

# A Comparison of Knee Strategies for Hierarchical Spatial Clustering

Brian J. Ross

Department of Computer Science  
Brock University  
St. Catharines, Ontario, Canada  
[bross@brocku.ca](mailto:bross@brocku.ca)

IEA-AIE 2018

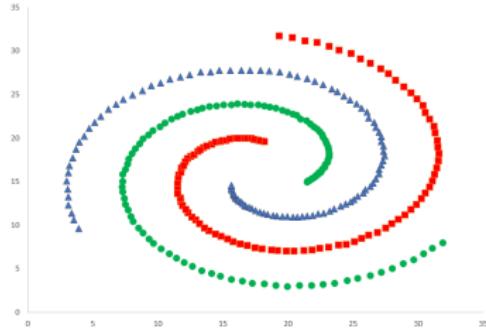
# Overview

- Introduction
- Setup
  - ▶ Clustering algorithms
  - ▶ Knee heuristics
  - ▶ Data sets
- Results
- Conclusion

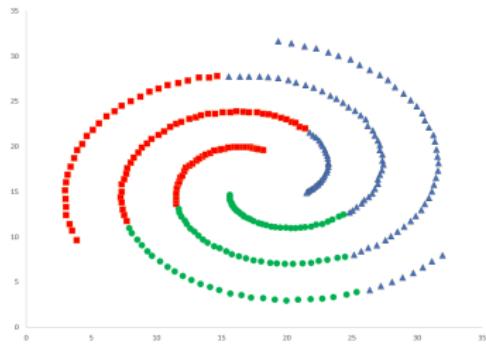
# Introduction

- **Hierarchical clustering:** automatic grouping of data into sets with similar characteristics
  - ▶ Incrementally build clusters, from  $K$  clusters of 1 point each, to 1 cluster of all  $K$  points.
  - ▶ Dendrogram represents incremental cluster creation by clustering algorithm.
  - ▶ Determine optimal clustering afterwards.
- **Clustering of 2D spatial data:** group planar points into sets
- **Computational limitations**
  - ▶ Clustering is in general NP-complete.
  - ▶ Optimality is often subjective and ill-defined.

# Introduction

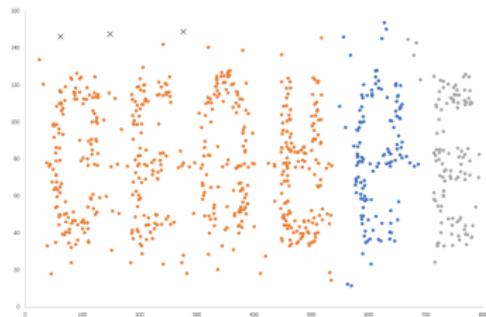


Spiral and single-linkage clustering, K=3.

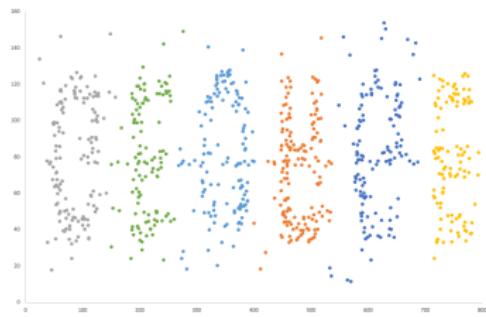


Spiral and group average clustering, K=3.

# Introduction



t5.8k and single-linkage, K=3.

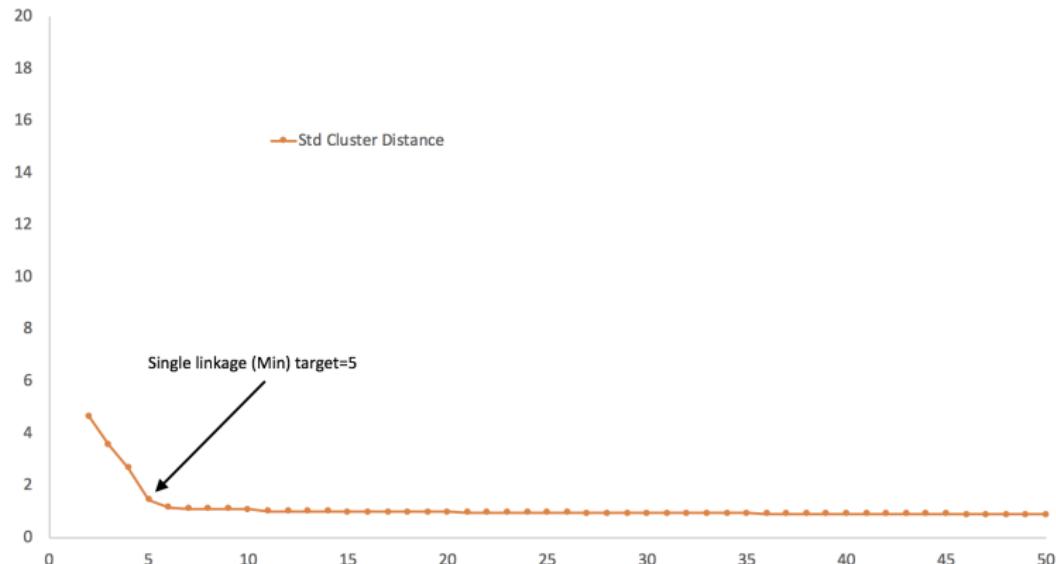


t5.8k and group average, K=6.

# Introduction

- **Knee:** heuristic for determining an optimal clustering
  - ▶ Conventional dendrogram denotes distance measures used during incremental clustering merging.
  - ▶ Typically, the knee is a "bend" in the dendrogram, that visually denotes the optimal clustering.
  - ▶ Knee shows point of *maximal marginal rate of return*. [Zhang et al. 2014]

# Introduction



Aggregation dataset, standard dendrogram, single-linkage clustering  
Successful knee heuristics: max magnitude, max ratio, 2nd derivative

# Introduction

- Issues

- ▶ Different clustering algorithms.
- ▶ Clustering is imperfect.
- ▶ Optimal clustering is ill-defined.
- ▶ Different datasets.
- ▶ Different ways to characterize a knee.
- ▶ Different ways to characterize *distance* in dendrogram.
- ▶ Knees don't always work, because they might not exist.

# Introduction

- Goal: Comparative study...
  - ▶ Knee strategies (9)
  - ▶ 2D spatial datasets (16)
  - ▶ Clustering algorithms (2)
  - ▶ Dendrogram distance measures (3)
  - ▶ ⇒ Total 756 cases.  
(Not all datasets used for one clustering algorithm.)
- How do knee strategies compare, given the above parameters?

# Setup: Hierarchical Clustering

**Table 1.** Hierarchical Clustering Algorithm

---

*Input:*  $S_i = (x_i, y_i), i = 1, \dots, K$

*Output:* Dendogram.

Initialization:

For  $i=1, \dots, K$

    Add new cluster  $C_{\{i\}}$

For all  $C_i, C_j$  ( $i \neq j$ )

$Distance(C_i, C_j) \leftarrow DistanceMeasure(S_i, S_j)$

Cluster generation:

For  $i=1, \dots, K \}$

    Find  $C_p, C_q$  with minimum  $d_{min} \leftarrow Distance(C_p, C_q)$ .

$Dendogram.Dist_i \leftarrow d_{min}$

$Dendogram.Clust_i \leftarrow (C_p, C_q)$

    Remove  $C_p$  and  $C_q$ .

    Add new cluster  $C_{p \cup q}$ .

    Update  $Distance$  table:

        Remove all distances referring to  $C_p$  and  $C_q$ .

        For all clusters  $C_w \neq C_{p \cup q}$

$Distance(C_w, C_{p \cup q}) \leftarrow DistanceMeasure(C_w, C_{p \cup q})$

}

---

Return *Dendogram*.

---

# Setup: Clustering Algorithms

- ① Single-linkage: when clusters  $p$  and  $q$  merged, distance table for other clusters  $w$  revised...

$$\begin{aligned} \text{Distance}(C_w, C_{p \cup q}) = \\ \text{minimum}(\text{Distance}(C_w, C_p), \text{Distance}(C_w, C_q)) \end{aligned}$$

- ② Group average:

$$\begin{aligned} \text{Distance}(C_w, C_{p \cup q}) = \\ \text{average}(\text{Distance}(C_w, C_p), \text{Distance}(C_w, C_q)) \end{aligned}$$

## Setup: Dendrogram Measures

- ① Standard distance (Std): distance used by clustering algorithm.
- ② Global average medoid distance (Avg Med):

$$\text{AvgMed} = \frac{\sum_{i=1}^K MD_i}{T}$$

where  $MD_i$  is avg distance of medoid to other elements in cluster  $i$ ,  
and  $T$  is total # clusters.

- ③ Global average centroid distnace (Avg Cent):

$$\text{AvgCent} = \frac{\sum_{i=1}^K CD_i}{T}$$

where  $CD_i$  is avg distance of centroid to other elements in cluster  $i$ .

# Setup: Knee Strategies

- ① Magnitude: maximum  $d_{i+1} - d_i$ .
- ② Ratio: maximum  $d_{i+1}/d_i$ .
- ③ Second derivative: maximum second derivative.
- ④ Minimum: minimum value.
- ⑤ L-method: [Salvador and Chan 2004] Fit 2 line segments to dendrogram with min RMSE. Node at intersection is knee.
  - ▶ 6. L-method D: If N points on LHS line, then use next N points for RHS line.
  - ▶ 7. L-method S: If N points on LHS line, then evenly sample N points on dendrogram for RHS line.

## Setup: Knee Strategies (cont.)

- F score: Based on F test of one-way ANOVA, applied at each node of dendrogram.

- ⑧ F score A: The highest  $i$  in which:

$$(f_{i+1} - f_i) > \delta_{1..i}^2$$

where  $\delta_{1..i}^2$  is the std dev of F scores 1 to  $i$ .

- ⑨ F score B: The highest  $i$  in which:

$$(f_{i+1} - f_i) > \delta_{1..k}^2$$

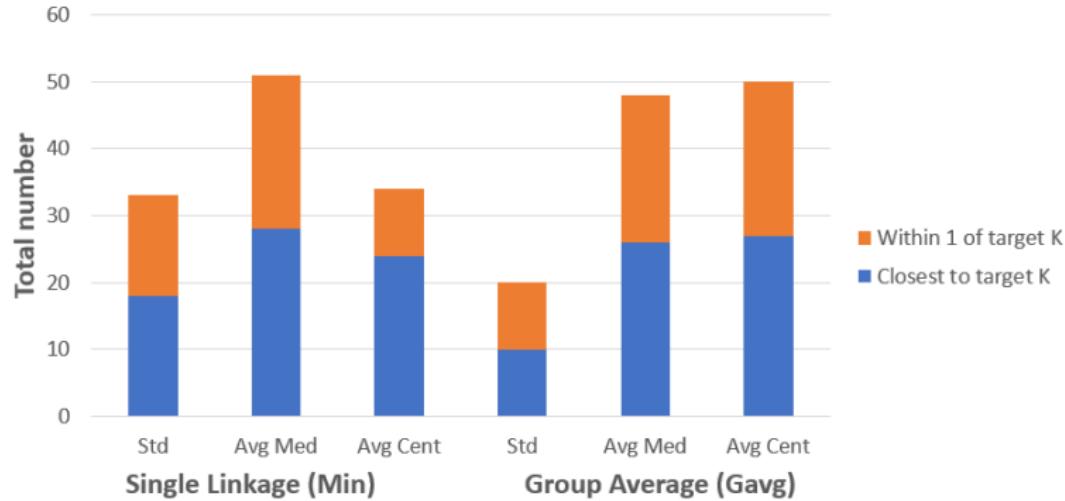
where  $\delta_{1..i}^2$  is the std dev of all F scores on dendrogram.

# Setup: 2D Spatial Datasets

Name	# Nodes		Target Cluster Size		
	Orig.	Reduced	Orig.	Min	Gavg
a1	3000	800	20	10	<b>20</b>
Aggregation	788	788	7	5	<b>7</b>
Birch3	10000	800	100	47	<b>69</b>
Compound	399	399	6	3	<b>6</b>
D31	3100	800	31	19	<b>31</b>
Flame	240	240	2	-	<b>2</b>
Jain	373	373	2	<b>2</b>	<b>2</b>
Pathbased	300	300	3	-	<b>3</b>
R15	600	600	15	11	<b>15</b>
RRR	54	54	3	<b>3</b>	<b>3</b>
Spiral	312	312	3	<b>3</b>	<b>3</b>
t4.8k	8000	800	6	-	<b>6</b>
t5.8k	8000	800	6	3	<b>6</b>
t7.10k	10000	800	9	-	<b>9</b>
t8.8k	8000	800	8	2	<b>8</b>
Unbalance	6500	800	8	<b>7</b>	<b>7</b>

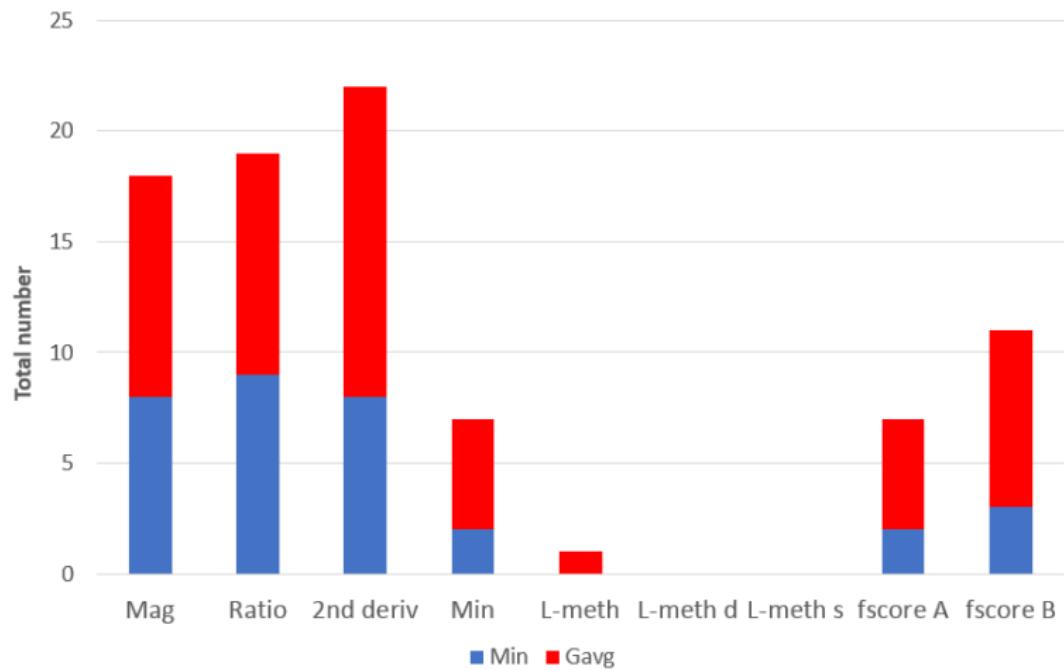
# Results

- Knee performance wrt distance metric



# Results

- Frequency that knee strategies were closest to target cluster size K



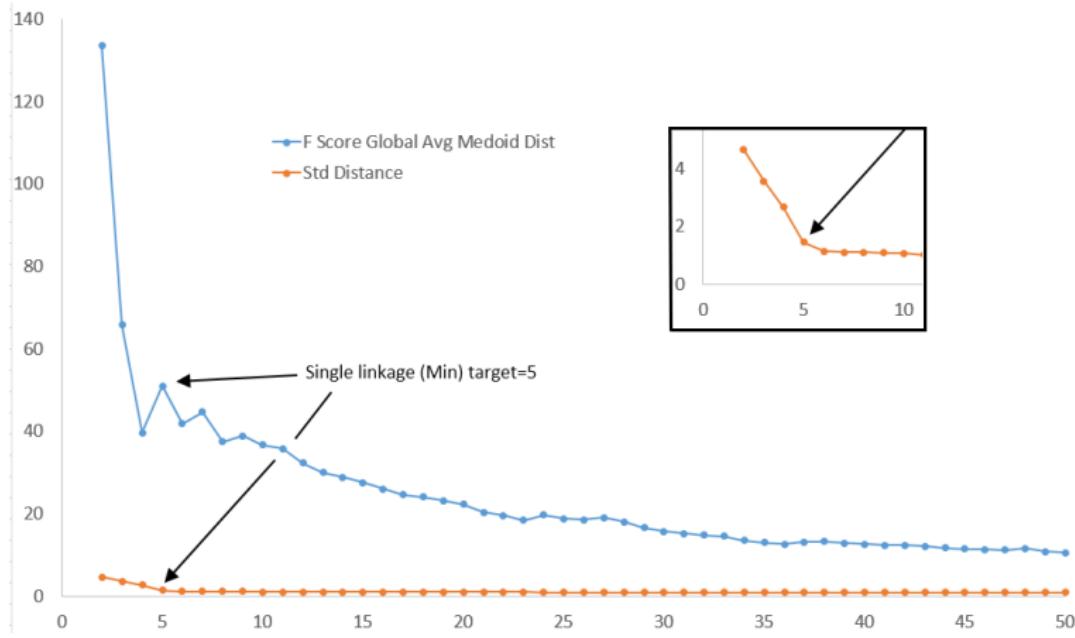
# Results

- Details of knee performance: closest to target

Knee	Std		Avg Med		Avg Cent		Total
	Min	Gavg	Min	Gavg	Min	Gavg	
Mag	5	4	3	1	0	4	18
Ratio	5	2	3	3	1	5	19
2nd deriv	5	4	3	5	0	5	22
Min	0	0	1	2	1	3	7
L-meth	0	0	0	1	0	0	1
L-meth D	0	0	0	0	0	0	0
L-meth S	0	0	0	0	0	0	0
F score A	-	-	1	2	1	3	7
F score B	-	-	2	5	1	3	11

# Results

- Comparing knees using different dendrogram distances  
(Aggregation DS, single linkage clustering)



# Results

- Knee found by F score A and B  
(Aggregation DS, group avg clustering)



- Knee shape is determined by *between group variance* term of ANOVA formula (see tech report).

# Conclusion

- Knee detection is a heuristic. It is not guaranteed to work.
- Many factors for success: data set, clustering algorithm, distance measure, knee strategy.
- **Serendipity.**

Future work:

- Could consider more datasets, clustering algorithms, knee strategies.  
But results will be the same.
- Interesting idea: Use machine learning to discover new knee strategies for different families of datasets.
- Another idea: Use machine learning to identify families of datasets conducive to different clustering optimization strategies.
- Maybe knees are not necessary.