# Clusters

**Sumit Mishra**

Department of Computer Science & Engineering
Indian Institute of Information Technology Guwahati

Place: IIT Patna
Date: 4 April 2019

# Outline

# Outline

# Number of Clustering Results Obtained from 2 Records

Given 2 points/records, number of possible clustering results?

# Number of Clustering Results Obtained from 2 Records

Given 2 points/records, number of possible clustering results?

## For 2 records

1. $\{\{r_1, r_2\}\}$
2. $\{\{r_1\}, \{r_2\}\}$

# Number of Clustering Results Obtained from 3 Records

Given 3 points/records, number of possible clustering results?

# Number of Clustering Results Obtained from 3 Records

Given 3 points/records, number of possible clustering results?

## For 3 records

1. $\{\{r_1, r_2, r_3\}\}$
2. $\{\{r_1\}, \{r_2, r_3\}\}$
3. $\{\{r_2\}, \{r_1, r_3\}\}$
4. $\{\{r_3\}, \{r_1, r_2\}\}$
5. $\{\{r_1\}, \{r_2\}, \{r_3\}\}$

# Number of Clustering Results Obtained from 4 Records

Given 4 points/records, number of possible clustering results?

# Number of Clustering Results Obtained from 4 Records

Given 4 points/records, number of possible clustering results?

## For 4 records

1. $\{\{r_1, r_2, r_3, r_4\}\}$
2. $\{\{r_1, r_2\}, \{r_3, r_4\}\}$
3. $\{\{r_1, r_3\}, \{r_2, r_4\}\}$
4. $\{\{r_1, r_4\}, \{r_2, r_3\}\}$
5. $\{\{r_1\}, \{r_2, r_3, r_4\}\}$
6. $\{\{r_2\}, \{r_1, r_3, r_4\}\}$
7. $\{\{r_3\}, \{r_1, r_2, r_4\}\}$
8. $\{\{r_4\}, \{r_1, r_2, r_3\}\}$
9. $\{\{r_1\}, \{r_2\}, \{r_3, r_4\}\}$
10. $\{\{r_1\}, \{r_3\}, \{r_2, r_4\}\}$
11. $\{\{r_1\}, \{r_4\}, \{r_2, r_3\}\}$
12. $\{\{r_2\}, \{r_3\}, \{r_1, r_4\}\}$
13. $\{\{r_2\}, \{r_4\}, \{r_1, r_3\}\}$
14. $\{\{r_3\}, \{r_4\}, \{r_1, r_2\}\}$
15. $\{\{r_1\}, \{r_2\}, \{r_3\}, \{r_4\}\}$

# Number of Clustering Results Obtained from $n$ Records

Given $n$ points/records, number of possible clustering results?

# Number of Clustering Results Obtained from $n$ Records

Given $n$ points/records, number of possible clustering results?

## Using Bell number [1]

$$B_{n+1} = \sum_{k=0}^{n} \binom{n}{k} B_k, \qquad B_0 = B_1 = 1 \tag{1}$$

# Number of Clustering Results Obtained from $n$ Records

Given $n$ points/records, number of possible clustering results?

## Using Bell number [1]

$$B_{n+1} = \sum_{k=0}^{n} \binom{n}{k} B_k, \qquad B_0 = B_1 = 1 \tag{1}$$

## Using Stirling numbers of the second kind [8]

$$B_n = \sum_{k=0}^{n} \left\{ n \atop k \right\} \tag{2}$$

Here, the Stirling number $\left\{ n \atop k \right\}$ is the number of ways to partition a set of cardinality $n$ into exactly $k$ nonempty subsets.

# Bell Triangle

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| **1** | | | | | | | | | |
| 1 | **2** | | | | | | | | |
| 2 | 3 | **5** | | | | | | | |
| 5 | 7 | 10 | **15** | | | | | | |
| 15 | 20 | 27 | 37 | **52** | | | | | |
| 52 | 67 | 87 | 114 | 151 | **203** | | | | |
| 203 | 255 | 322 | 409 | 523 | 674 | **877** | | | |
| 877 | 1080 | 1335 | 1657 | 2066 | 2589 | 3263 | **4140** | | |
| 4140 | 5017 | 6097 | 7432 | 9089 | 11155 | 13744 | 17007 | **21147** | |
| 21147 | 25287 | 30304 | 36401 | 43833 | 52922 | 64077 | 77821 | 94828 | **115975** |

Figure 1 : Bell triangle for 10 records.

# Outline

# Arrangements of $n$ Records in Clustering Results

Given $n$ points/records, how the records are arranged in the clustering result?

# Arrangements of $n$ Records in Clustering Results

Given $n$ points/records, how the records are arranged in the clustering result? Bell Polynomial [1, 2]

$$B_n(x_1, \ldots, x_n) = \sum_{k=1}^{n} B_{n,k}(x_1, x_2, \ldots, x_{n-k+1}) \qquad (3)$$

$B_{n,k}(x_1, x_2, \ldots, x_{n-k+1})$ is the partial Bell polynomial and is given by Equation (4).

$$B_{n,k}(x_1, x_2, \ldots, x_{n-k+1}) = \sum \frac{n!}{j_1! j_2! \cdots j_{n-k+1}!} \left(\frac{x_1}{1!}\right)^{j_1} \left(\frac{x_2}{2!}\right)^{j_2} \cdots$$
$$\left(\frac{x_{n-k}}{(n-k)!}\right)^{j_{n-k}} \left(\frac{x_{n-k+1}}{(n-k+1)!}\right)^{j_{n-k+1}} \qquad (4)$$

where the sum is taken over all sequences $j_1, j_2, j_3, \ldots, j_{n-k+1}$ of non-negative integers such that $j_1 + j_2 + \cdots = k$ and $j_1 + 2j_2 + 3j_3 + \cdots = n$.

# Arrangements of 3 Records in Clustering Results

$3^{rd}$ complete Bell polynomial is given by Equation (5).

$$\begin{aligned}
B_3(x_1, x_2, x_3) &= \sum_{k=1}^{3} B_{3,k}(x_1, x_2, \ldots, x_{3-k+1}) \\
&= B_{3,1}(x_1, x_2, x_3) + B_{3,2}(x_1, x_2) + B_{3,3}(x_1) \\
&= (x_3) + (3x_1 x_2) + \left(x_1^3\right)
\end{aligned} \tag{5}$$

In simple terms Equation (5) can be written as Equation (6).

$$B_3(x) = x + 3x^2 + x^3 \tag{6}$$

# Arrangements of 3 Records in Clustering Results . . .

- One way to group 3 records in single cluster,
  1. $\{\{r_1, r_2, r_3\}\}$
- Three ways to group 3 records in 2 clusters,
  1. $\{\{r_1\}, \{r_2, r_3\}\}$
  2. $\{\{r_2\}, \{r_1, r_3\}\}$
  3. $\{\{r_3\}, \{r_1, r_2\}\}$
- One way to group 3 records in 3 clusters.
  1. $\{\{r_1\}, \{r_2\}, \{r_3\}\}$

# Arrangements of 4 Records in Clustering Results

$4^{th}$ complete Bell polynomial is given by Equation (7).

$$
\begin{aligned}
B_4(x_1, x_2, x_3, x_4) &= \sum_{k=1}^{4} B_{4,k}(x_1, x_2, \ldots, x_{4-k+1}) \\
&= B_{4,1}(x_1, x_2, x_3, x_4) + B_{4,2}(x_1, x_2, x_3) + B_{4,3}(x_1, x_2) + B_{4,4}(x_1) \\
&= (x_4) + \left(3x_2^2 + 4x_1x_3\right) + \left(6x_1^2x_2\right) + \left(x_1^4\right) \qquad (7)
\end{aligned}
$$

In simple terms Equation (7) can be written as Equation (8).

$$
B_4(x) = x + 7x^2 + 6x^3 + x^4 \qquad (8)
$$

# Arrangements of 4 Records in Clustering Results . . .

- One way to group 4 records in single cluster,
  1. $\{\{r_1, r_2, r_3, r_4\}\}$
- Seven $(3 + 4)$ ways to group 4 records in 2 clusters,
  1. $\{\{r_1, r_2\}, \{r_3, r_4\}\}$
  2. $\{\{r_1, r_3\}, \{r_2, r_4\}\}$
  3. $\{\{r_1, r_4\}, \{r_2, r_3\}\}$
  4. $\{\{r_1\}, \{r_2, r_3, r_4\}\}$
  5. $\{\{r_2\}, \{r_1, r_3, r_4\}\}$
  6. $\{\{r_3\}, \{r_1, r_2, r_4\}\}$
  7. $\{\{r_4\}, \{r_1, r_2, r_3\}\}$
- Six ways to group 4 records in 3 clusters,
  1. $\{\{r_1\}, \{r_2\}, \{r_3, r_4\}\}$
  2. $\{\{r_1\}, \{r_3\}, \{r_2, r_4\}\}$
  3. $\{\{r_1\}, \{r_4\}, \{r_2, r_3\}\}$
  4. $\{\{r_2\}, \{r_3\}, \{r_1, r_4\}\}$
  5. $\{\{r_2\}, \{r_4\}, \{r_1, r_3\}\}$
  6. $\{\{r_3\}, \{r_4\}, \{r_1, r_2\}\}$
- One way to group 4 records in 4 clusters.
  1. $\{\{r_1\}, \{r_2\}, \{r_3\}, \{r_4\}\}$

# Outline

# Integer Partition

A partition of a positive integer $n$ is defined to be a sequence of positive integers whose sum is $n$ [3].

## Integer partition for $n = 3$

1. 3
2. $1 + 2$
3. $1 + 1 + 1$

## Integer partition for $n = 4$

1. 4
2. $2 + 2$
3. $1 + 3$
4. $1 + 1 + 2$
5. $1 + 1 + 1 + 1$

# Integer Partition . . .

An asymptotic expression for number of partitions of an integer n is given by Equation (9) [5, 7].

$$p(n) \sim \frac{1}{4n\sqrt{3}} \exp\left(\pi\sqrt{\frac{2n}{3}}\right) \text{ as } n \to \infty. \tag{9}$$

# Integer Partition and Clustering Results for 3 Records

| Serial No. | Number of Cluster | Possible Partition | Partial Bell Polynomial | No. of Clustering Results Corresponding to a Partition |
|---|---|---|---|---|
| 1 | 1 | $\{3\}$ | $x_3$ | 1 |
| 2 | 2 | $\{1, 2\}$ | $3x_1 x_2$ | 3 |
| 3 | 3 | $\{1, 1, 1\}$ | $x_1^3$ | 1 |
| Total number of clustering results $B_3$ | | | | 5 |

Table 1 : Possible partition of 3 records along with the number of clustering results corresponding to each partition.

# Integer Partition and Clustering Results for 4 Records

| Serial No. | Number of Cluster | Possible Partition | Partial Bell Polynomial | No. of Clustering Results Corresponding to a Partition |
|---|---|---|---|---|
| 1 | 1 | $\{4\}$ | $x_4$ | 1 |
| 2 | 2 | $\{2,2\}$ | $3x_2^2 + 4x_1x_2$ | 3 |
| 3 | | $\{1,3\}$ | | 4 |
| 4 | 3 | $\{1,1,2\}$ | $6x_1^2x_2$ | 6 |
| 5 | 4 | $\{1,1,1,1\}$ | $x_1^4$ | 1 |
| Total number of clustering results $B_4$ | | | | 15 |

Table 2 : Possible partition of 4 records along with the number of clustering results corresponding to each partition.

# Outline

# Motivation



## Observation

- There are two different authors with the same name.

# Problem Statement

## Motivation

- DBLP contains more than 1300 papers published by authors having the name "Wei Wang".
- How many authors having the same name?
- What is the categorization of the papers?

## DEFINITION: Entity Matching [6]

Given a name 'pName' and a set of records $\mathbb{R} = \{r_1, r_2, \ldots, r_n\}$ corresponding to name 'pName', the Entity Matching is to divide the records in $\mathbb{R}$ into different clusters such that the following holds.

- All the records in a cluster belong to an entity.
- All the records by an entity should be in a single cluster.

# Outline

# Evolutionary Algorithm based Approach

### Motivation

- Modeled entity matching problem as an optimization problem.
- Used evolutionary algorithm as an optimization framework.
- Single as well as multiple objectives are considered.
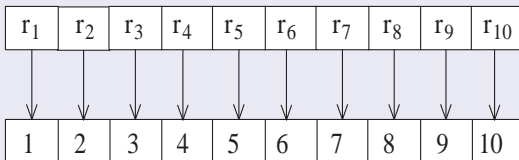
---

**Algorithm 1** GENETIC ALGORITHM

---

1: Initialize population $P$
2: Evaluate the fitness of all individuals
3: Select fitter individuals for reproduction
4: Apply recombination among individuals
5: Mutate individuals
6: Evaluate the fitness of the modified individuals
7: Generate a new population

---

# Chromosome Initialization

## Record Encoding

| $r_1$ | $r_2$ | $r_3$ | $r_4$ | $r_5$ | $r_6$ | $r_7$ | $r_8$ | $r_9$ | $r_{10}$ |
|---|---|---|---|---|---|---|---|---|---|

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|

## Chromosome Representation

- Represents the representative of the cluster, *i.e.*, one of the elements from the cluster.
- The $K$ length of *Chromosome* means that there are $K$ clusters.
- The elements in the *Chromosome* are distinct.

| 2 | 5 | 8 | 10 |
|---|---|---|---|

**Cluster representative**

# Population Initialization and Assignment of Records

## Population Initialization

- The number of clusters is unknown.
- Size of *chromosome* varies between 2 and $n - 1$.
- To initialize population, each *chromosome* in the population is initialized.

## Assignment of Records

- $\mathbb{R} = \{r_1, r_2, \ldots, r_{10}\}$
- Chromosome $= \{r_2, r_5, r_8, r_{10}\}$
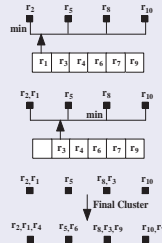- $\mathbb{UR} = \{r_1, r_3, r_4, r_6, r_7, r_9\}$
- Distance Measure



Figure 2 : Assignment of records.

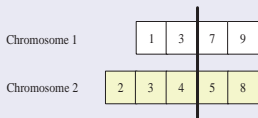# Crossover and Mutation

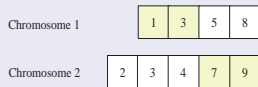## Crossover Operation on Chromosomes



Figure 3 : Before crossover

Figure 4 : After crossover

## Mutation Operation on Chromosome

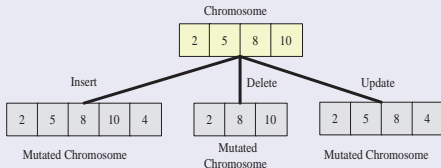Three mutation operations are considered.

i. Insert
ii. Delete
iii. Update



Figure 5 : Mutation operation

# Outline

Figure 6 : Flow of MOO-EMT

# NSGA-II



Figure 7 : NSGA-II Procedure [4]

# NSGA-II . . .

---

**Algorithm 2** NSGA-II

---

**Input:** $\mathbb{P}_t$: Population for $t^{th}$ generation
**Output:** $\mathbb{P}_{t+1}$: Population for $t + 1^{th}$ generation

1: $\mathbb{Q}_t \leftarrow$ Generate offspring population after crossover and mutation operations
2: $\mathbb{R}_t \leftarrow \mathbb{P}_t \cup \mathbb{Q}_t$      // Combine parent and offspring populations
3: $\mathcal{F} \leftarrow$ NON-DOMINATED-SORT$(\mathbb{R}_t)$      // $\mathcal{F} = \{F_1, F_2, \ldots, F_K\}$, set of non-dominated fronts in the decreasing order of their dominance nature
4: $\mathbb{P}_{t+1} \leftarrow \Phi$      // Initialize population for the next generation
5: $k \leftarrow 1$
6: **while** $|\mathbb{P}_{t+1}| + |F_k| \leq N$ **do**
7:     $\mathbb{P}_{t+1} \leftarrow P_{t+1} \cup F_k$      // Include $k^{th}$ non-dominated front in $\mathbb{P}_{t+1}$
8:     $k \leftarrow k + 1$      // Check the next front for inclusion in $\mathbb{P}_{t+1}$

    // Number of solutions to be included in population $\mathbb{P}_{t+1}$
9: $T = N - (|F_1| + |F_2| + \ldots + |F_{k-1}|)$
10: CROWDING-DISTANCE-ASSIGNMENT$(F_k)$     // Calculate crowding distance in $F_k$
11: Sort the solutions in $F_k$ based on crowding distance
12: $P_{t+1} \leftarrow \mathbb{P}_{t+1} \cup F_k [1 : T]$      // Choose the first $T$ solutions from $F_k$

# NSGA-II . . .

---

**Algorithm 3** CROWDING-DISTANCE-ASSIGNMENT($I$)

---

**Input:** $I$ : Non-dominated front

**Output:** $I$ : Crowded distance assignment to each solution in $I$

1:   $l \leftarrow |I|$        // Number of solutions in non-dominated front $I$

2:   **for** $i \leftarrow 1$ to $l$ **do**

3:     $I[i]_{\text{distance}} \leftarrow 0$        // Initialize the crowding distance

4:   **for each** objective m **do**

5:     $I \leftarrow$ SORT $(I, m)$ // Sort the solutions in $I$ in descending order using $m^{th}$ objective

6:     $I[1]_{\text{distance}} \leftarrow \infty$        // Set the value for boundary points

7:     $I[l]_{\text{distance}} \leftarrow \infty$        // Set the value for boundary points

8:     **for** $i \leftarrow 2$ to $l - 1$ **do**

9:       $I[i]_{\text{distance}} \leftarrow I[i]_{\text{distance}} + (I[i+1].m - I[i-1].m) / (f_m^{\max} - f_m^{\min})$

---

## Conclusions

- Discussed the number of clusters and how it can be obtained.
- Also discussed entity matching problem and how it can be solved using multi-objective clustering.

# References I

[1] HW Becker and John Riordan.
The arithmetic of bell and stirling numbers.
*American journal of Mathematics*, pages 385–394, 1948.

[2] Eric Temple Bell.
Partition polynomials.
*Annals of Mathematics*, pages 38–46, 1927.

[3] David M Burton.
*Elementary number theory*.
Tata McGraw-Hill Education, 2006.

[4] Kalyanmoy Deb, Amrit Pratap, Sameer Agarwal, and T. Meyarivan.
A Fast and Elitist Multiobjective Genetic Algorithm: NSGA-II.
*IEEE Transactions on Evolutionary Computation*, 6(2):182–197, April 2002.

[5] Godfrey H Hardy and Srinivasa Ramanujan.
Asymptotic formulæ in combinatory analysis.
*Proceedings of the London Mathematical Society*, 2(1):75–115, 1918.

[6] Shaohua Li, Gao Cong, and Chunyan Miao.
Author name disambiguation using a new categorical distribution similarity.
In *Machine learning and knowledge discovery in databases*, pages 569–584. Springer, 2012.

# References II

[7] Hans Rademacher.
On the expansion of the partition function in a series.
*Annals of Mathematics,* pages 416–422, 1943.

[8] Basil Cameron Rennie and Annette Jane Dobson.
On Stirling Numbers of the Second Kind.
*Journal of Combinatorial Theory,* 7(2):116–121, 1969.

# Thank you!