

# Методы детектирования аномалий.

## Лекция 2: Кластеризация

Иван Шанин  
ivan.shanin@gmail.com

ИПИ РАН

18.02.2019

# Локальный интегральный показатель корреляции (LOCI)<sup>1</sup>

Для объекта  $\bar{X}$  определим

- ▶  $M(\bar{X}, \epsilon)$  - плотность данных в  $\epsilon$ -окрестности объекта  $\bar{X}$
- ▶  $AM(\bar{X}, \epsilon, \delta)$  - среднее значение  $M(\bar{X}, \epsilon)$  по всем объектам  $\delta$ -окрестности объекта  $\bar{X}$

$$AM(\bar{X}, \epsilon, \delta) = MEAN_{\{\bar{Y}: \rho(\bar{X}, \bar{Y}) \leq \delta\}} M(\bar{Y}, \epsilon)$$

- ▶ Мультигранулярный показатель отклонения

$$MDEF(\bar{X}, \epsilon, \delta) = 1 - \frac{M(\bar{X}, \epsilon)}{AM(\bar{X}, \epsilon, \delta)}$$

Большое значение  $MDEF$  является признаком аномальности  $\bar{X}$ .

---

<sup>1</sup><http://www.cs.cmu.edu/~christos/PUBLICATIONS/icde03-loci-tr.pdf>

# Диаграмма LOCI

Пусть  $\epsilon = \frac{1}{2}\delta$ .

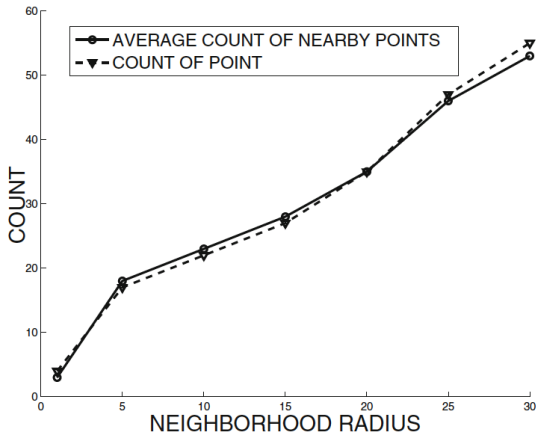


Рис. 1: Объект внутри кластера

# Диаграмма LOCI

Пусть  $\epsilon = \frac{1}{2}\delta$ .

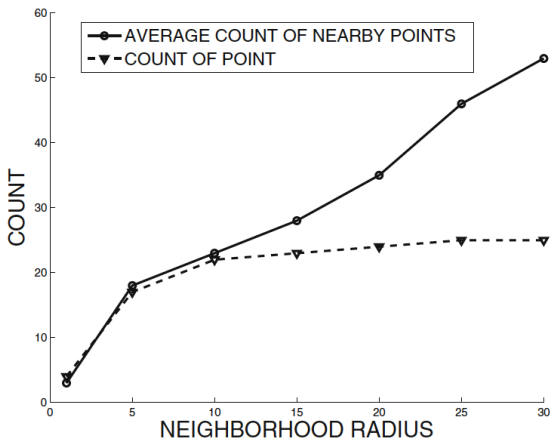


Рис. 2: Небольшой кластер аномальных объектов

# Диаграмма LOCI

Пусть  $\epsilon = \frac{1}{2}\delta$ .

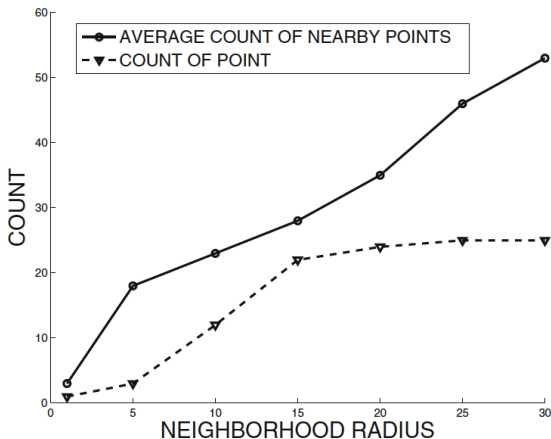
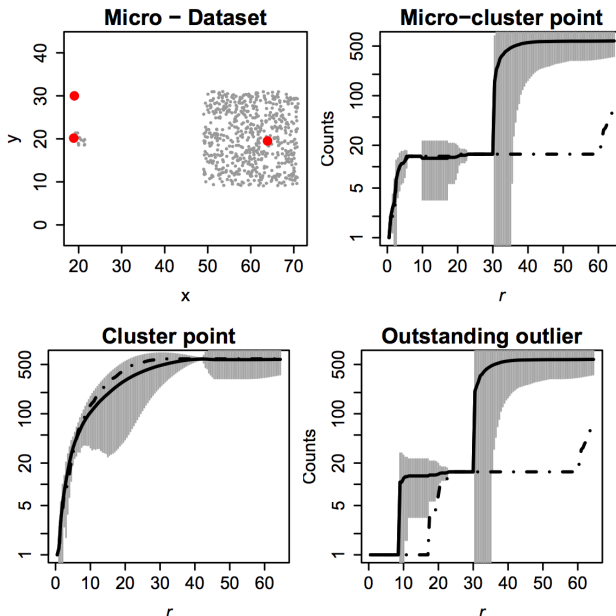


Рис. 3: Аномальный объект

# Диаграмма LOCI



## Смесь вероятностных распределений

Определим порождающую модель данных  $M$ . Предположим, что данные сгенерированы из вероятностных распределений  $G_1, G_2, \dots, G_k$  по следующему стохастическому процессу:

- ▶ с вероятностью  $\alpha_r$  выбирается вероятностное распределение  $G_r$
- ▶ из  $G_r$  генерируется объект выборки

где  $\alpha_1, \dots, \alpha_r$  - набор **априорных** вероятностей.

Функция плотности распределения объектов:

$$p(\bar{X}|M) = \sum_{i=1}^k \alpha_i p_i(\bar{X}|\theta_i)$$

Функция правдоподобия выборки  $D$ :

$$Likelihood(D|M) = \prod_{\bar{X} \in D} p(\bar{X}|M)$$

## Задача

По выборке  $D$  и заранее заданному параметру  $k$  оценить параметры модели  $M(\bar{\alpha}, \bar{\theta})$ .

Требуется максимизировать логарифм правдоподобия

$$L(D|M) = \log \left[ \prod_{\bar{X} \in D} p(\bar{X}|M) \right] = \sum_{\bar{X} \in D} \log \left[ \sum_{i=1}^k \alpha_i p_i(\bar{X}|\theta_i) \right]$$

Данная задача решается чередованием двух этапов:

- ▶ **Е-шаг (expectation):** По текущей оценке параметров  $M$  оценить  $P(I_r|\bar{X})$  – апостериорную вероятность того, что объект  $\bar{X}$  был сгенерирован  $r$ -й компонентной (по формуле Байеса)
- ▶ **М-шаг (maximization):** По полученной разметке выборки построить новые значения параметров  $M(\bar{\alpha}, \bar{\theta})$



## ЕМ-алгоритм

Приведем формулы<sup>2</sup> расчета скрытых и явных параметров на каждом шаге.

► **Е-шаг (формула Байеса):**

$$\hat{P}(I_r|\bar{X}, \bar{\alpha}, \bar{\theta}) = \frac{\alpha_r p_r(\bar{X}|\theta_r)}{\sum_{i=1}^k \alpha_i p_r(\bar{X}|\theta_i)}$$

► **М-шаг (максимизация правдоподобия):**

$$\alpha_r = \frac{1}{|D|} \sum_{\bar{X} \in D} \hat{P}(I_r|\bar{X}, \bar{\alpha}, \bar{\theta});$$

$$\theta_r = \arg \max_{\theta} \sum_{\bar{X} \in D} \hat{P}(I_r|\bar{X}, \bar{\alpha}, \bar{\theta}) \log p_r(\bar{X}|\theta)$$

---

<sup>2</sup><http://www.machinelearning.ru/wiki/images/6/6d/Voron-ML-1.pdf>

# Метод $k$ средних

Выберем  $k$  случайных точек - нулевое приближение к центрам масс кластеров. Затем будем повторять следующие шаги:

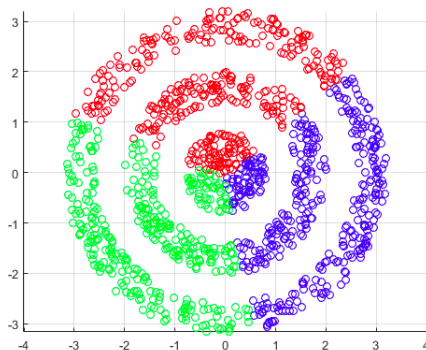
- ▶ Отнести каждый объект выборки к одному из  $k$  классов (ближайшему по некоторой функции расстояния)
- ▶ По получившимся кластерам рассчитать новые центры масс

Таким образом на каждом шаге будет монотонно убывать суммарное квадратичное отклонение точек кластеров от центров этих кластеров. Критерий останова - шаг, на котором данное значение не изменилось.

- ▶ Метод является упрощенной версией EM-алгоритма

## Проблемы метода k-средних

- ▶ Не гарантируется достижение глобального минимума суммарного квадратичного отклонения, а только одного из локальных минимумов.
- ▶ Результат зависит от выбора исходных центров кластеров, их оптимальный выбор неизвестен.
- ▶ Число кластеров надо знать заранее.



# DBSCAN

Поделим выборку на следующие классы объектов:

- ▶ Элемент ядра, если хотя бы  $n$  точек находятся в  $\epsilon$ -окрестности
- ▶ Объект  $p$  доступен из элемента ядра  $q$  **напрямую**, если находится в его  $\epsilon$ -окрестности
- ▶ Объект  $p$  доступен из элемента ядра  $q$ , если существует путь  $p_1, \dots, p_n$  из  $p$  в  $q$ , где все  $(p_i, p_{i+1})$  связаны напрямую и являются элементами ядра (возможно, за исключением  $q$ )

