

SS-ZG548: ADVANCED DATA MINING

09

AR Mining On Streams



Dr. Kamlesh Tiwari

Assistant Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

Sept 18, 2021

(WILP @ BITS-Pilani July-Dec 2021)

ONLINE

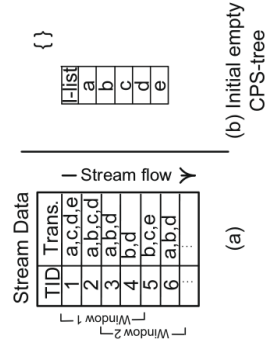
<http://kti.wari.in/adm>

Frequent pattern mining over data streams

- Applications involves retail market data analysis, network monitoring, web usage mining, and stock market prediction.
- Using sliding window
- Efficiently remove the obsolete, old stream data
- Compact Pattern Stream tree (CPS-tree) ¹
- Highly compact frequency-descending tree structure at runtime
- Efficient in terms of memory and time complexity
- Pane and window
- Insertion and restructuring

¹Traber, Syed Khaliuzzaman and Ahmed, Chowdhury Faizan and Jeong, Byeong-Soo and Lee, Young-Koo, "Sliding window-based frequent pattern mining over data streams", in Information sciences, 179(22), pages 3843-3865, Elsevier, 2009

CPS-tree construction



Recap: Data streams

Stream of data arrive in rapid succession, no Re-scan is, storage space is insufficient to accommodate all data points.

Without storing all the data one wish to estimate

- Set of frequent items
- Number of distinct items
- Frequent itemsets
- etc

[Constraints are on memory and processing power]

CPS-tree construction

Algorithm 1 (Construction of a CPS-tree for a data stream)

Input: Stream_data, Pane_size, Window_size, Initial_Sort_Order

Output: T_{curr} is CPS-tree for the current window

Method:

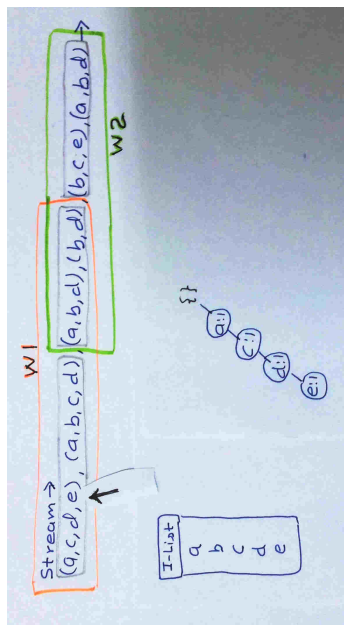
```

1:  $w = w_0$ ;
2:  $T =$  a prefix-tree with null initialization;
3:  $Current\_Sort\_Order \leftarrow Initial\_Sort\_Order$ ;
4: //For each item in the current window
   While (for  $w$  Window_size) do
5:   Call Insert_Pane(T);
6:    $Current\_Sort\_Order \leftarrow$  Frequency-descending sort order; // Reorganizing Phase
7:    $w = w + 1$ ;
8:    $w = w - 1$ ;
9: End While
10: //At each slide of Window
   Begin
11:   Delete the oldest pane information from T;
12:   Call Insert_Pane(T);
13:    $Current\_Sort\_Order \leftarrow$  Frequency-descending sort order; // Reorganizing Phase
14:   End
15: End

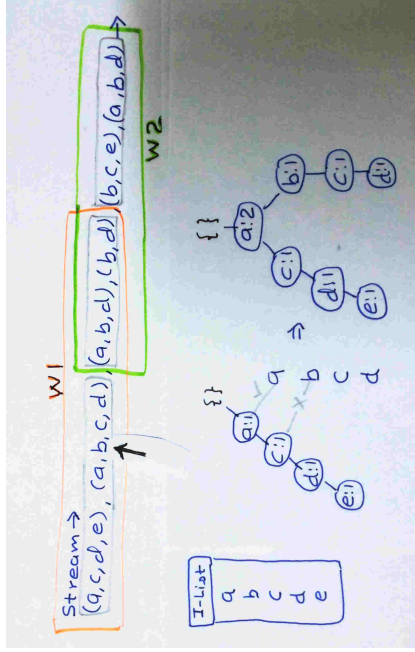
Insert_Pane(T)
Begin
1:  $p = a_i$ ;
2: While (for  $p$  Pane_size) do
3:   Scan transaction from the current location in Stream_data;
4:   Insert the scanned transaction into T according to Current_Sort_Order;
5:    $p = p + 1$ ;
6: End While
End
    
```

Fig. 3. The CPS-tree construction algorithm.

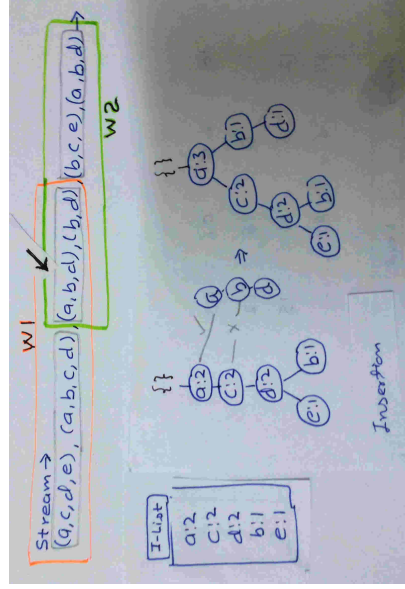
CPS-tree construction



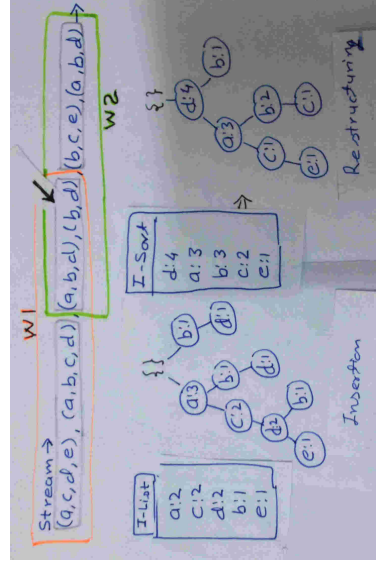
CPS-tree construction



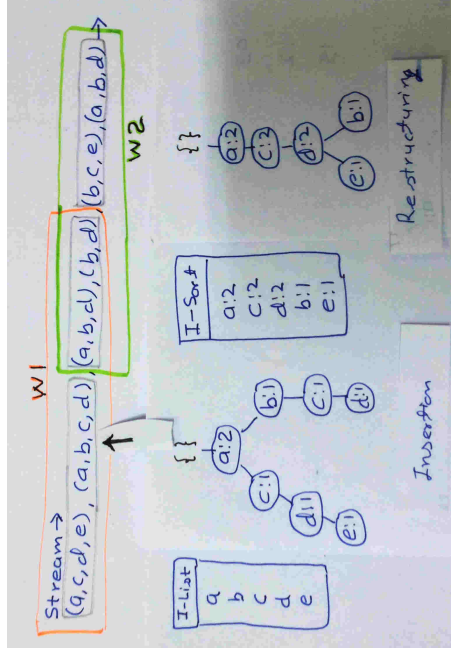
CPS-tree construction



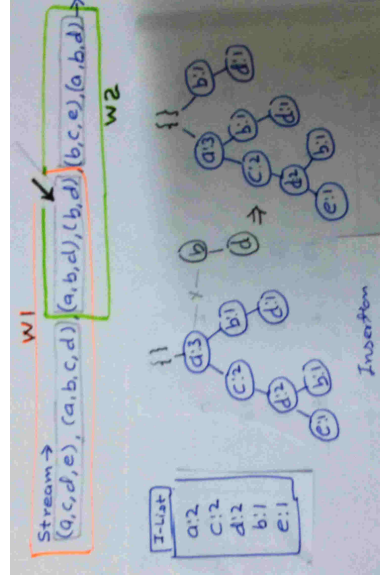
CPS-tree construction



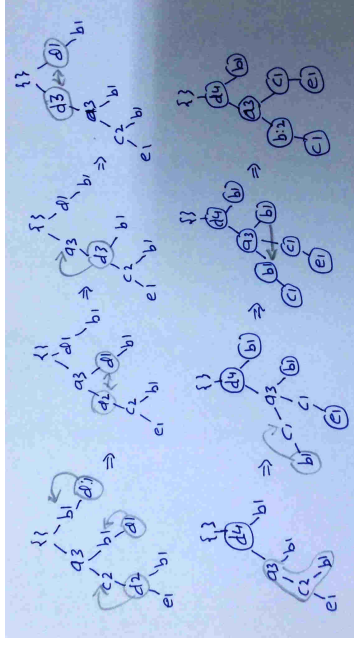
CPS-tree construction



CPS-tree construction

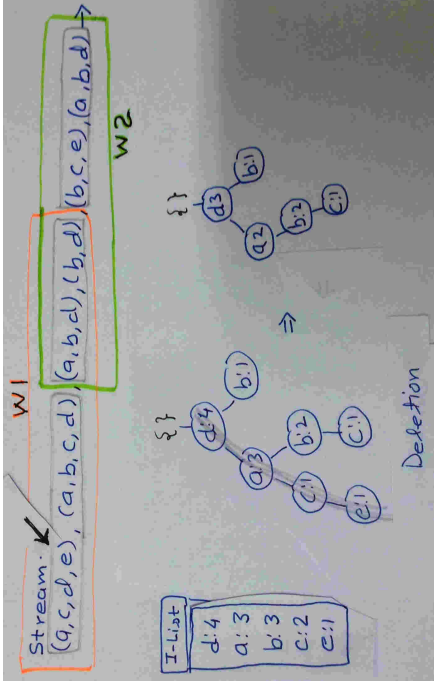


CPS-tree construction

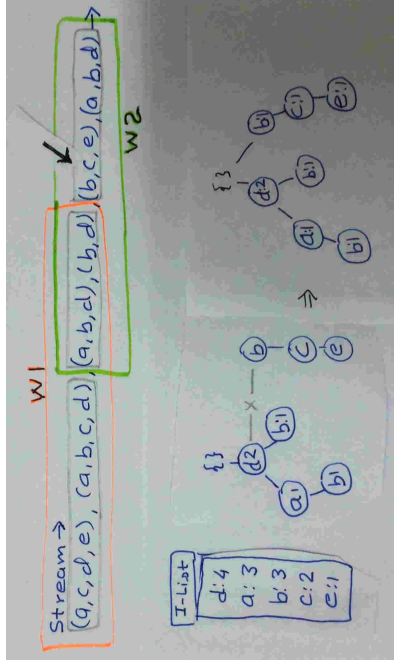


I-Sorted: d, a, b, c, e

CPS-tree construction

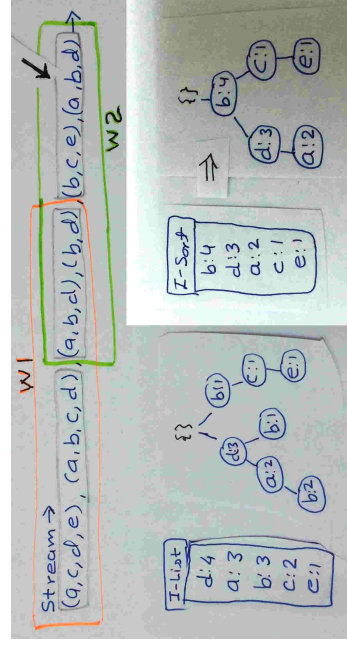


CPS-tree construction



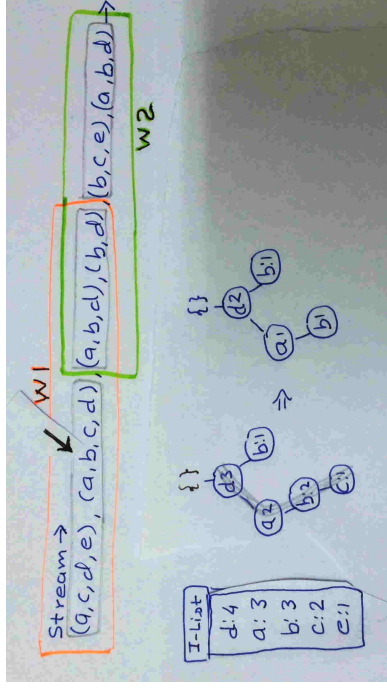
I-Sorted: d, a, b, c, e

CPS-tree construction



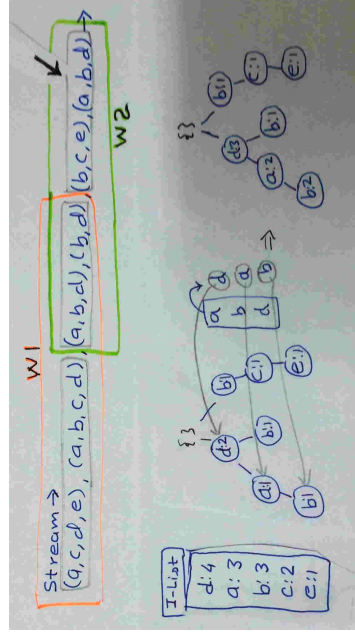
I-Sorted: b, d, a, c, e

CPS-tree construction



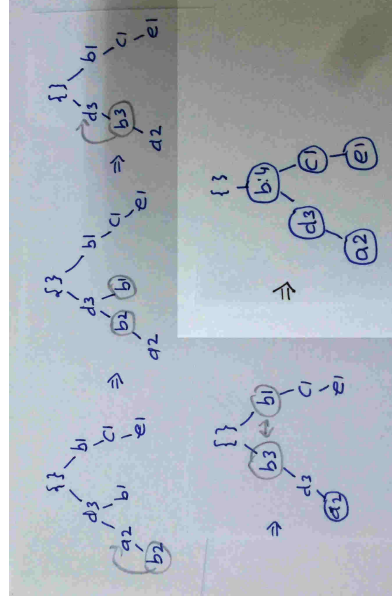
I-Sorted: d, a, b, c, e

CPS-tree construction



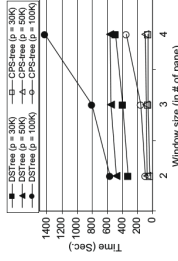
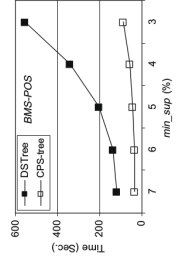
I-Sorted: d, a, b, c, e

CPS-tree construction



I-Sorted: b, d, a, c, e

CPS-tree Performance



Clustering over Evolving Data Stream

- Beside limited memory and one-pass constraints, we require
 - ▶ No assumption on the number of clusters,
 - ▶ Discovery of clusters with arbitrary shape and
 - ▶ Ability to handle outliers
- DenStream³, is a new approach for discovering clusters in an evolving data stream.
- Uses core-micro-cluster to summarize the clusters
- Along with potential core-micro-cluster and outlier micro-cluster structures
- Designed a pruning strategy that guarantees the precision of the weights of the micro-clusters
- User damped window model. Weights with time t are $f(t) = 2^{-\lambda \cdot t}$, (other models landmark and sliding window)

Clustering over Evolving Data Stream

Definition (potential-micro-cluster) A potential-micro-cluster (or p-micro-cluster) at time t for a group of close points p_1, p_2, \dots, p_{i_n} with time stamp $T_{t_1}, T_{t_2}, \dots, T_{t_n}$ is defined as $\{CF^1, CF^2, w\}$, where $w = \sum_{j=1}^n f(t - T_{t_j})$ is the weight, $\beta, 0 < \beta \leq 1$, is the parameter to determine the threshold of outlier relative to c-micro-clusters. $CF^1 = \sum_{j=1}^n f(t - T_{t_j})p_{t_j}$ is the weighted linear sum of the points. $CF^2 = \sum_{j=1}^n f(t - T_{t_j})p_{t_j}^2$ is the weighted squared linear sum of the points. Centre of p-micro-cluster is $c = \frac{CF^1}{w}$ and the radius of a p-micro-cluster $r = \sqrt{\frac{CF^2}{w} + (\frac{CF^1}{w})^2}$, $r \leq \epsilon$

Clustering over Evolving Data Stream (DenStream)

Clustering over Evolving Data Stream

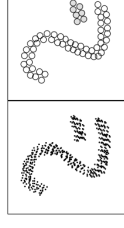
With limited memory and one-pass constraint one want to determine arbitrary number of clusters of arbitrary shape by efficiently handling outliers.

Use DenStream²

- Damped window model: Weights with time t are $f(t) = 2^{-\lambda \cdot t}$, (other models landmark and sliding window)
- core-micro-cluster, potential-micro-cluster and outlier-micro-cluster structures
- Guarantees the precision of the weights of the micro-clusters

Clustering over Evolving Data Stream

- **Definition (core object)** It is an object in whose ϵ neighborhood the overall weight of data points is at least μ
- **Definition (density-area)** A density area is defined as the union of the ϵ neighborhoods of core objects
- **Definition (core-micro-cluster)** At time t it is defined as $CMC(w, c, r)$ for a group of close points p_1, p_2, \dots, p_{i_n} with time stamp $T_{t_1}, T_{t_2}, \dots, T_{t_n}$, $w = \sum_{j=1}^n f(t - T_{t_j})$, $w \geq \mu$ is the weight. $c = \frac{\sum_{j=1}^n f(t - T_{t_j})p_{t_j}}{w}$ is the center, $r = \frac{\sum_{j=1}^n f(t - T_{t_j})dist(p_{t_j}, c)}{w}$, $r \leq \epsilon$ is the radius, where $dist(p_{t_j}, c)$ denotes the Euclidean distance between point p_{t_j} and center c
- When a clustering request arrives, each c-micro-cluster will be labeled to get the final result.



Clustering over Evolving Data Stream

Definition (outlier-micro-cluster) An outlier-micro-cluster (or o-micro-cluster) at time t for a group of close points p_1, p_2, \dots, p_{i_n} with time stamp $T_{t_1}, T_{t_2}, \dots, T_{t_n}$ is defined as $\{CF^1, CF^2, w, t_0\}$. The definition of w, CF^1, CF^2 center and radius are the same as p-micro-cluster. $t_0 = T_{t_i}$, denotes the creation time of o-micro-cluster which is used to define the life span of o-micro-cluster. However $w < \beta \mu$.

Note: p-micro-cluster and o-micro-cluster can be maintained incrementally.

Clustering Algorithm has two parts:

- Online part of micro-cluster maintenance
- Offline part of generation of final clusters, on demand of user

Merging of P

p-micro-clusters and o-micro-clusters are maintained in an online way.

Algorithm 1 Merging (p)

```

1: Try to merge  $p$  into its nearest p-micro-cluster  $c_p$ ;
2: if  $r_p$  (the new radius of  $c_p$ )  $\leq \epsilon$  then
3:   Merge  $p$  into  $c_p$ ;
4: else
5:   Try to merge  $p$  into its nearest o-micro-cluster  $c_o$ ;
6:   if  $r_o$  (the new radius of  $c_o$ )  $\leq \epsilon$  then
7:     Merge  $p$  into  $c_o$ ;
8:     if  $w$  (the new weight of  $c_o$ )  $> \beta\mu$  then
9:       Remove  $c_o$  from outlier-buffer and create a
       new p-micro-cluster by  $c_o$ ;
10:    end if
11:   else
12:     Create a new o-micro-cluster by  $p$  and insert it
       into the outlier-buffer;
13:   end if
14: end if

```

Performance

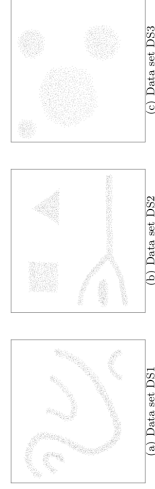


Figure 4. Synthetic data sets

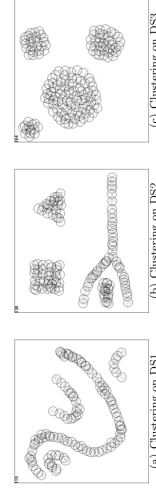


Figure 5: Clustering on DSI, DS2 and DS3

Performance

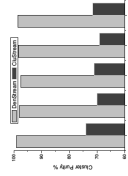


Figure 8: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

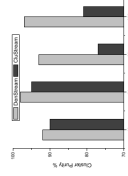


Figure 9: Clustering quality(EDS) data stream, horizon=10, stream speed=1000

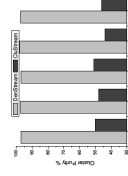


Figure 10: Clustering quality(Network Intrusion data set, horizon=1, stream speed=1000)

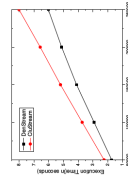


Figure 11: Clustering quality(Network Intrusion data set, horizon=5, stream speed=1000)

DenStream Algorithm

```

Algorithm 2 DenStream ( $DS, \epsilon, \beta, \mu, \lambda$ )
1:  $T_p = \lceil \lambda \log(\frac{\beta\mu}{\beta\mu-1}) \rceil$ ;
2: Get the next point  $p$  at current time  $t$  from data
   stream  $DS$ ;
3: Merging( $p$ );
4: if  $(t \bmod T_p) = 0$  then
5:   for each p-micro-cluster  $c_p$  do
6:     if  $w_p$  (the weight of  $c_p$ )  $< \beta\mu$  then
7:       Delete  $c_p$ ;
8:   end if
9:   end for
10:  for each o-micro-cluster  $c_o$  do
11:     $\xi = \frac{1}{2 \times (1 + \cos(\frac{2\pi}{T_p} - 1))}$ ;
12:    if  $w_o$  (the weight of  $c_o$ )  $< \xi$  then
13:      Delete  $c_o$ ;
14:    end if
15:  end for
16: end if
17: if a clustering request arrives then
18:   Generating clusters;
19: end if

```

```


```

Performance

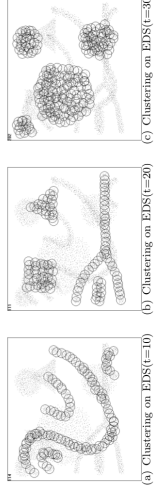
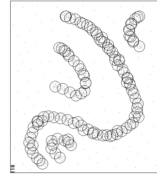
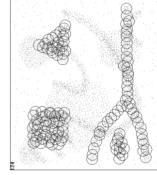


Figure 6: Clustering on the evolving data stream EDS



(a) Clustering on DSI with 5% noise



(b) Clustering on the evolving stream EDS with 5% noise (t=20)

Figure 7: Clustering on data streams with noise

Performance

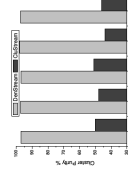


Figure 12: Clustering quality(EDS data stream with 1% noise, horizon=2, stream speed=2000)

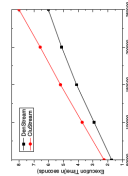


Figure 13: Execution time vs. length of stream(Network Intrusion data set)

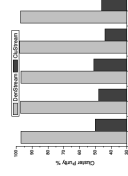


Figure 14: Execution time vs. length of stream(Charitable Donation data set)

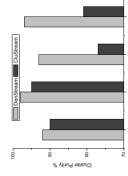


Figure 15: Clustering quality(EDS data stream with 5% noise, horizon=10, stream speed=1000)

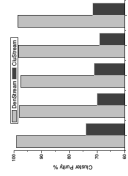


Figure 16: Clustering quality(EDS data stream, horizon=10, stream speed=1000)

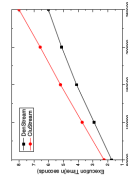


Figure 17: Clustering quality(EDS data stream, horizon=2, stream speed=2000)

Performance

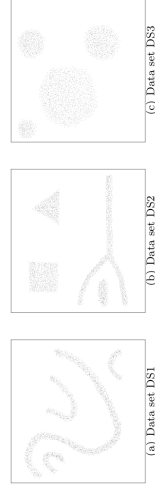


Figure 8: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

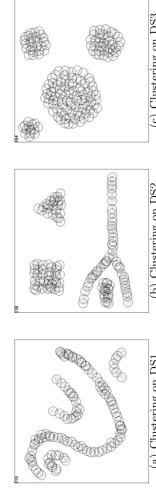


Figure 9: Clustering quality(EDS) data stream, horizon=10, stream speed=1000

Performance

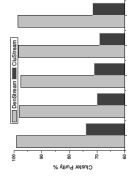


Figure 10: Clustering quality(Network Intrusion data set, horizon=1, stream speed=1000)

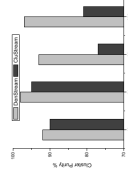


Figure 11: Clustering quality(Network Intrusion data set, horizon=5, stream speed=1000)

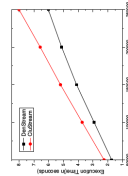


Figure 12: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

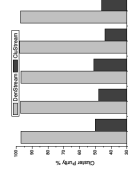


Figure 13: Execution time vs. length of stream(Charitable Donation data set)

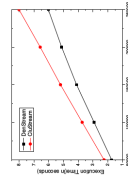


Figure 14: Execution time vs. length of stream(Network Intrusion data set)

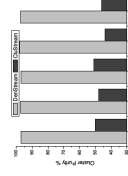


Figure 15: Clustering quality(EDS) data stream with 5% noise, horizon=10, stream speed=1000

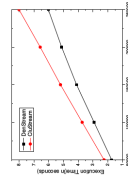


Figure 16: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

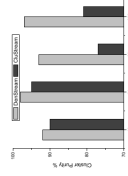


Figure 17: Clustering quality(EDS) data stream, horizon=10, stream speed=1000

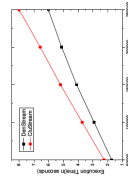


Figure 18: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

Performance

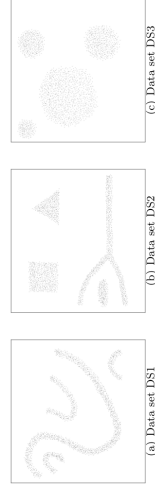


Figure 8: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

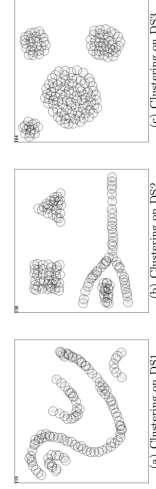


Figure 9: Clustering quality(EDS) data stream, horizon=10, stream speed=1000

Performance

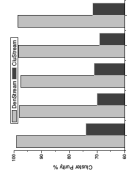


Figure 10: Clustering quality(Network Intrusion data set, horizon=1, stream speed=1000)

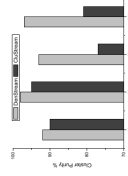


Figure 11: Clustering quality(Network Intrusion data set, horizon=5, stream speed=1000)

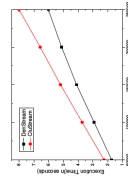


Figure 12: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

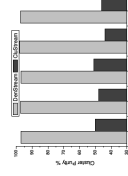


Figure 13: Execution time vs. length of stream(Charitable Donation data set)

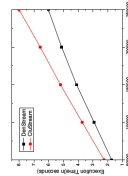


Figure 14: Execution time vs. length of stream(Network Intrusion data set)

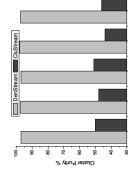


Figure 15: Clustering quality(EDS) data stream with 5% noise, horizon=10, stream speed=1000

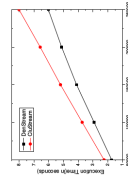


Figure 16: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

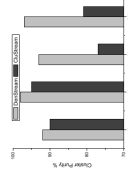


Figure 17: Clustering quality(EDS) data stream, horizon=10, stream speed=1000

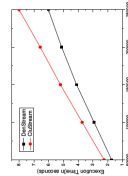


Figure 18: Clustering quality(EDS) data stream, horizon=2, stream speed=2000

Performance

Thank You!

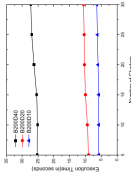


Figure 16: Execution time vs. dimensionality

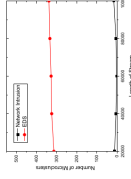


Figure 18: Memory usage

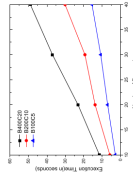


Figure 17: Execution time vs. number of clusters

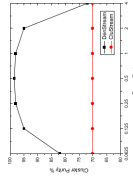


Figure 19: Clustering quality vs. decay factor λ

Thank you very much for your attention!
Queries ?