# The Wisdom of Crowds
## Cluster Ensemble

Muhammad Yousefnezhad

Member of ParNec & DQ groups

Department of Computer Science & Technology

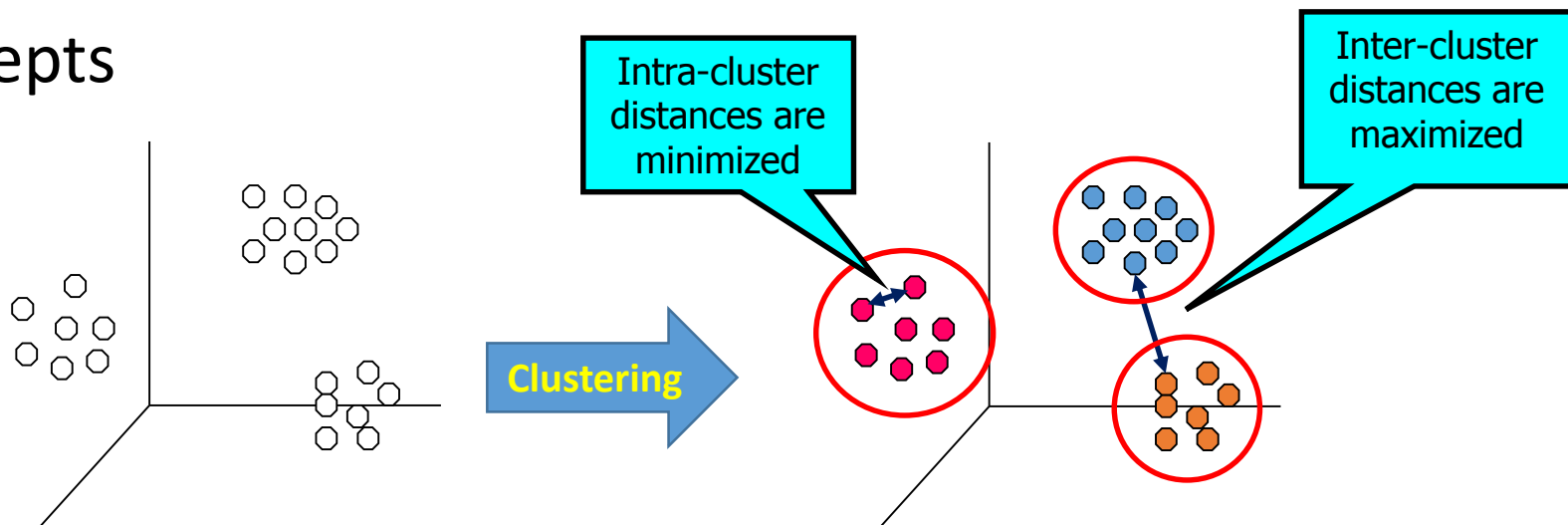Nanjing University of Aeronautics and Astronautics

email: myousefnezhad@nuaa.edu.cn
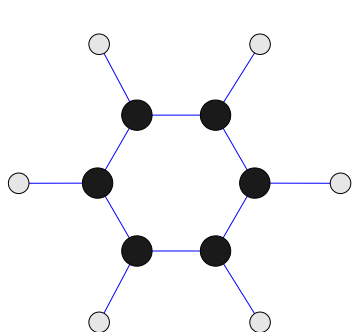
# Outlines

- Concepts of Clustering

- Basic Clustering Algorithms

- Cluster Ensemble Selection

- Wisdom of Crowds Theory

- Wisdom of Crowds Cluster Ensemble (WOCCE)

- Experimental Results

- Feature Works
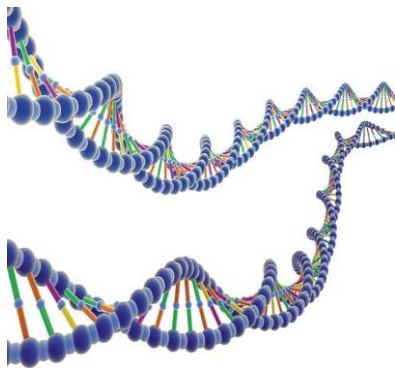
# Concepts of Clustering

- ## Concepts



Intra-cluster distances are minimized

Inter-cluster distances are maximized

Clustering

- ## Applications



**Chemical Elements**

**Gens**

**Brain Function Analysis**

Pictures

Voices

**Web & Social Networks**

# Basic Clustering Algorithms



MiniBatchKMeans  AffinityPropagation  MeanShift  SpectralClustering  Ward  AgglomerativeClustering  DBSCAN

# Cluster Ensemble Selection



**Challenges**

- **Strategy of Selection**
- **Generating Results**
- **Evaluation metrics**
- **Aggregation Methods**

Shannon $Entropy(X) = -\sum P(X)\log P(X)$

**Normal Mutual Information**
$$NMI(P^a, P^b) = \frac{-2\sum_{i=1}^{k_a}\sum_{j=1}^{k_b} n_{ij}^{ab} \log\left(\frac{n_{ij}^{ab} \times n}{n_i^a \times n_j^b}\right)}{\sum_{i=1}^{k_a} n_i^a \log\left(\frac{n_i^a}{n}\right) + \sum_{j=1}^{k_b} n_j^b\left(\frac{n_j^b}{n}\right)}$$

**The Wisdom of Crowds Cluster Ensemble Selection**

# The Wisdom of Crowds (WOC) Theory

**Surowiecki J. (2004), "The Wisdom of Crowds", Cover of mass market edition by Anchor, *ISBN: 978-0385503860*.**

❑ **Background**

- ○ **Condorcet, (1785), "Jury Theorem", Essay on the Application of Analysis to the Probability of Majority Decisions.**
- ○ **Galton, (1907), "Vox Populi", Nature, vol. 75, pp. 450-451.**

❑ **Criteria**

- ○ **Diversity of opinion**
  - ▪ **Each person should have private information even if it's just an eccentric interpretation of the known facts.**
- ○ **Independence**
  - ▪ **People's opinions aren't determined by the opinions of those around them.**
- ○ **Decentralization**
  - ▪ **People are able to specialize and draw on local knowledge.**
- ○ **Aggregation**
  - ▪ **Some mechanism exists for turning private judgments into a collective decision.**

❑ **Application**

- ○ **Delphi Method, Crowd Sourcing, Machine Learning, and etc.**

# Wisdom of Crowds Cluster Ensemble

The WOCCE framework (with feedback mechanism)

# WOCCE – Diversity

$$APMM(C_i^a, P^{b*}) = \frac{-2\log\left(\frac{n}{n_i^a}\right)\sum_{j=1}^{k_{b*}} n_j^{b*}}{n_i^a \log\left(\frac{n_i^a}{n}\right) + \sum_{j=1}^{k_{b*}} n_j^{b*} \log\left(\frac{n_j^{b*}}{n}\right)}$$

$$AAPMM(P) = \frac{1}{M}\sum_{j=1}^{M} APMM\left(C, P_j^{b*}\right)$$

$$A3(P) = \frac{1}{n}\sum_{i=1}^{k} n_i \times AAPMM(P_i)$$

$$Diversity(P) = 1 - A3(P)$$

**Diversity Condition** $\qquad Diversity(P) \geq dT$

# WOCCE – Independency

| Data | Alg1 | Alg2 | Class |
|------|------|------|-------|
| X1 | 1 | 1 | 1 |
| X2 | 2 | 2 | 2 |
| X3 | 1 | 1 | 1 |
| X4 | 2 | 1 | 1 |
| X5 | 2 | 2 | 2 |
| X6 | 2 | 2 | 2 |
| X7 | 1 | 1 | 1 |

Diversity(Alg1,Alg2) = 1 – NMI(Alg1,Alg2) ~ 0

**IF** Type(Alg1) <> Type(Alg2) **THEN** Accept

**ELSIF** [Basic Parameter(Alg1) - Basic Parameter(Alg2)] >> 0 **THEN** Accept

**ELSE** omit the results

**Independency Condition** $\Rightarrow$ $Independence(A\lg_i) \geq iT$

$$Independence(P) = \frac{1}{M} \sum_{i=1}^{M} \text{BPI}(P, P_i)$$

**Function BPI (P1, P2) Return Result**
    **If (Algorithm-Type (P1) == Algorithm-Type (P2) then**
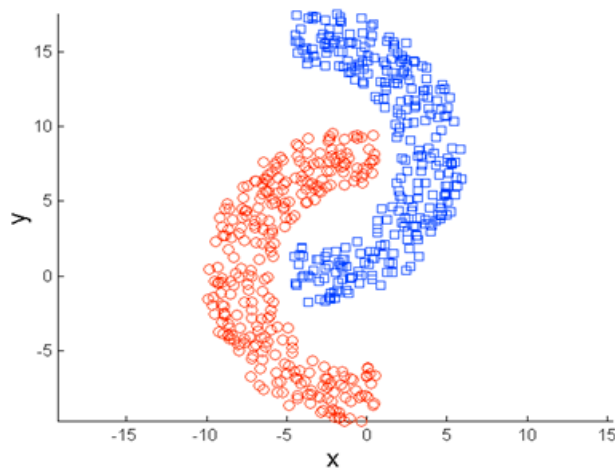        **Result = 1 - Likeness (Basic-Parameter (P1), Basic-Parameter (P2))**
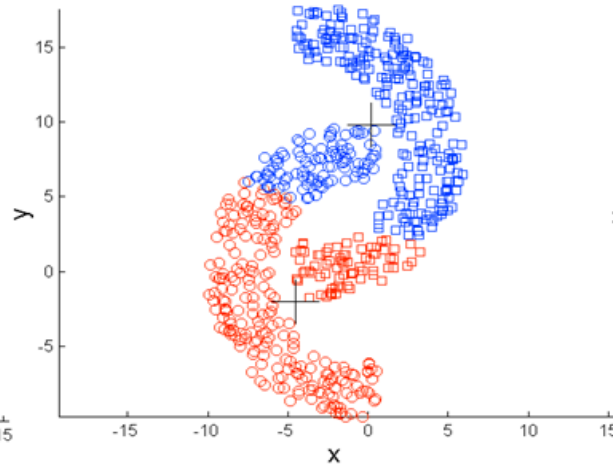    **Else**
        **Result = 1**
    **End if**
**End Function**

$$Likeness = 1 - (\frac{1}{MaxDis} \sum_{t=0}^{n} Sim_t)$$

# WOCCE – Decentralization



(a) Original Point     (b) 2-means Clustering     (c) 10-means Clustering

❑ Decentralization, a quality metric

o The number of basic algorithms participating in the crowd should be greater than one.

o The method of entering a primary algorithm into the crowd should ensure that the final result will not be affected by its errors. In other words, the decision making of the final ensemble should not be centralized.

o The threshold parameter **cT**, which we call the coefficient of decentralization, is a coefficient which is multiplied in the number of clusters. Every base algorithm clusters the dataset into at most **cT×k** clusters. i.e. it clusters the dataset into a number of clusters between **cT** to **cT×k**.

# WOCCE – Aggregation

○ **Evidence Accumulation Clustering (EAC)**



**Selected basic results (Wised Crowds)**

| Data | Alg(1) | ... | Alg(m) |
|------|--------|-----|--------|
| X1 | | | |
| X2 | | | |
| X3 | | | |
| ⋮ | | | |
| Xn | | | |

$$C(i, j) = \frac{n_{i,j}}{m_{i,j}}$$

**Co-Association Matrix**

$$\begin{array}{c c} & \begin{matrix} X1 & X2 & \cdots & Xn \end{matrix} \\ \begin{matrix} X1 \\ X2 \\ \vdots \\ Xn \end{matrix} & \begin{bmatrix} c_{1,1} & c_{1,2} & \cdots & c_{1,n} \\ c_{2,1} & c_{2,2} & \cdots & c_{2,n} \\ \vdots & \vdots & \cdots & \vdots \\ c_{n,1} & c_{n,2} & \cdots & c_{n,n} \end{bmatrix} \end{array}$$

**Final Result**

| Class | Data |
|-------|------|
| Lx | X1 |
| Ly | X2 |
| ... | ... |
| Lx | Xn |

**Dendrogram**

(X1, X2, X5, X3, X4)
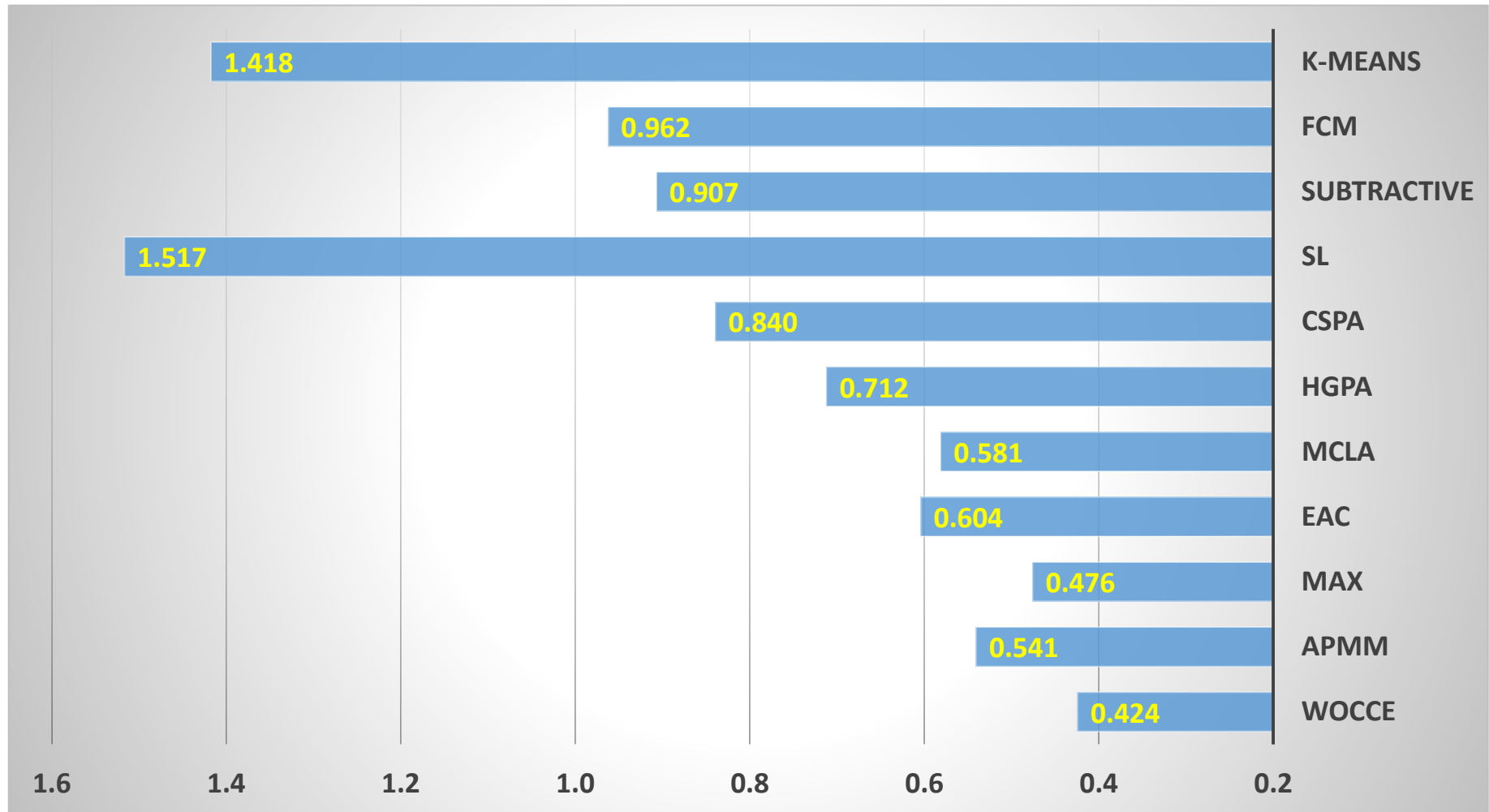
# Experimental Results (Performance)

o Average of Performance: Calculated based on 10-times independence runtime on 14 standard data sets from UCI

# Experimental Results (Standard error)

o Average of Error: Calculated based on 10-times independence runtime on 14 standard data sets from UCI



Horizontal bar chart (axis from 1.6 on the left to 0.2 on the right):

| Method | Value |
|---|---|
| K-MEANS | 1.418 |
| FCM | 0.962 |
| SUBTRACTIVE | 0.907 |
| SL | 1.517 |
| CSPA | 0.840 |
| HGPA | 0.712 |
| MCLA | 0.581 |
| EAC | 0.604 |
| MAX | 0.476 |
| APMM | 0.541 |
| WOCCE | 0.424 |

**The Wisdom of Crowds Cluster Ensemble Selection**

# Feature Works

❑ **Publications**

o **M. Yousefnezhad**, H. Alizadeh, B. Minaei-Bidgoli, (2013), "New cluster ensemble selection method based on Diversity and Independent metrics", The 5th Conference on Information and Knowledge Technology (IKT'13), Shiraz, Iran, May 22-24/2013, Index by I.S.C (in Persian)

o H. Alizadeh, **M. Yousefnezhad**, B. Minaei-Bidgoli, (2015), "Wisdom of Crowds Cluster Ensemble", Intelligent Data Analysis, IOS Press, Vol. 19(3).

• **M. Yousefnezhad**, A. Reihanian, B. Minaei-Bidgoli, "A new selection strategy for selective cluster ensemble based on Diversity and Independency", Submitted in Apply Soft Computing (ASOC).

• **M. Yousefnezhad**, D. Zhang, A. Reihanian, "A new framework for Cluster Ensemble Selection based on Diversity and Graph based Independency", Submitted in PA-KDD 2015.

• **M. Yousefnezhad**, D. Zhang, "Wised Semi-Supervised Cluster Ensemble Selection: a new framework for selecting and combing multiple partitions based on prior knowledge", Submitted in IEEE Transaction on Pattern Analysis and Machine Intelligence (TPAMI).

❑ **Other ideas**

   o Introducing mathematician concept of Independency in CES
      ▪ **It can be proved by probability concept in statistics**
      ▪ **It can be proved by learning automata**
   o Analyzing the effects of other basic algorithms in WOCCE framework
   o Analyzing the effects of feedback mechanism in other CES techniques
   o Developing WOCCE for semi-supervised learning applications

# Thank You

# Q&A