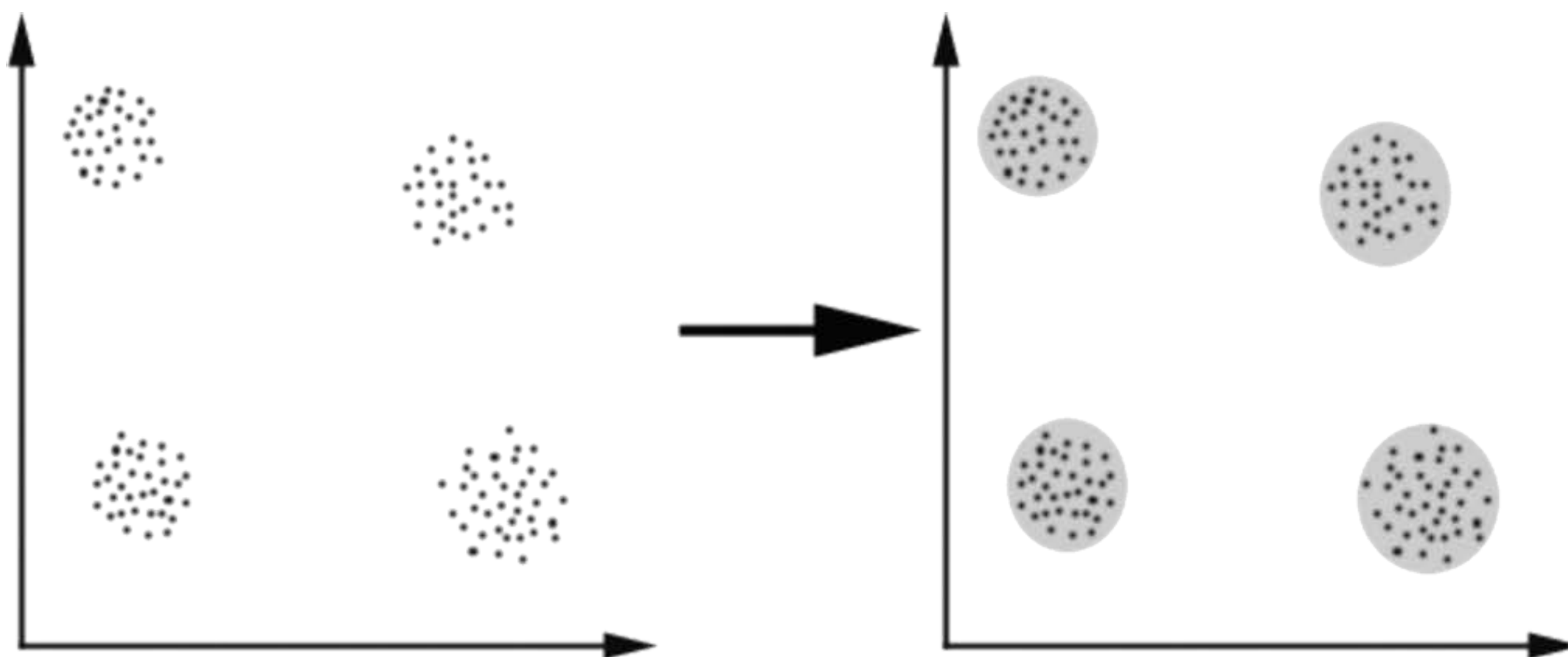
$x_{\ell+1}, \dots, x_{\ell+u}$

» В регрессии: y_i — прогнозируемая величина
В классификации: y_i — метка класса

КЛАСТЕРИЗАЦИЯ

- » «Обучающая» выборка:
 x_1, \dots, x_ℓ — объекты
- » Она же и тестовая
- » Нужно поставить метки y_1, \dots, y_ℓ так, чтобы объекты с одной и той же меткой были похожи, а с разными метками — не очень похожи

КАК ЭТО ВЫГЛЯДИТ



СРЕДНЕЕ ВНУТРИКЛАСТЕРНОЕ РАССТОЯНИЕ

$$F_0 = \frac{\sum_{i < j} [y_i = y_j] \rho(x_i, x_j)}{\sum_{i < j} [y_i = y_j]} \rightarrow \min$$

СРЕДНЕЕ МЕЖКЛАСТЕРНОЕ РАССТОЯНИЕ

$$F_1 = \frac{\sum_{i < j} [y_i \neq y_j] \rho(x_i, x_j)}{\sum_{i < j} [y_i \neq y_j]} \rightarrow \max$$

ПРИДУМЫВАЕМ МЕТРИКУ КАЧЕСТВА

$$F_0 = \frac{\sum_{i < j} [y_i = y_j] \rho(x_i, x_j)}{\sum_{i < j} [y_i = y_j]} \quad F_1 = \frac{\sum_{i < j} [y_i \neq y_j] \rho(x_i, x_j)}{\sum_{i < j} [y_i \neq y_j]}$$

$$\frac{F_0}{F_1} \rightarrow \min$$

ПРИМЕРЫ ЗАДАЧ КЛАСТЕРИЗАЦИИ

ЗАЧЕМ НУЖНЫ РАЗНЫЕ АЛГОРИТМЫ КЛАСТЕРИЗАЦИИ

- Каждые данные в чём-то «особенные»
- Каждая задача кластеризации тоже
- В разных задачах кластеризации могут быть отличия:
 - ▶ Форма кластеров
 - ▶ Необходимость делать кластеры вложенными друг в друга
 - ▶ Размер кластеров
 - ▶ Кластеризация — основная задача или побочная
 - ▶ «Жёсткая» или «мягкая» кластеризация
- В задачах с разными особенностями могут быть уместны разные методы

ФОРМА КЛАСТЕРОВ



ФОРМА КЛАСТЕРОВ



ФОРМА КЛАСТЕРОВ



ФОРМА КЛАСТЕРОВ



ФОРМА КЛАСТЕРОВ



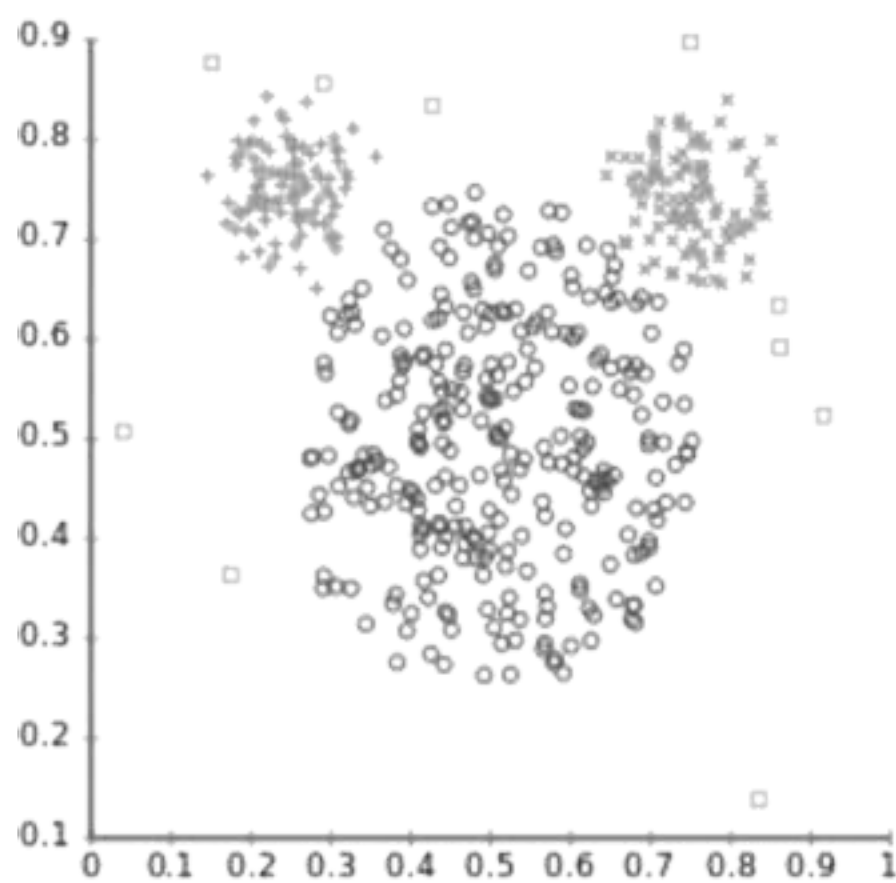
ФОРМА КЛАСТЕРОВ



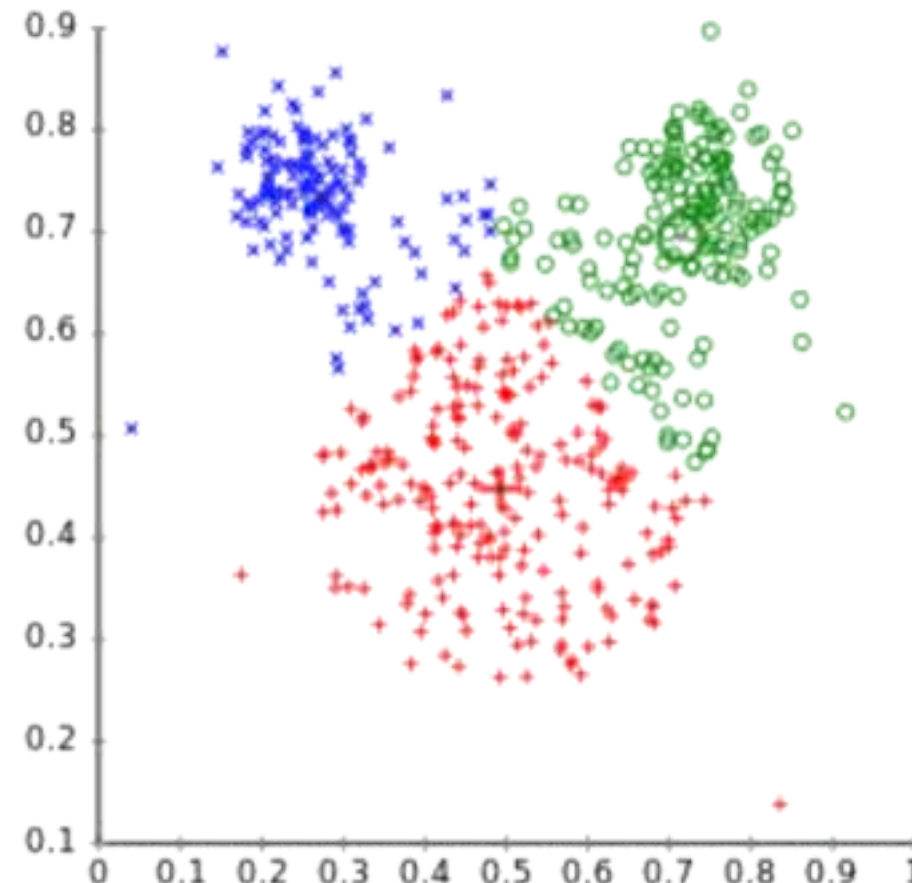
ФОРМА КЛАСТЕРОВ



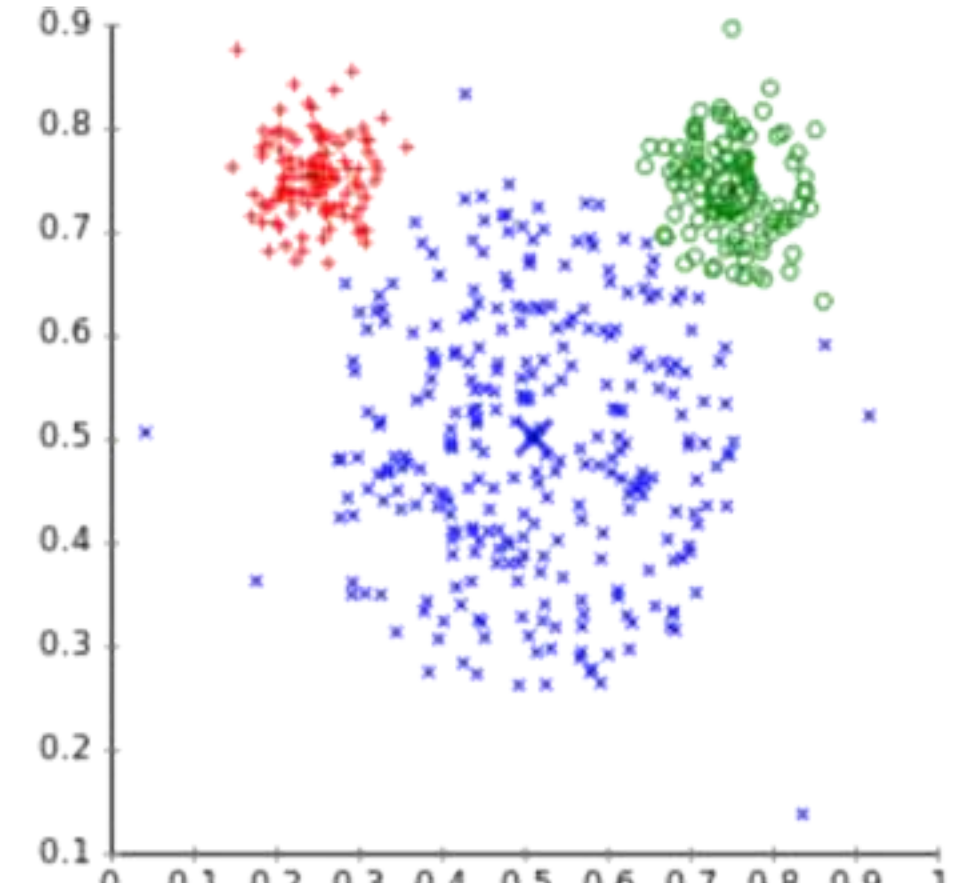
РАЗЛИЧИЯ В РЕЗУЛЬТАТАХ РАБОТЫ



Исходная выборка
("Mouse" dataset)



Метод k средних
(K-Means)



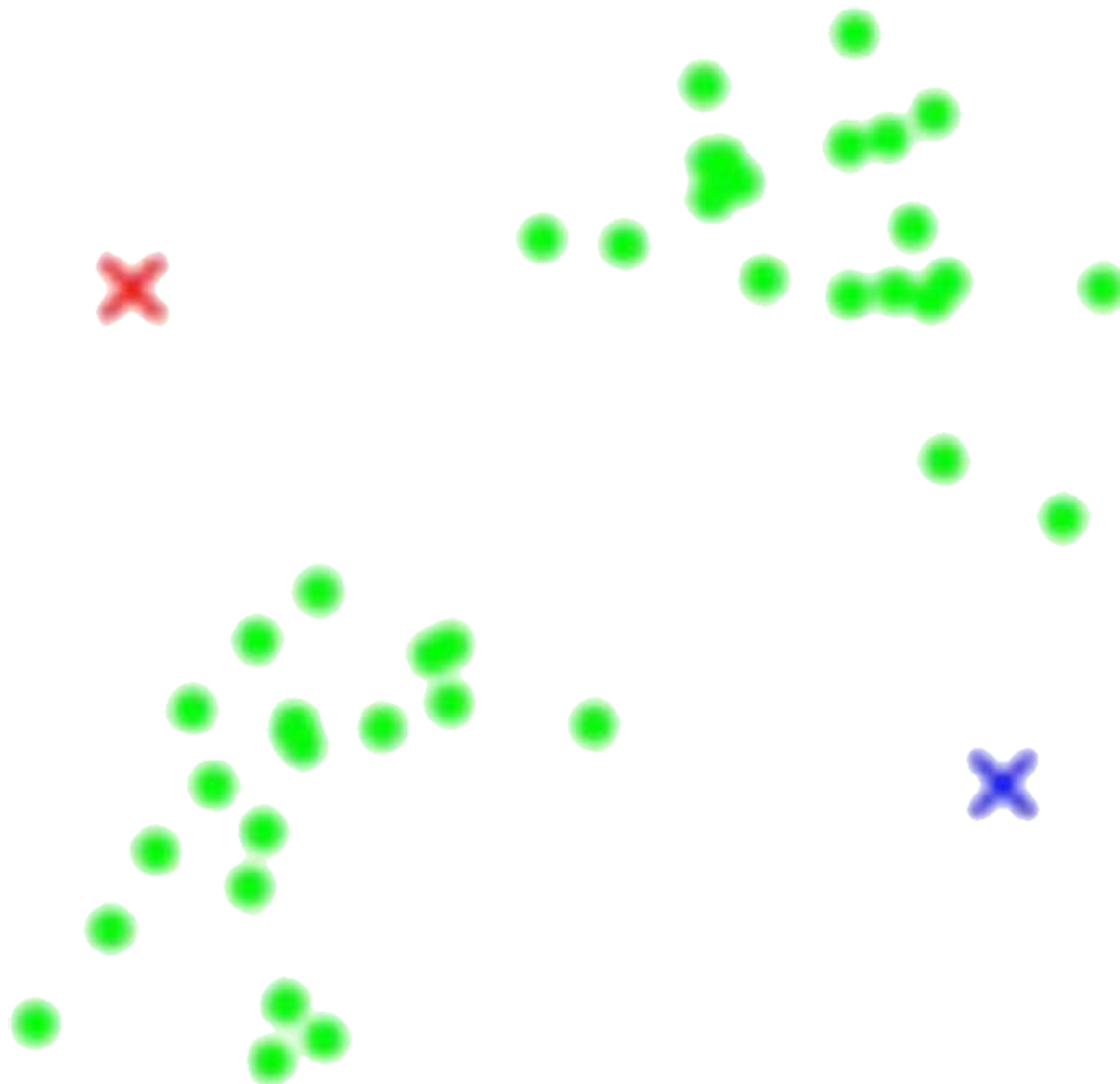
ЕМ-алгоритм

МЕТОД К СРЕДНИХ (K MEANS)

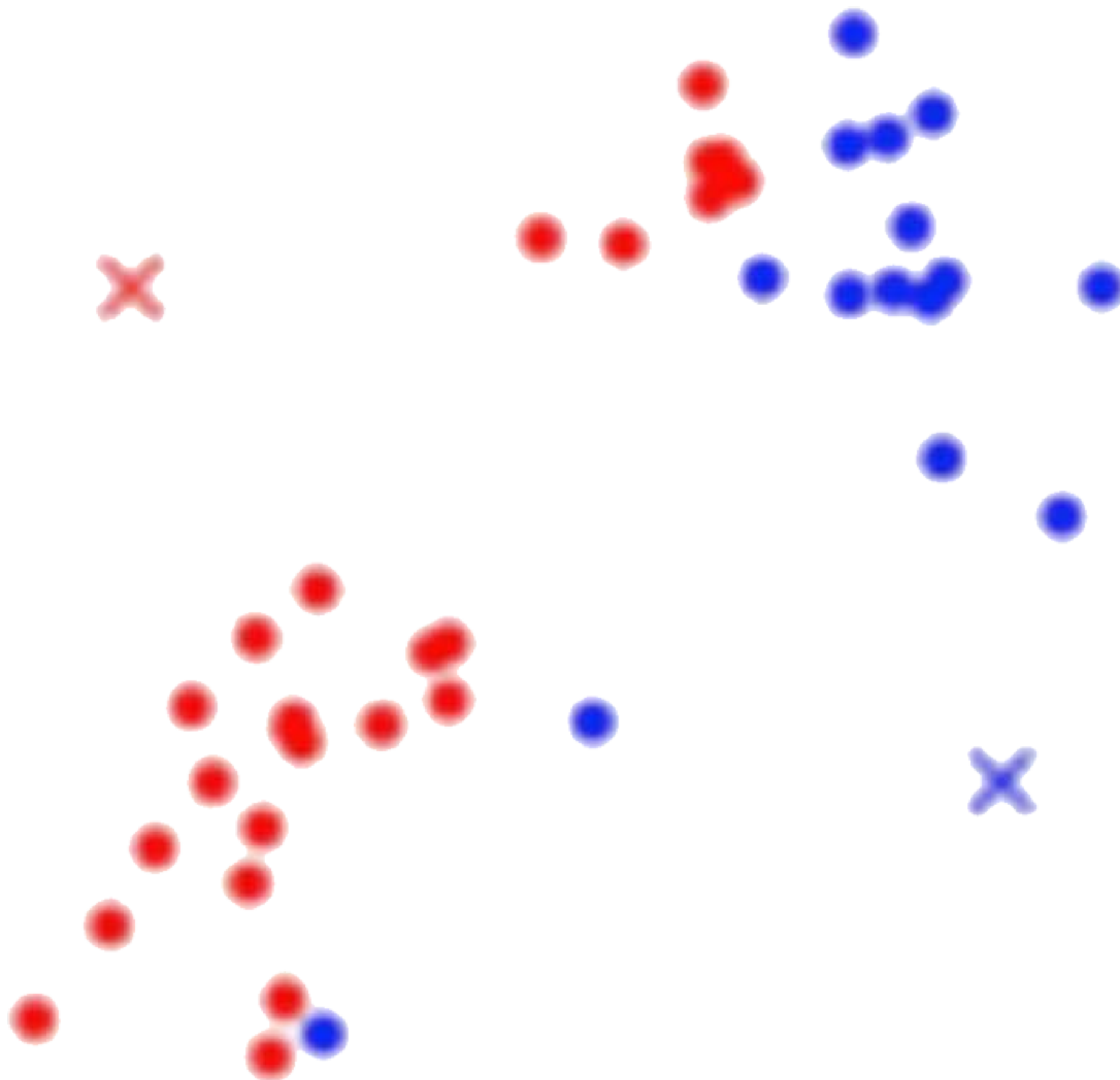
КАК РАБОТАЕТ K MEANS



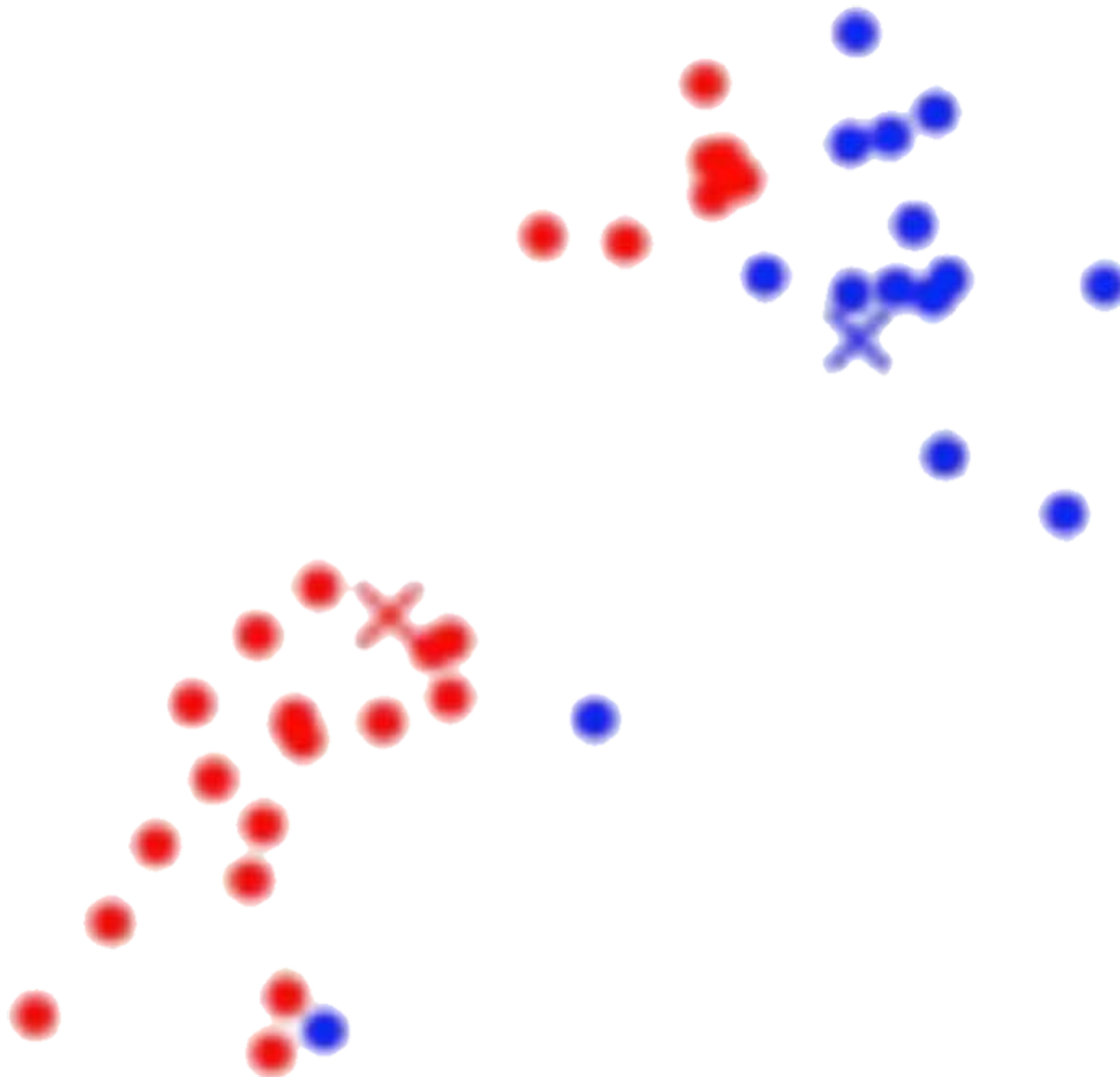
КАК РАБОТАЕТ K MEANS



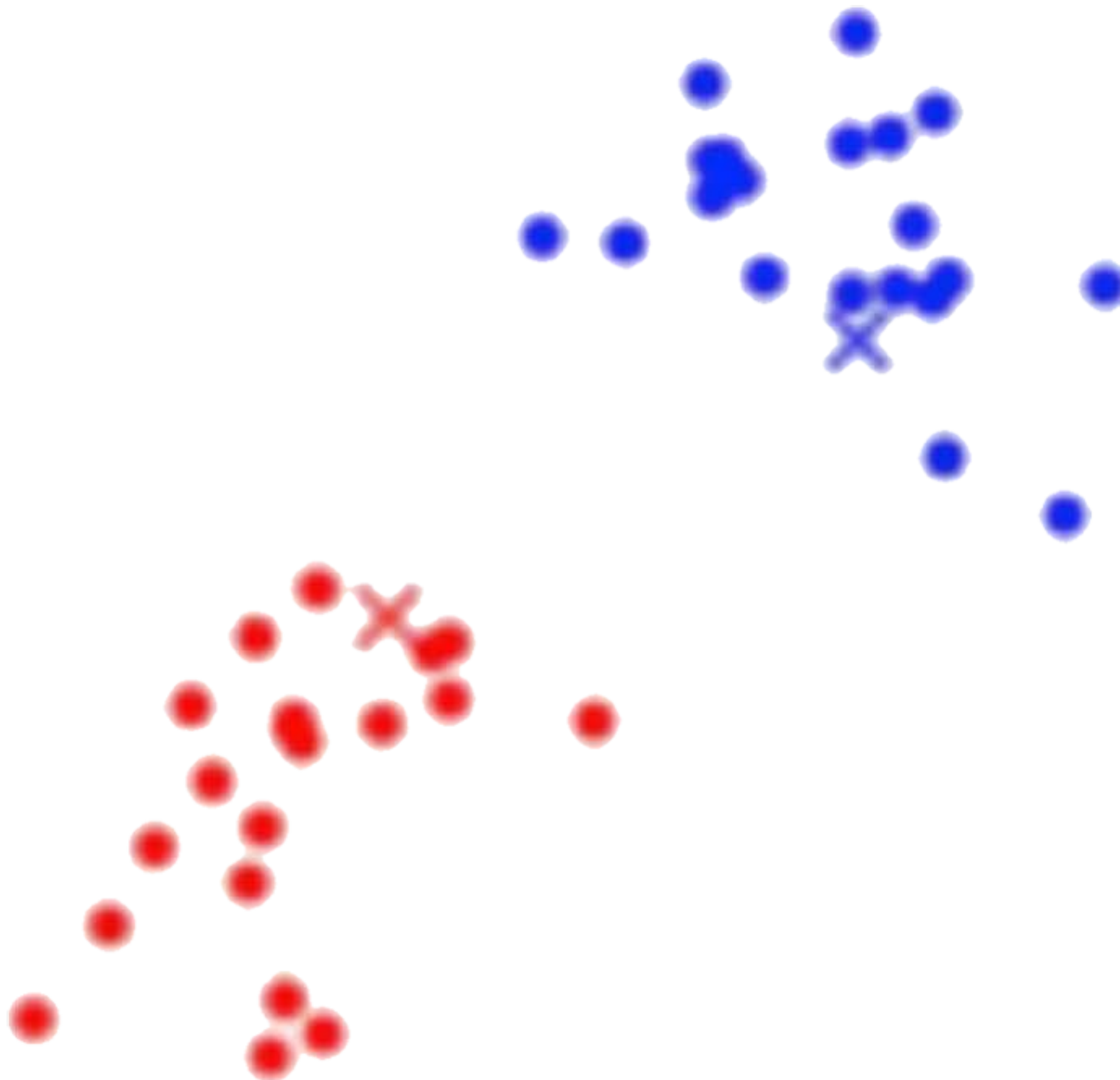
КАК РАБОТАЕТ K MEANS



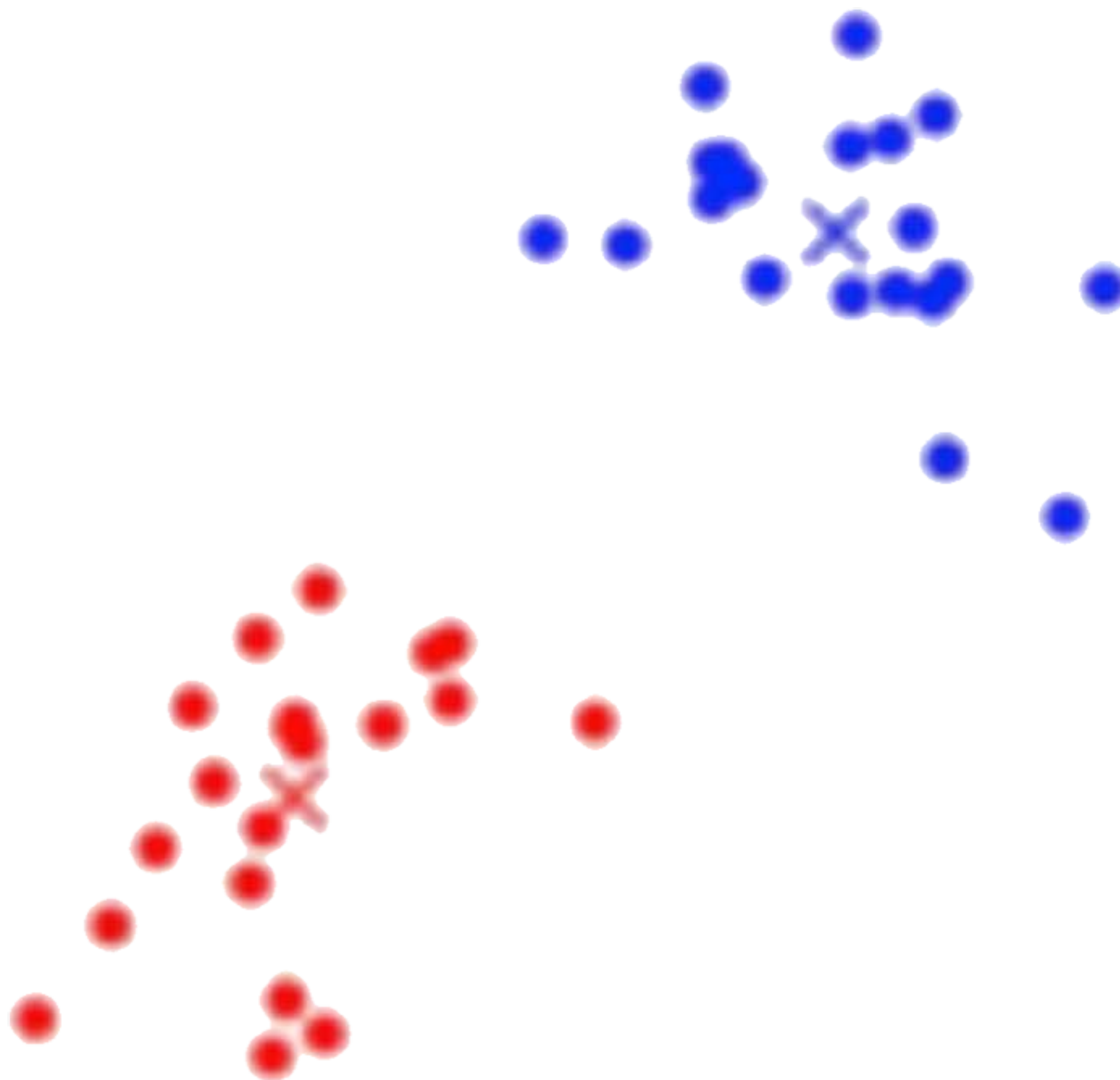
КАК РАБОТАЕТ K MEANS



КАК РАБОТАЕТ K MEANS



КАК РАБОТАЕТ K MEANS



ПРИМЕР: КВАНТИЗАЦИЯ ИЗОБРАЖЕНИЙ

Original image (96,615 colors)



ПРИМЕР: КВАНТИЗАЦИЯ ИЗОБРАЖЕНИЙ

Quantized image (64 colors, Random)

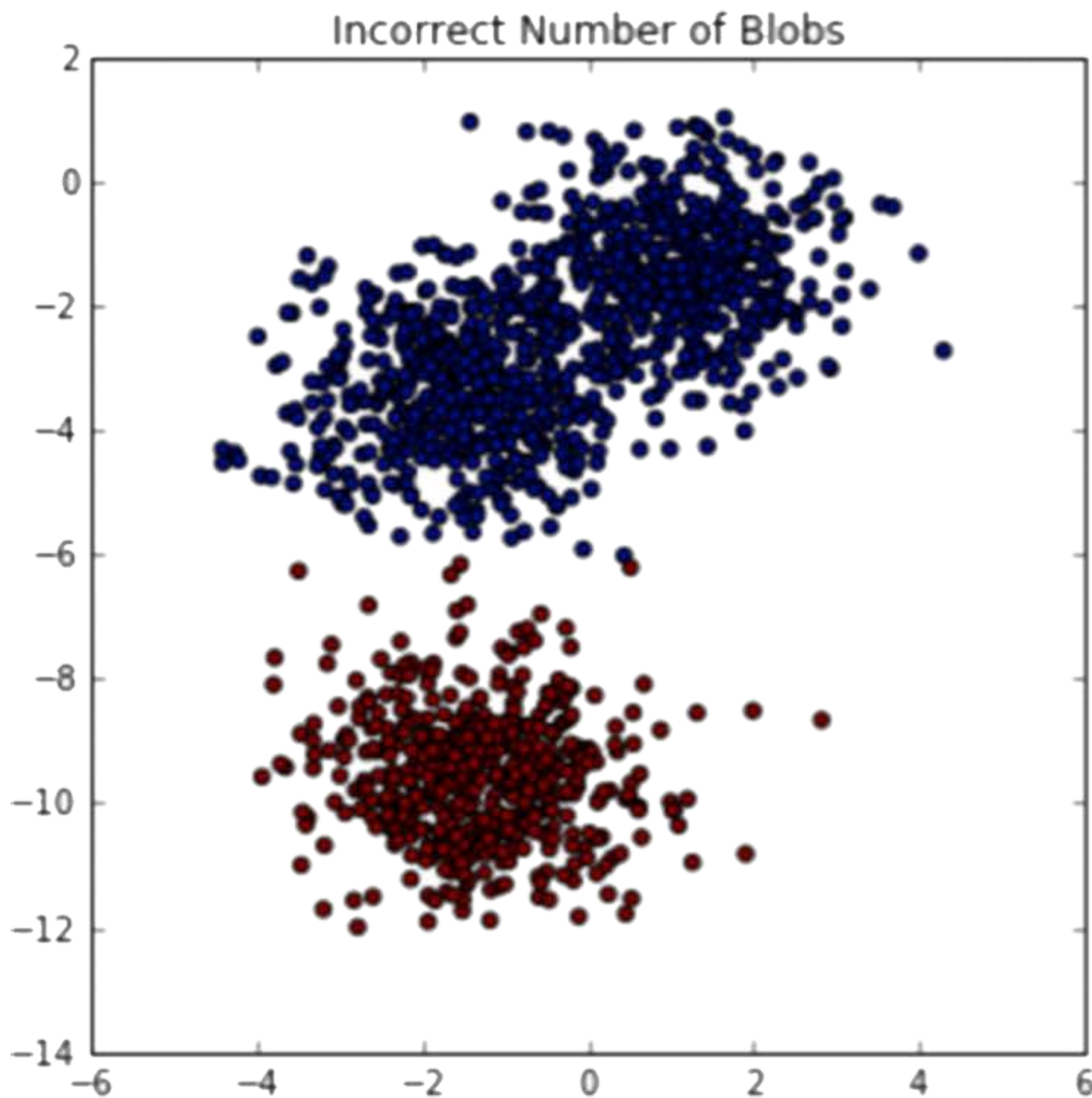


ПРИМЕР: КВАНТИЗАЦИЯ ИЗОБРАЖЕНИЙ

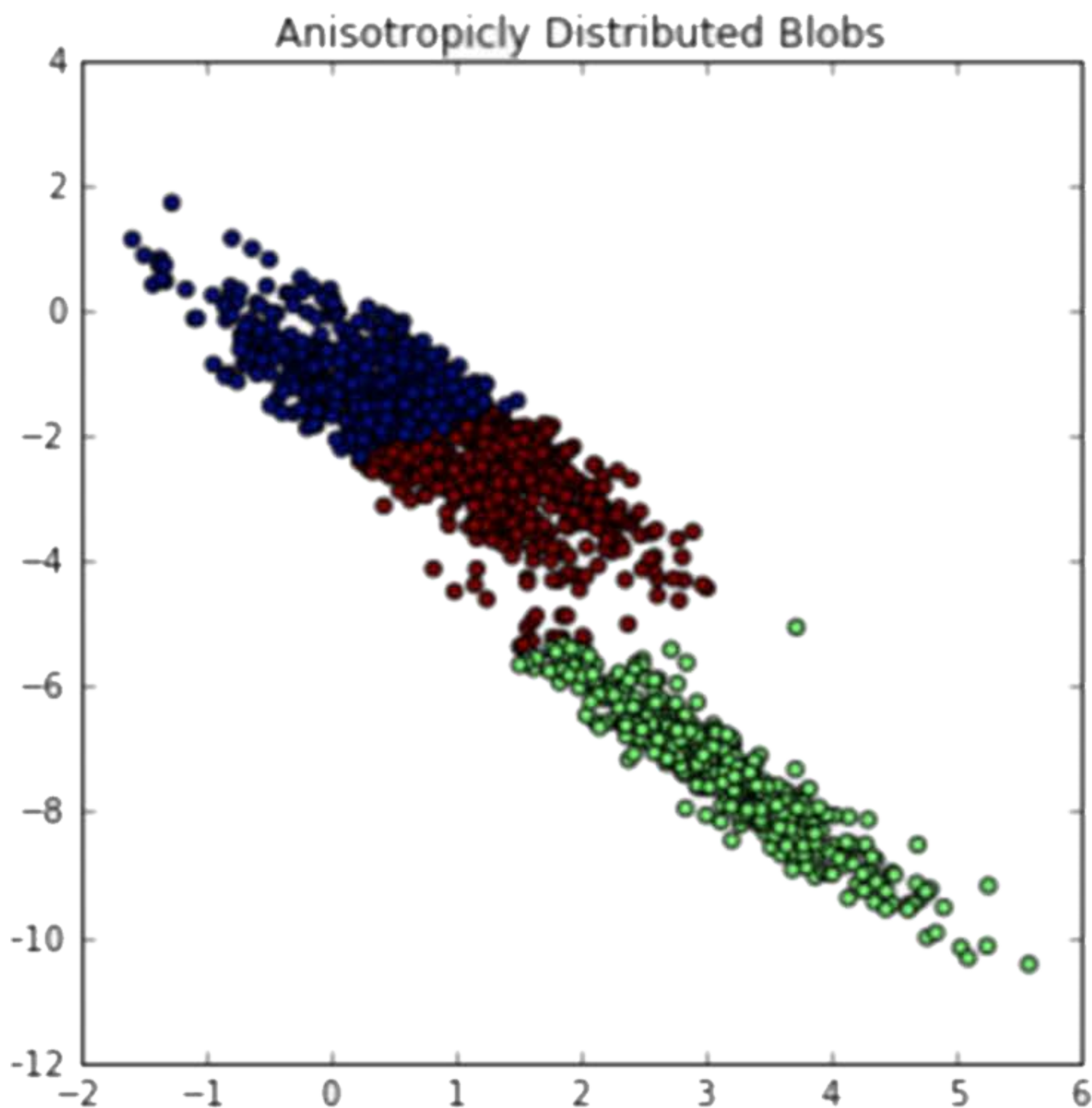
Quantized image (64 colors, K-Means)



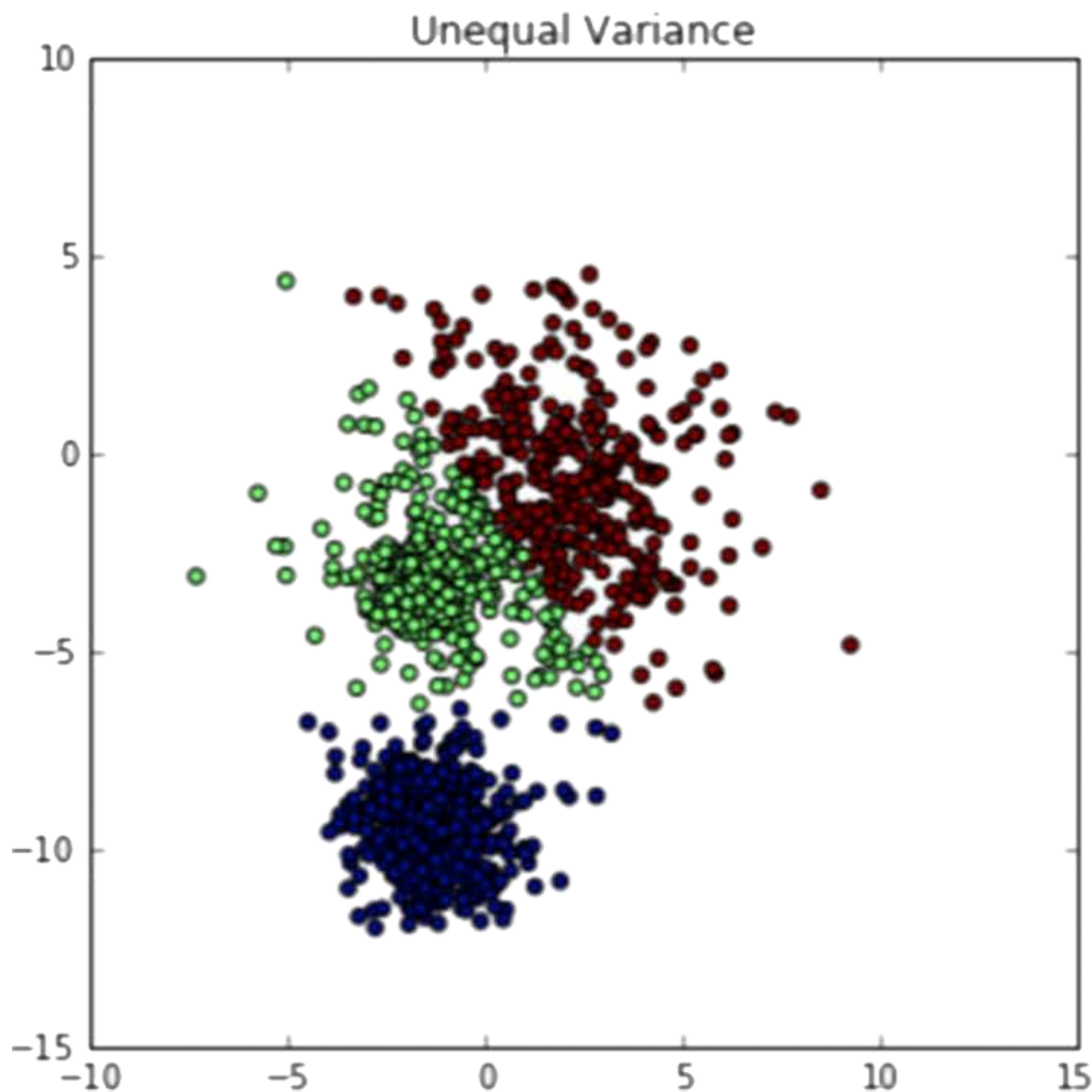
K MEANS И РАЗНЫЕ ФОРМЫ КЛАСТЕРОВ



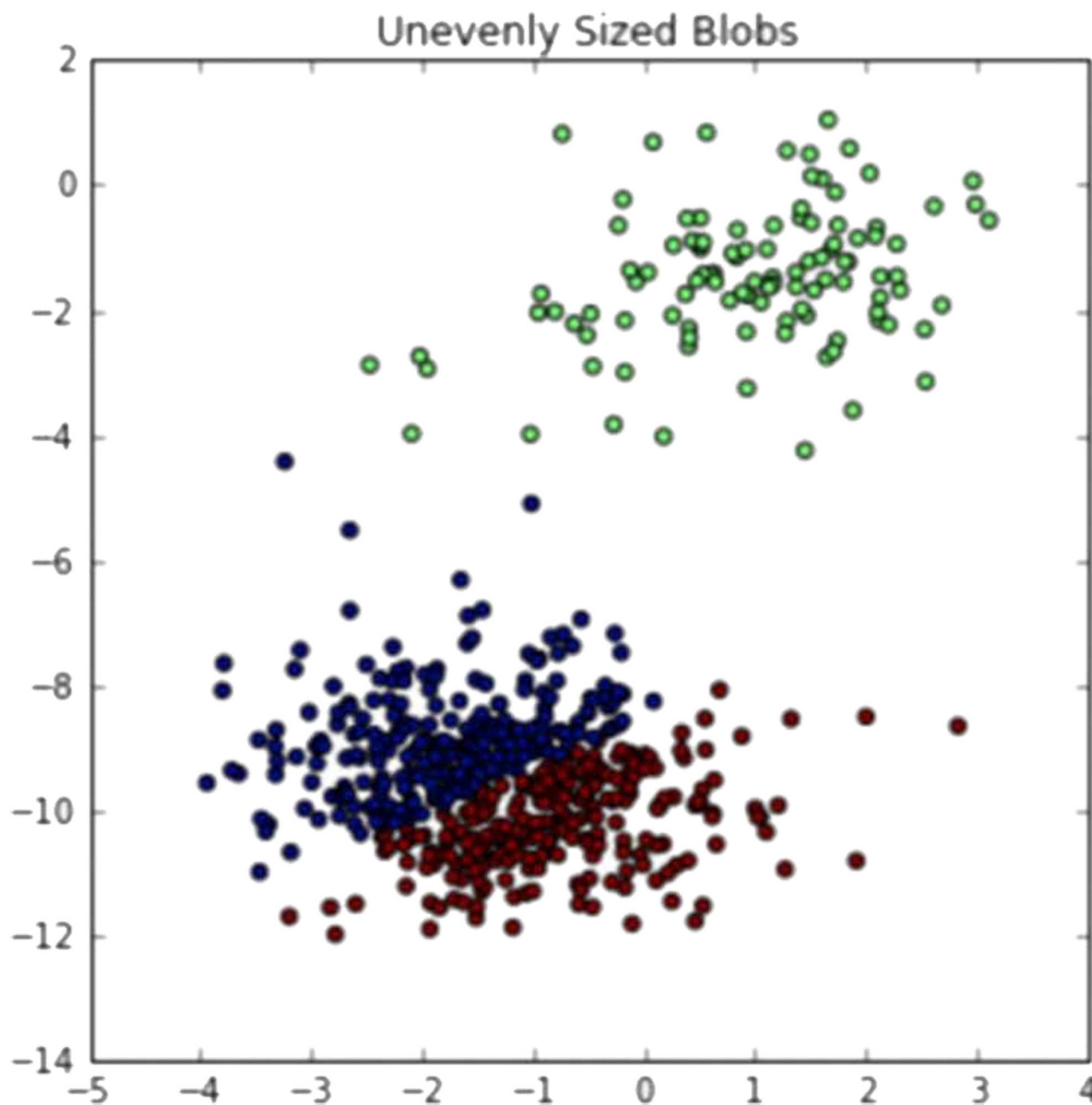
К MEANS И РАЗНЫЕ ФОРМЫ КЛАСТЕРОВ



K MEANS И РАЗНЫЕ ФОРМЫ КЛАСТЕРОВ



K MEANS И РАЗНЫЕ ФОРМЫ КЛАСТЕРОВ



ЧТО ОПТИМИЗИРУЕТ K MEANS

› Среднее внутрикластерное расстояние:

$$F_0 = \frac{\sum_{i < j} [y_i = y_j] \rho(x_i, x_j)}{\sum_{i < j} [y_i = y_j]} \rightarrow \min$$

ЧТО ОПТИМИЗИРУЕТ K MEANS

- › Среднее внутрикластерное расстояние:

$$F_0 = \frac{\sum_{i < j} [y_i = y_j] \rho(x_i, x_j)}{\sum_{i < j} [y_i = y_j]} \rightarrow \min$$

- › Альтернативный вариант, если есть центры кластеров:

$$\Phi_0 = \sum_{y \in Y} \frac{1}{|K_y|} \sum_{i: y_i = y} \rho^2(x_i, \mu_y) \rightarrow \min$$