



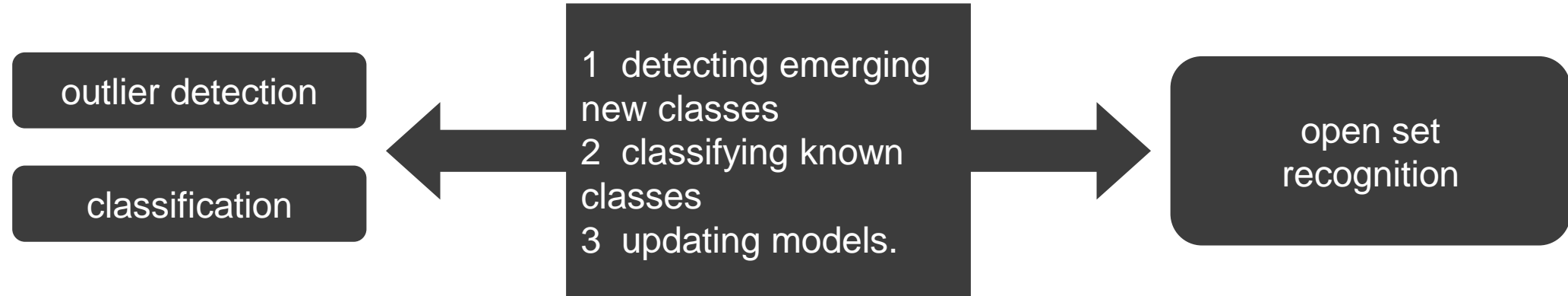
outlier

2018.4.26



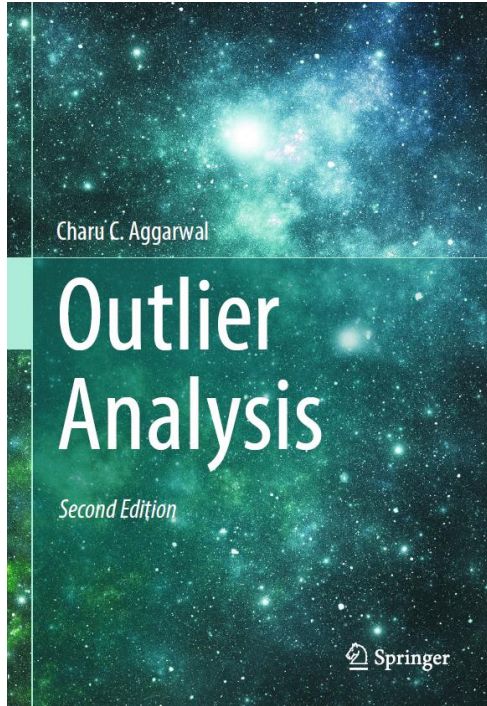
## Question

Difference between outlier detection and open set recognition?



将已有异常检测方法“\*\*”为既能检测离群值又同时可完成多分类任务

**Example:** Classification Under Streaming Emerging New Classes: A Solution Using Completely-Random Trees  
Xin Mu, Kai Ming Ting, and Zhi-Hua Zhou TKDE2017



Semi-Supervised

Transductive



outlier

- 1.3 The Basic Outlier Detection Models . . . . .
  - 1.3.1 Feature Selection in Outlier Detection . . . . .
  - 1.3.2 **Extreme-Value Analysis** . . . . .
  - 1.3.3 Probabilistic and Statistical Models . . . . .
  - 1.3.4 Linear Models . . . . .
    - 1.3.4.1 Spectral Models . . . . .
  - 1.3.5 Proximity-Based Models . . . . .
  - 1.3.6 Information-Theoretic Models . . . . .
  - 1.3.7 High-Dimensional Outlier Detection . . . . .
- 1.4 Outlier Ensembles . . . . .
  - 1.4.1 Sequential Ensembles . . . . .
  - 1.4.2 Independent Ensembles . . . . .
- 1.5 The Basic Data Types for Analysis . . . . .
  - 1.5.1 Categorical, Text, and Mixed Attributes . . . . .
  - 1.5.2 When the Data Values have Dependencies . . . . .
    - 1.5.2.1 Times-Series Data and Data Streams . . . . .
    - 1.5.2.2 Discrete Sequences . . . . .
    - 1.5.2.3 Spatial Data . . . . .
    - 1.5.2.4 Network and Graph Data . . . . .
- 1.6 Supervised Outlier Detection . . . . .
- 1.7 Outlier Evaluation Techniques . . . . .

# Extreme-Value Analysis

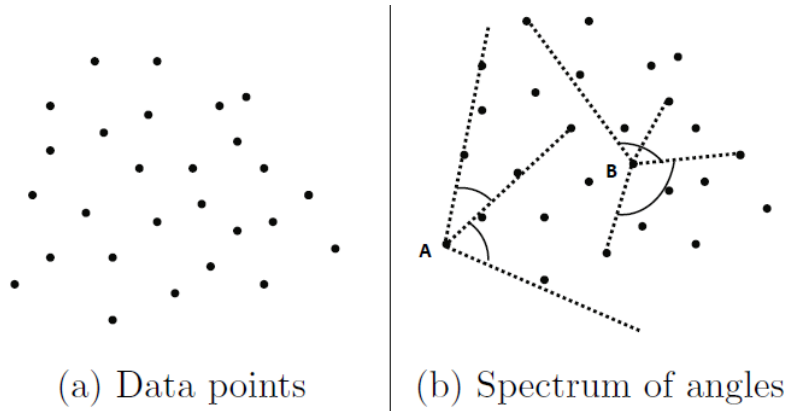
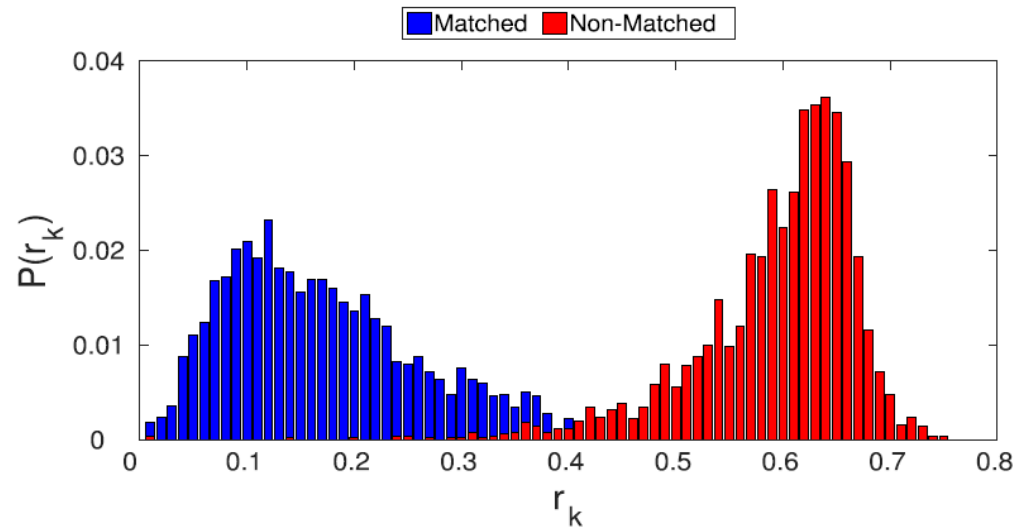
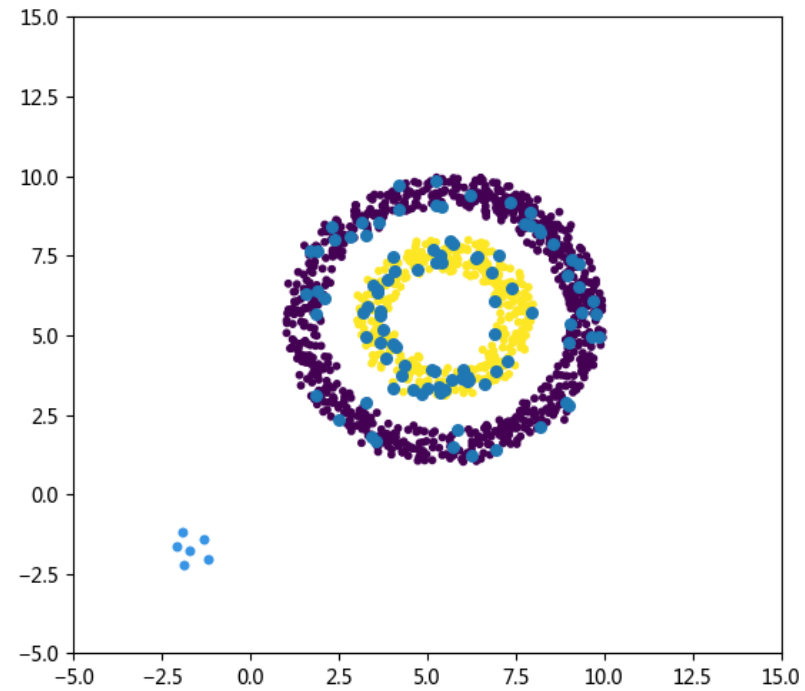
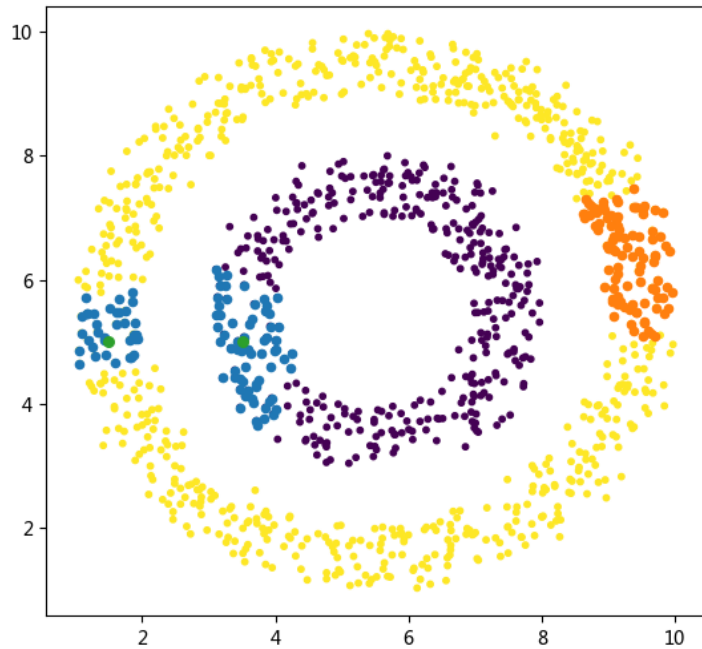


Figure 2.6: Angle-based outlier detection



(PAMI 2017)

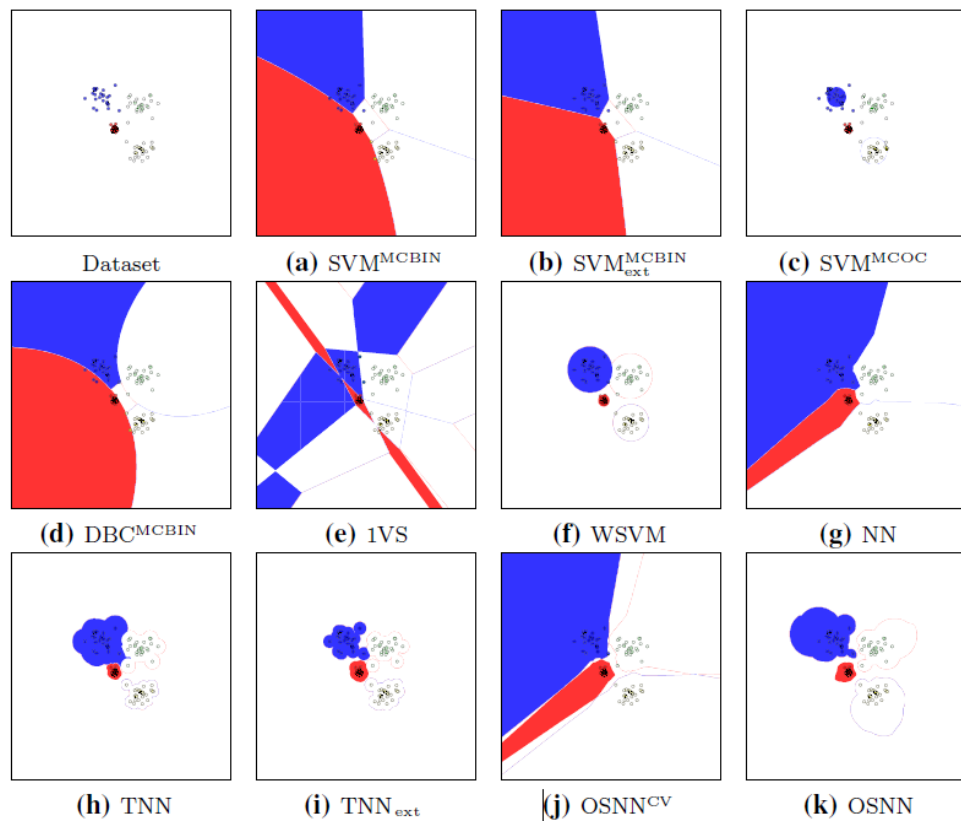


# Two worries

## 1. How to balance specific and generalization?

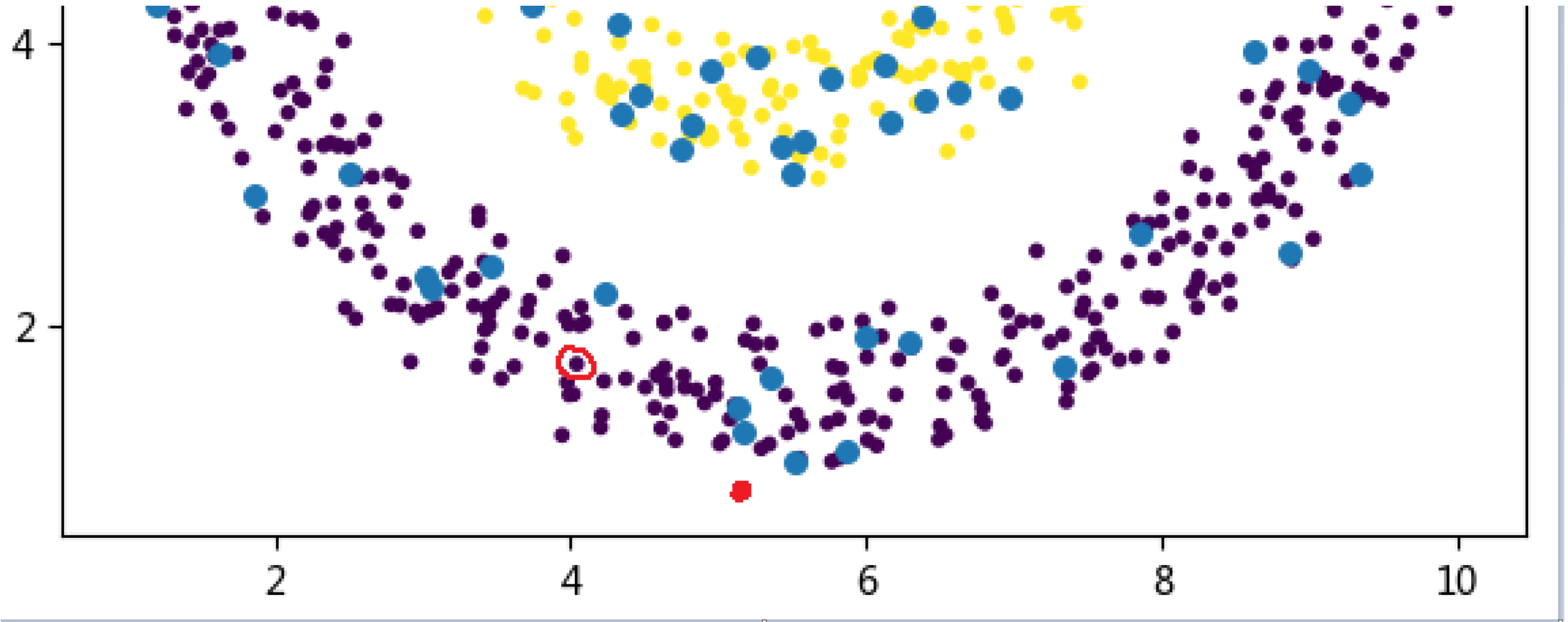


(Machine Learning 2017)

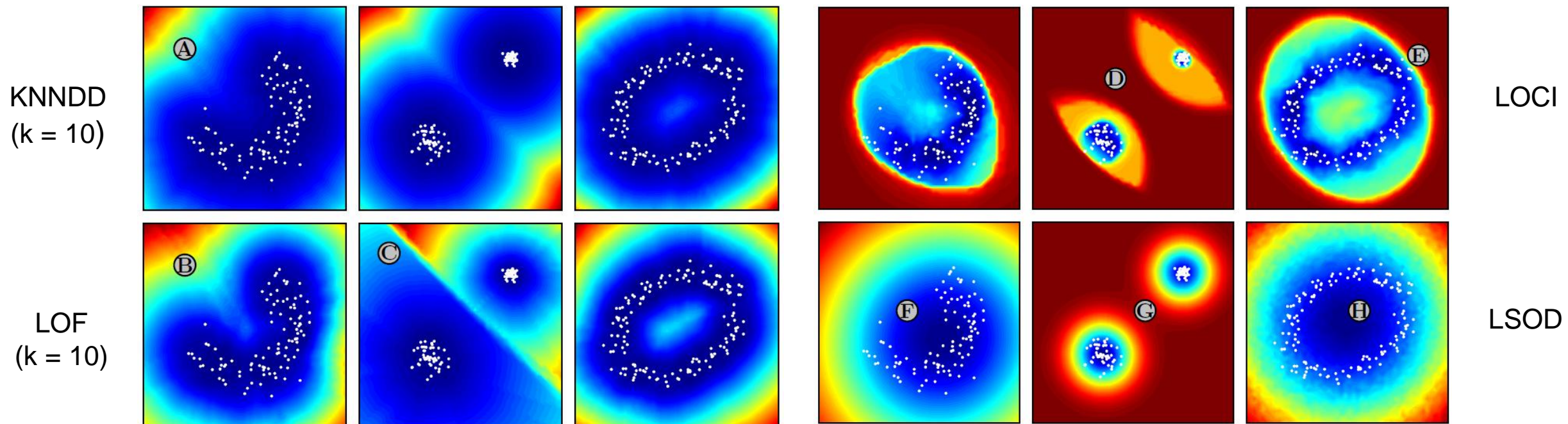
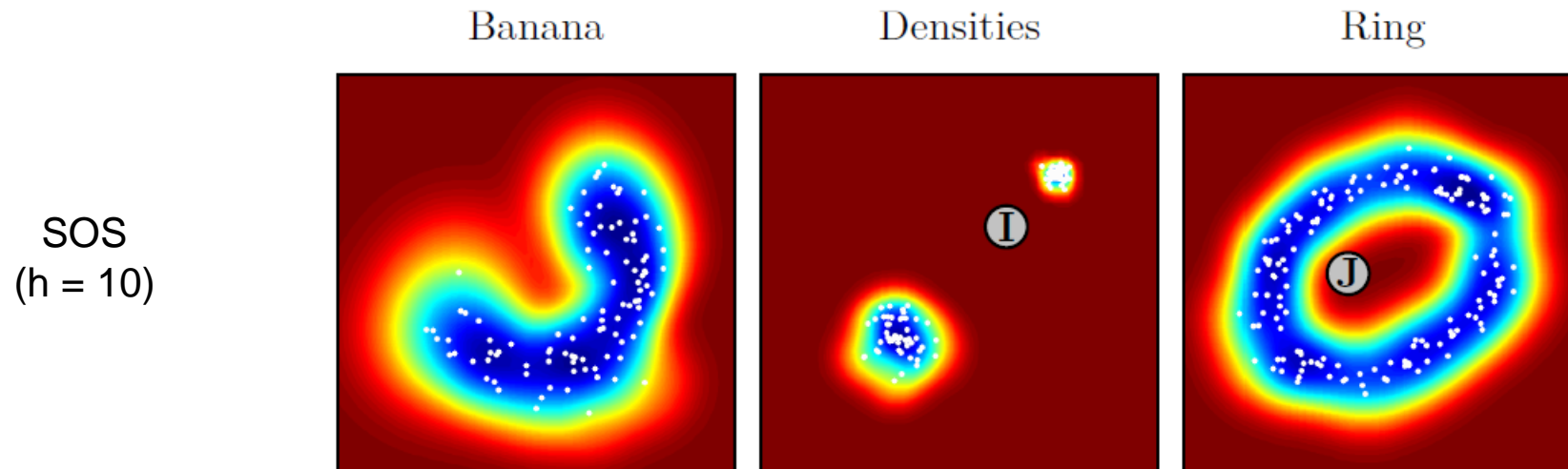


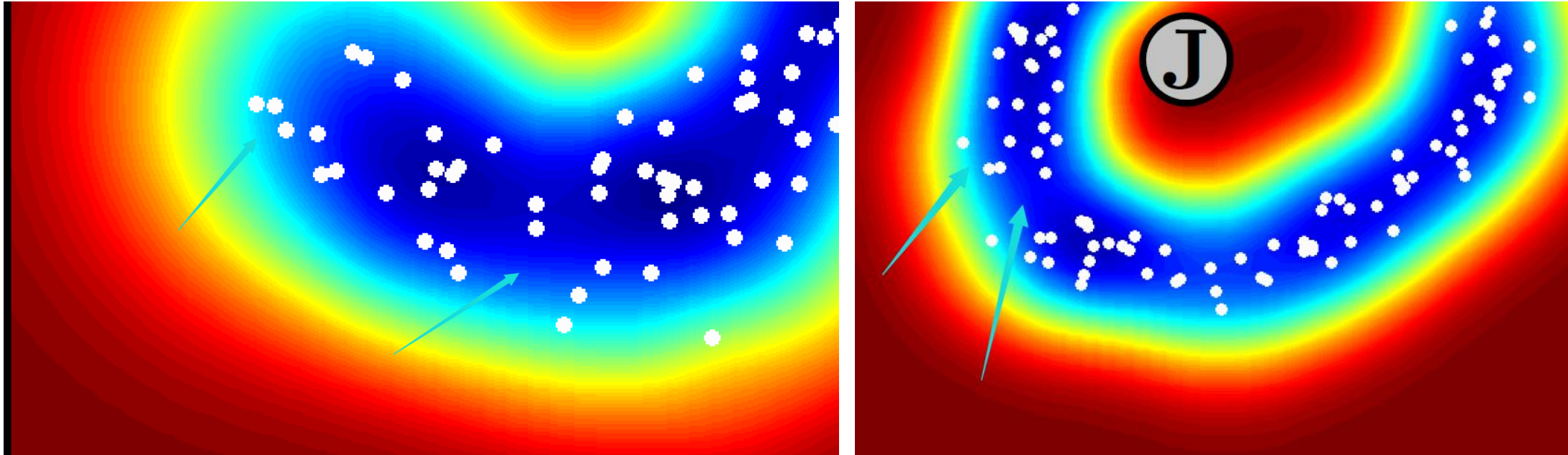


2. Smoothness assumption used in semi-supervised classification tasks, is often both local and global.



# Stochastic Outlier Selection





Visualizing Data using t-SNE JMLR 2011

The similarity of datapoint  $x_j$  to datapoint  $x_i$  is the conditional probability, that  $x_i$  **would pick  $x_j$  as its neighbor** if neighbors were picked in proportion to their probability density under a Gaussian centered at  $x_i$

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

$$p_{i|i} = 0$$

$$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$



$$C = \sum_i KL(P_i||Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}} \quad P_g(x) \rightarrow 0 \quad P_r \rightarrow 1 \quad P_g(x) \log \frac{P_g(x)}{P_r(x)} \rightarrow 0$$

$$\frac{\delta C}{\delta y_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j) \quad P_g(x) \rightarrow 1 \quad P_r \rightarrow 0 \quad P_g(x) \log \frac{P_g(x)}{P_r(x)} \rightarrow +\infty$$

the SNE cost function focuses on retaining the local structure of the data in the map

appropriate  $\sigma_i$

Any particular value of  $\sigma_i$  induces a probability distribution,  $P_i$ , over all of the other datapoints. This distribution has an entropy which increases as  $\sigma_i$  increases

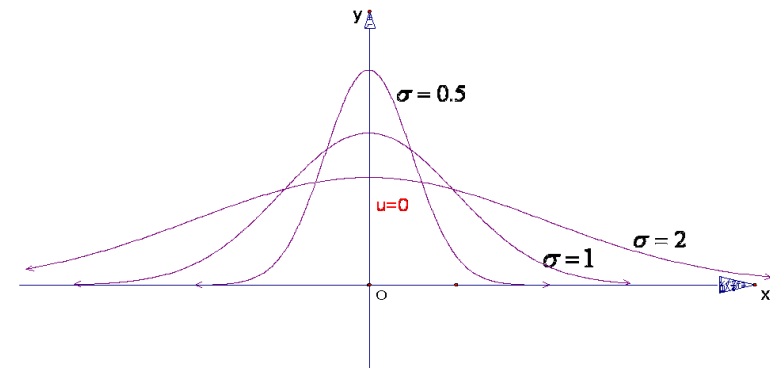
perplexity

$$Perp(P_i) = 2^{H(P_i)} \quad H(P_i) = - \sum_j p_{j|i} \log_2 p_{j|i}$$

The perplexity can be interpreted as a smooth measure of the effective number of neighbors.

In dense regions, a smaller value of  $\sigma_i$  is usually more appropriate than in sparser regions.

short version



$$a_{ij} = \begin{cases} \exp(-d_{ij}^2 / 2\sigma_i^2) & \text{if } i \neq j \\ 0 & \text{if } i = j \end{cases}$$

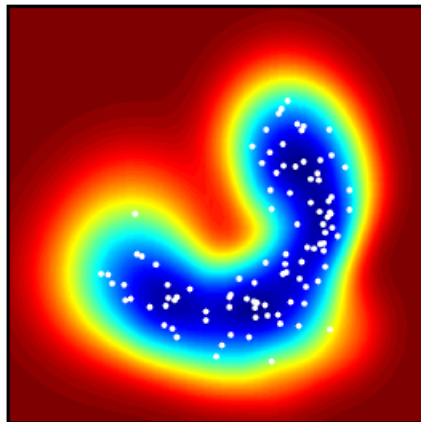
binding matrix **B**

$$b_{ij} = \frac{a_{ij}}{\sum_{k=1}^n a_{ik}}$$

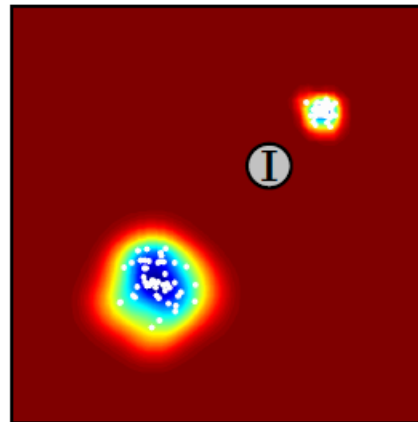
$$p(\mathbf{x}_i \in C_O) = \prod_{j \neq i} (1 - b_{ji})$$

$$f_{\text{SOS}}(\mathbf{x}) = \begin{cases} \text{outlier} & \text{if } \varphi_{\text{SOS}}(\mathbf{x}) > \theta, \\ \text{inlier} & \text{if } \varphi_{\text{SOS}}(\mathbf{x}) \leq \theta. \end{cases}$$

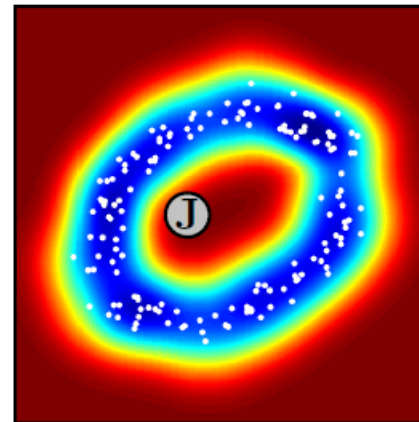
Banana



Densities



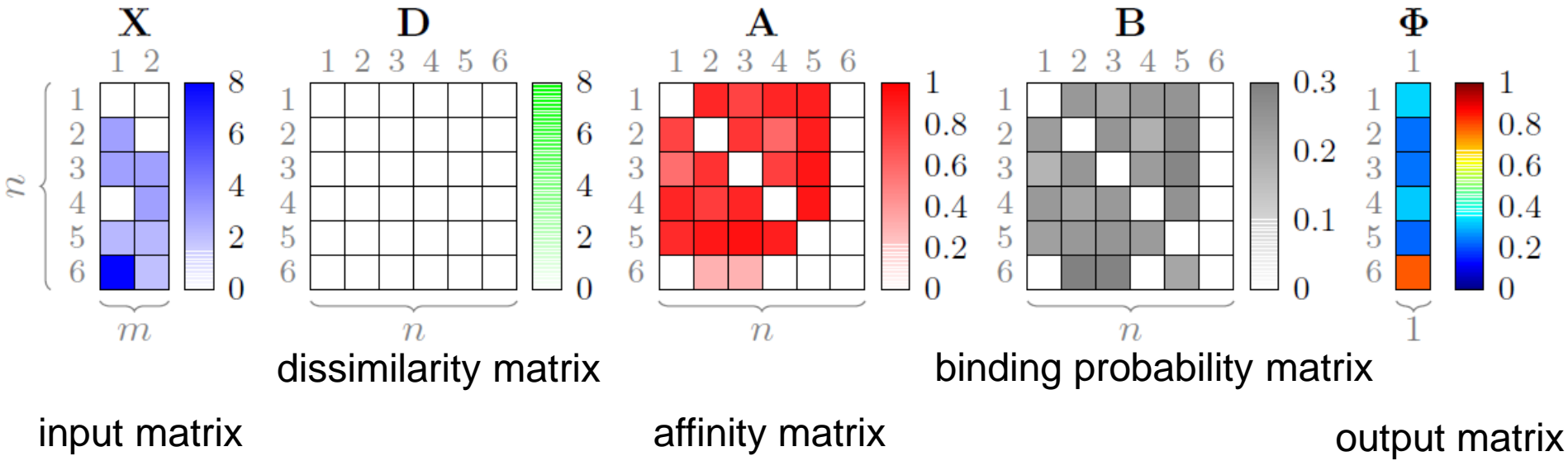
Ring



n. 密切关系; 吸引力; 姻亲关系; 类同

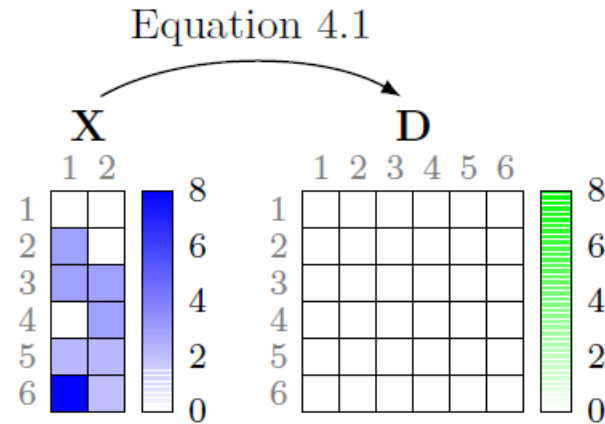
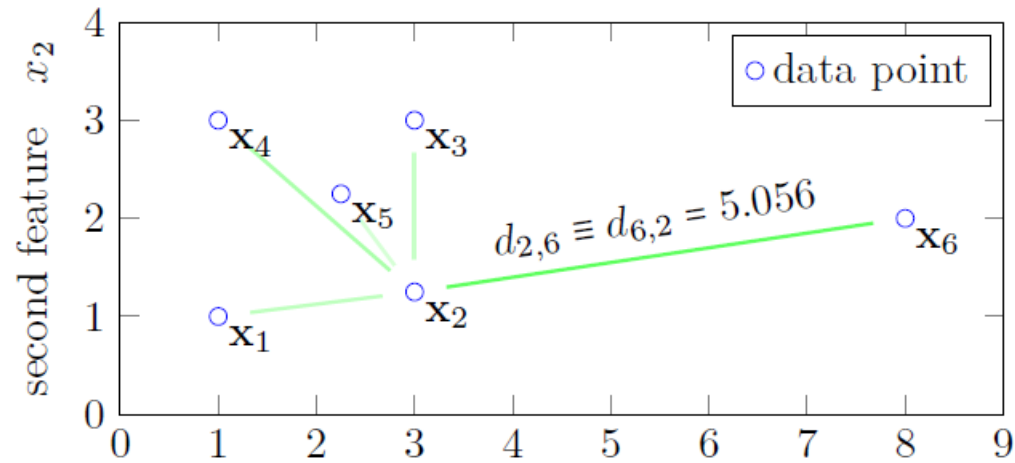
SOS employs the concept of **affinity** to quantify the relationship from one data point to another data point. Affinity is **proportional to the similarity between two data points**.

**Affinity:** the problem of clustering and dimension reduction





The features of the data points are used to measure the dissimilarity between pairs of data points.

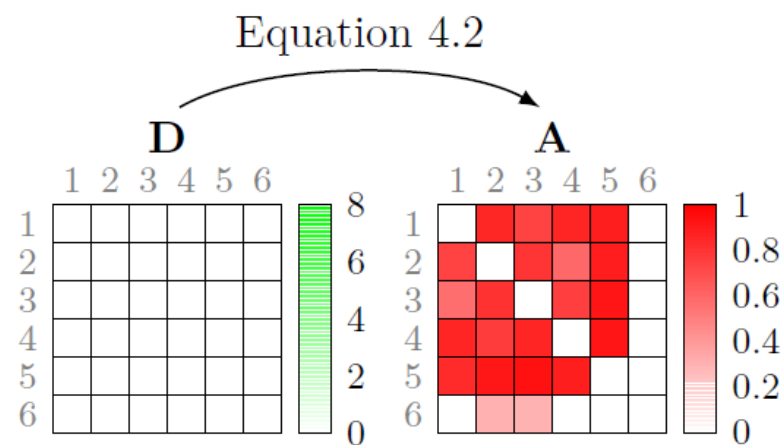
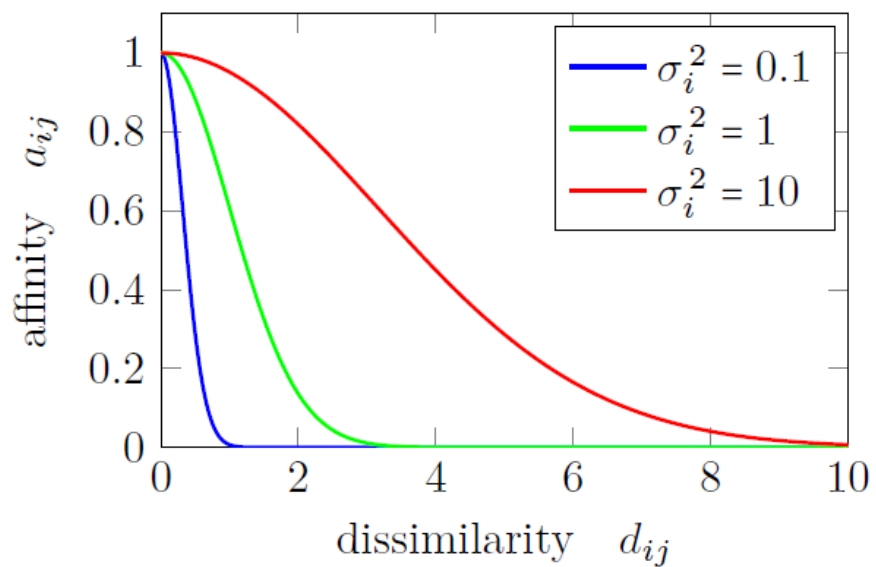


$$d_{ij} = \sqrt{\sum_{k=1}^m (x_{jk} - x_{ik})^2}$$

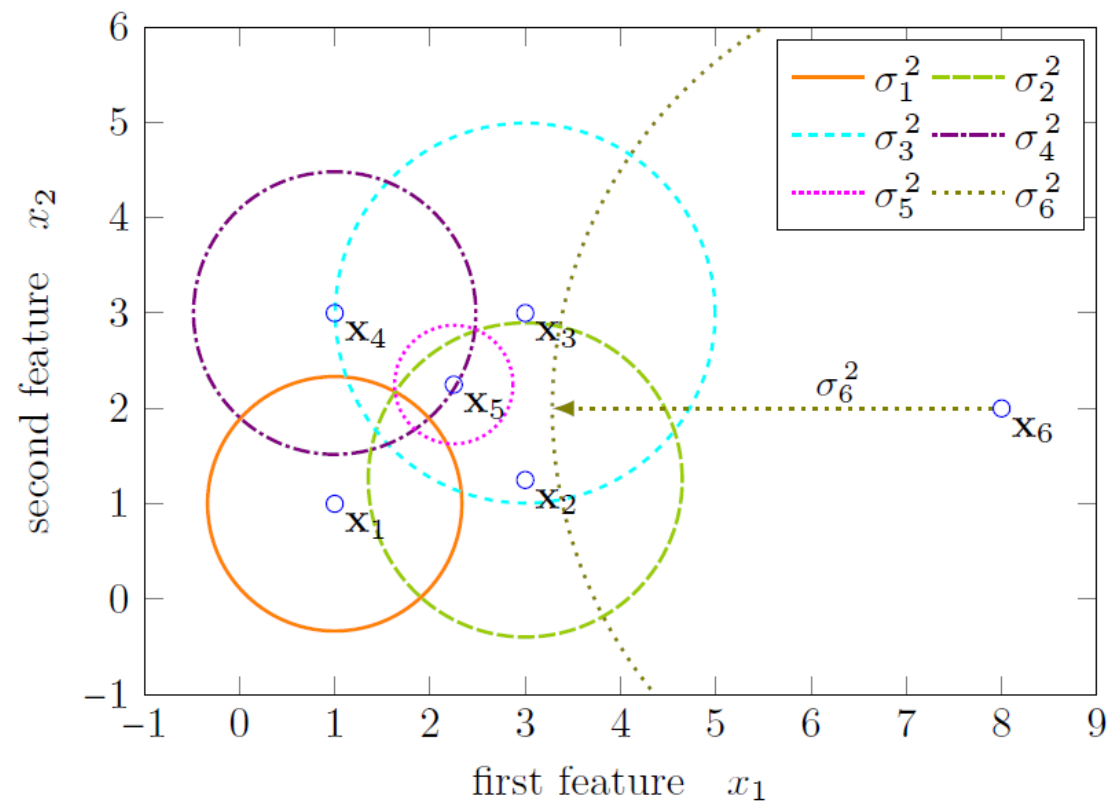
Transforming dissimilarity into affinity

**Definition (Affinity)**

$$a_{ij} = \begin{cases} \exp(-d_{ij}^2 / 2\sigma_i^2) & \text{if } i \neq j \\ 0 & \text{if } i = j \end{cases}$$



The radii of the circles correspond to the variance for the data points



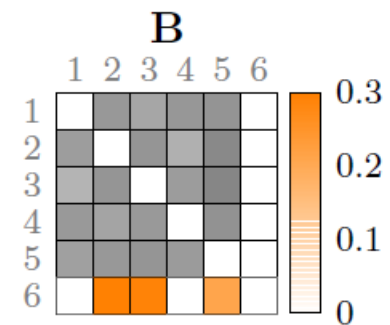
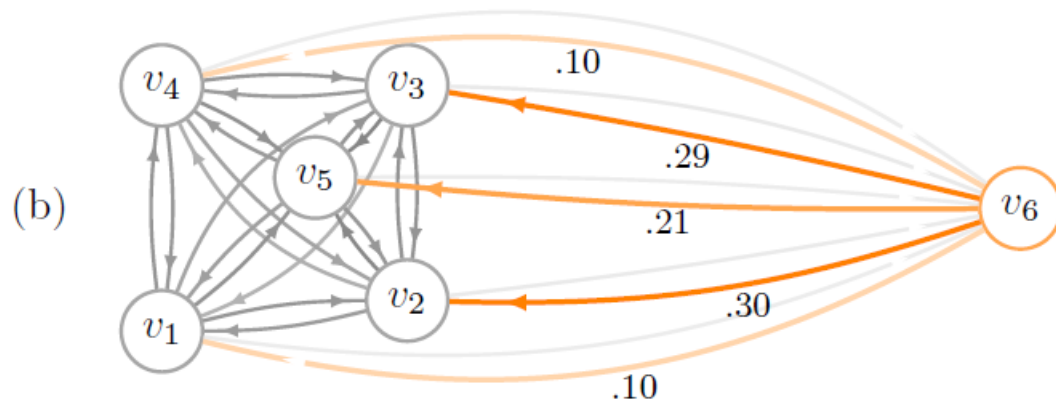
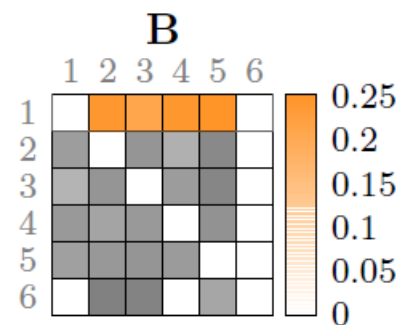
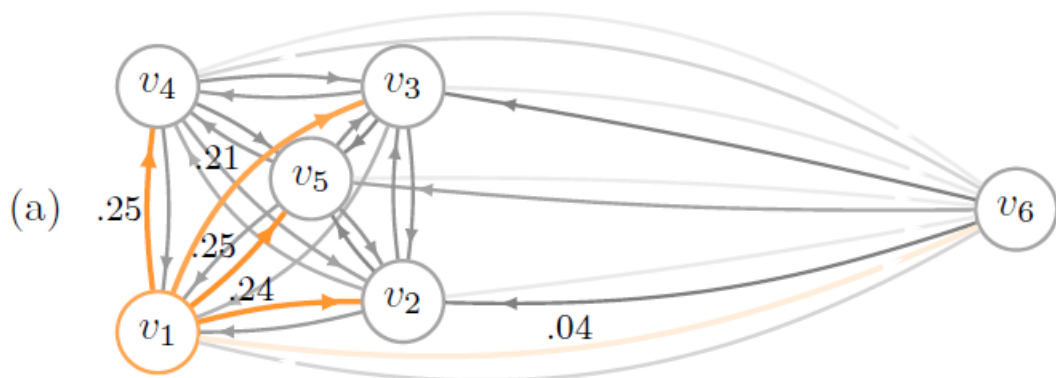
# Stochastic Neighbour Graphs based on affinities

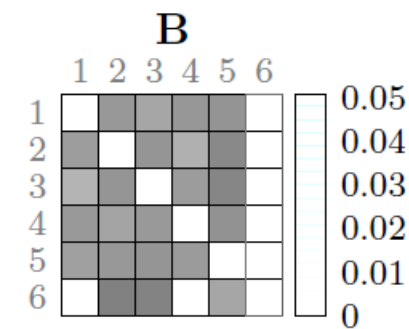
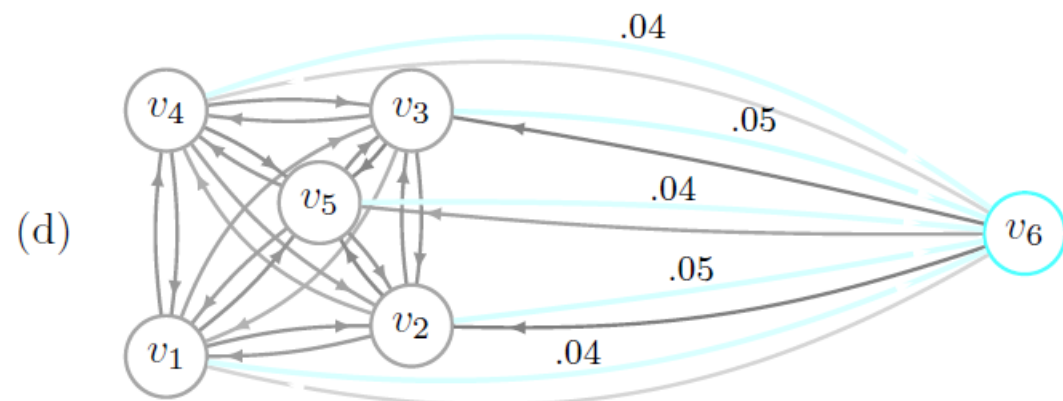
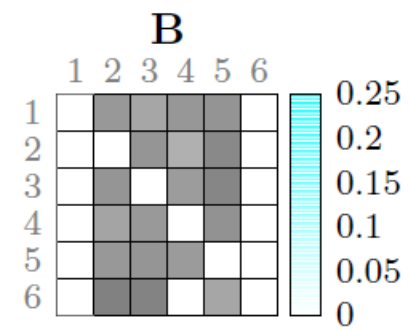
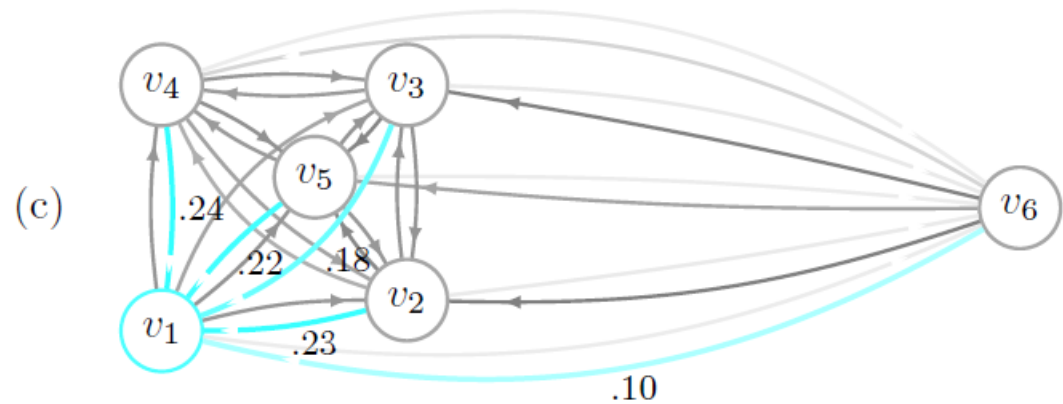


## Binding probabilities

$$b_{ij} \equiv p(i \rightarrow j \in \mathcal{E}_G) \propto a_{ij}$$

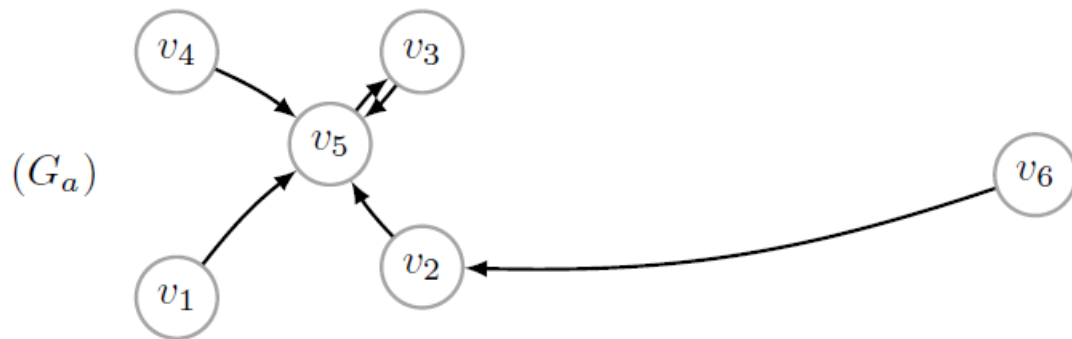
$$b_{ij} = \frac{a_{ij}}{\sum_{k=1}^n a_{ik}}$$







# Being an outlier given one SNG

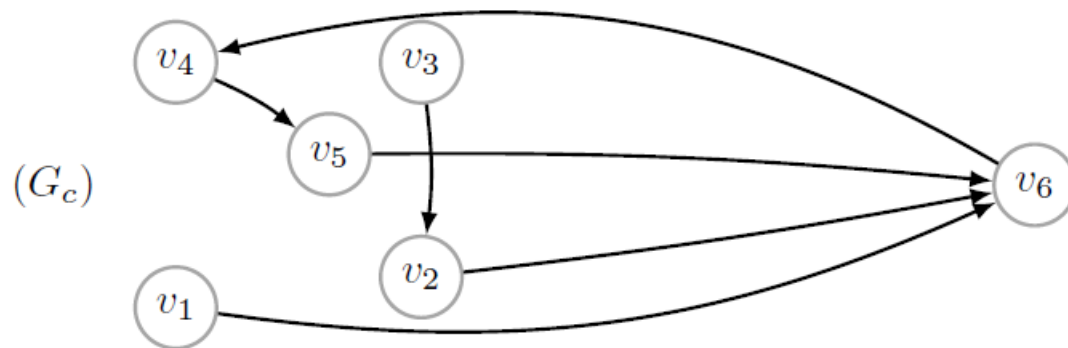


$$p(G_a) = 3.931 \cdot 10^{-4}$$

$$\mathcal{C}_O | G_a = \{x_1, x_4, x_6\}$$

$$\mathcal{C}_O | G = \{x_i \in \mathbf{X} \mid \text{deg}^-_G(v_i) = 0\}$$

$$\mathcal{C}_O | G = \{x_i \in \mathbf{X} \mid \nexists v_j \in \mathcal{V} : j \rightarrow i \in \mathcal{E}_G\}$$

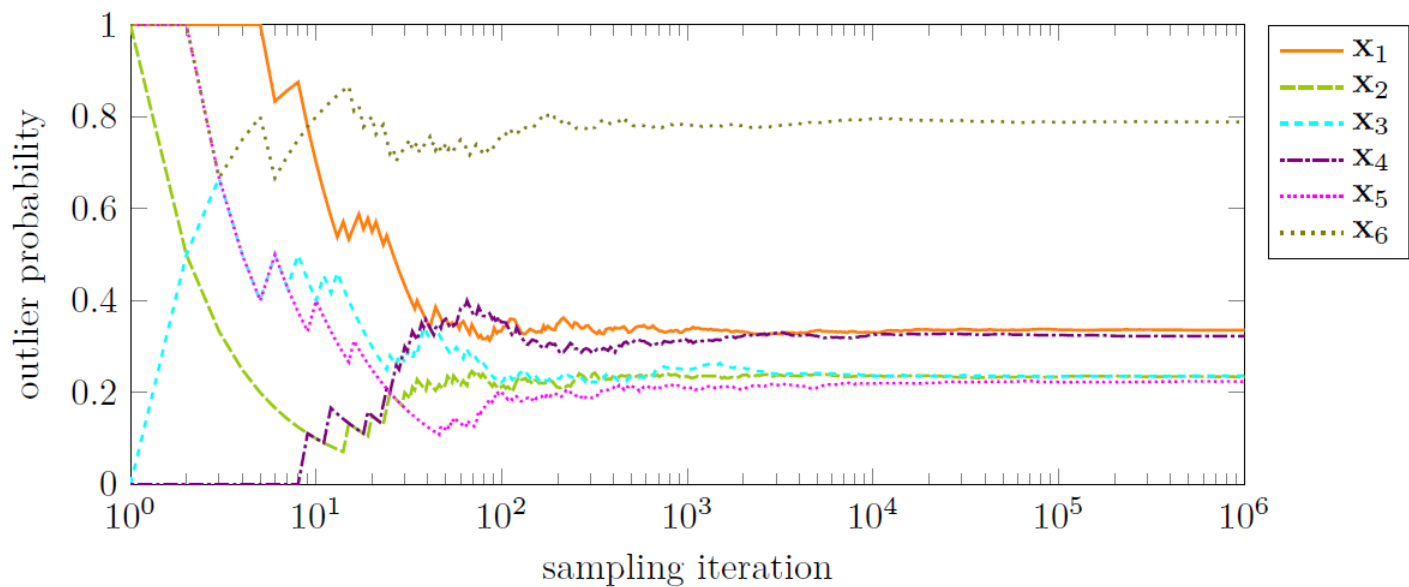


$$p(G_c) = 5.950 \cdot 10^{-7}$$

$$\mathcal{C}_O | G_c = \{x_1, x_3\}$$



$$p(\mathbf{x}_i \in \mathcal{C}_O) = \lim_{S \rightarrow \infty} \frac{1}{S} \sum_{s=1}^S \mathbb{1}\{\mathbf{x}_i \in \mathcal{C}_O | G^{(s)}\} \quad , \quad G^{(s)} \sim P(\mathcal{G}) \quad p(G) = \prod_{i \rightarrow j \in \mathcal{E}_G} b_{ij}$$



$$\begin{aligned} p(\mathbf{x}_i \in \mathcal{C}_O) &= \sum_{G \in \mathcal{G}} \mathbb{1}\{\mathbf{x}_i \in \mathcal{C}_O | G\} \cdot p(G) \\ &= \sum_{G \in \mathcal{G}} \mathbb{1}\{\mathbf{x}_i \in \mathcal{C}_O | G\} \cdot \prod_{q \rightarrow r \in \mathcal{E}_G} b_{qr} \end{aligned}$$

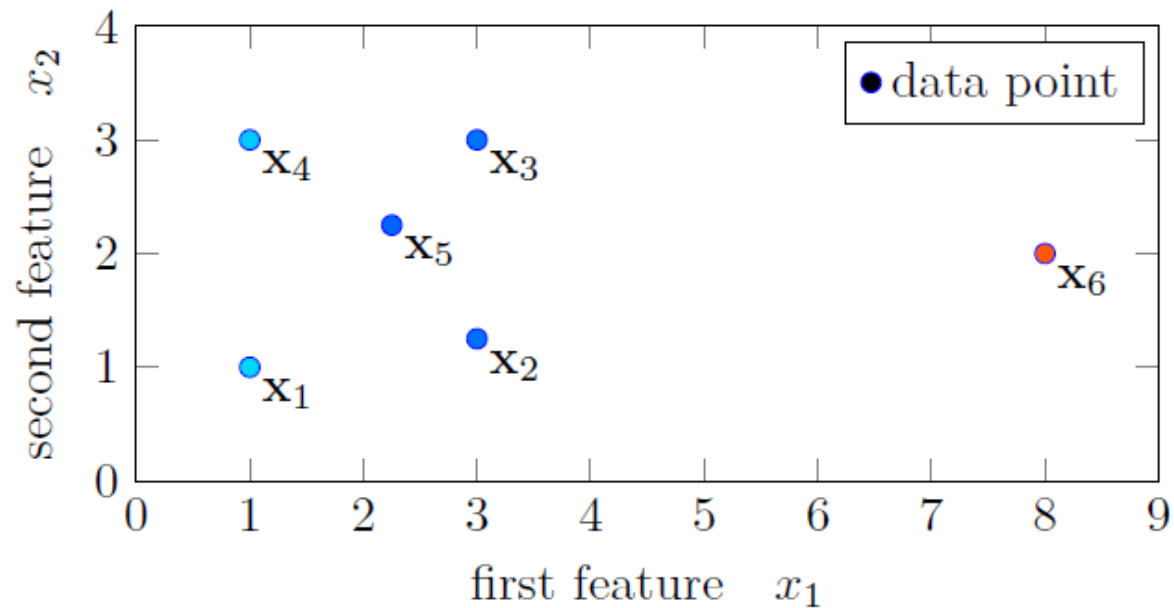
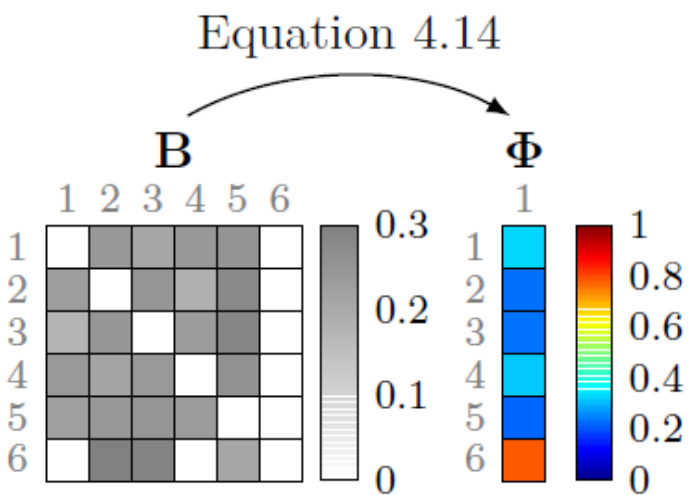
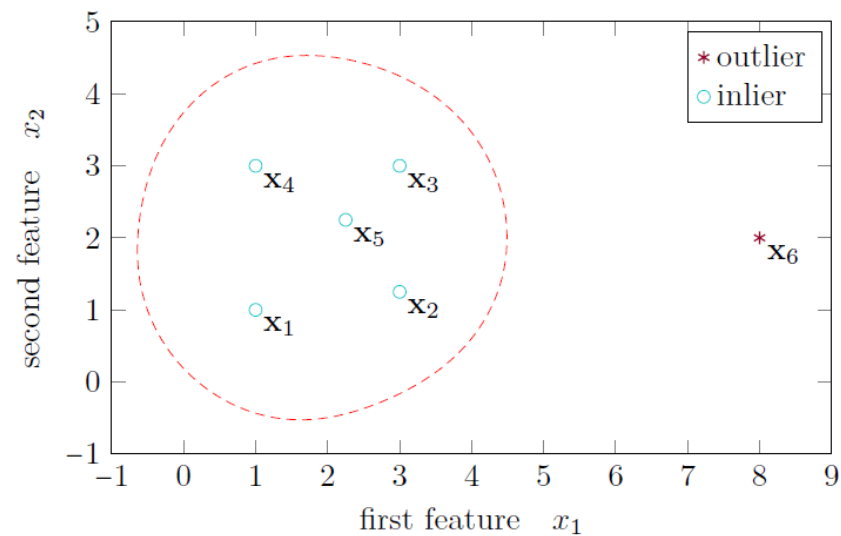
$5^6 = 15,625$  SNGs.

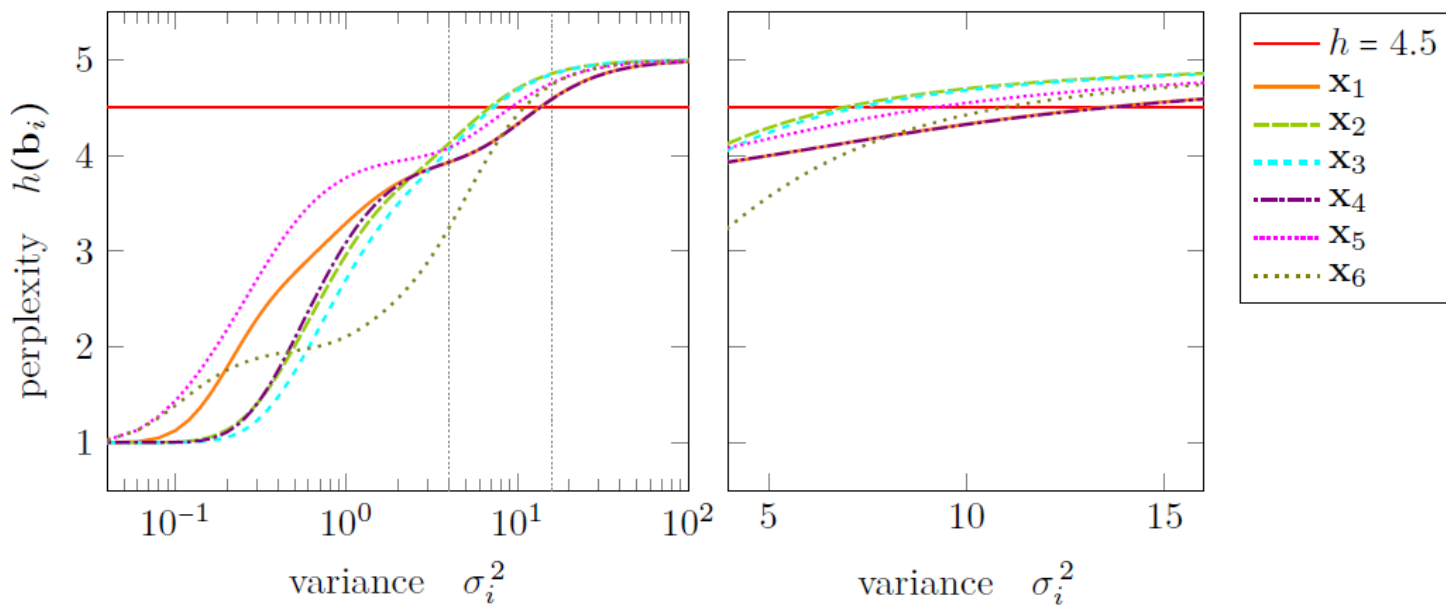
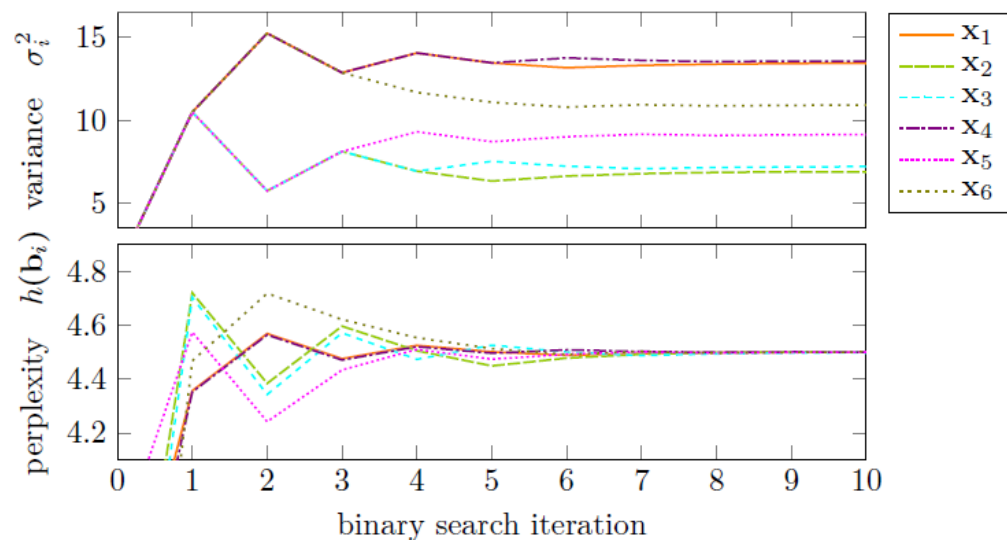
0:335; 0:235; 0:237; 0:323; 0:224, 0:788



$$p(\mathbf{x}_i \in \mathcal{C}_O) = \prod_{j \neq i} (1 - b_{ji}) \quad \varphi_{\text{SOS}}(\mathbf{x}_i) \equiv p(\mathbf{x}_i \in \mathcal{C}_O)$$

$$f_{\text{SOS}}(\mathbf{x}) = \begin{cases} \text{outlier} & \text{if } \varphi_{\text{SOS}}(\mathbf{x}) > \theta \\ \text{inlier} & \text{if } \varphi_{\text{SOS}}(\mathbf{x}) \leq \theta \end{cases}$$







$$p(\mathbf{x}_i \in \mathcal{C}_O) = \prod_{j \neq i} (1 - b_{ji}) \quad 0.9^{1000} = 1.748 \times 10^{-46}$$

classification

outlier detection is an unsupervised problem

not only outliers

Matrix Sketch

better expression of the outlier probability?

Sparse Representation

