

Scaling Geometric Monitoring Over Distributed Streams

Alexandros D. Keros

June 23, 2016

Supervised by: Prof. V.Samoladas

Table of contents

Introduction

Theoretical Background

- The Geometric Monitoring Method

- Theoretical Tools

- Related Work

Problem Statement & Implementation

- Problem Statement

- Implementation

Experimental Results

- Data & Setup

- Experiments

Conclusions & Future Work

- Conclusion

- Future Work

Introduction

Theoretical Background

- The Geometric Monitoring Method
- Theoretical Tools
- Related Work

Problem Statement & Implementation

- Problem Statement
- Implementation

Experimental Results

- Data & Setup
- Experiments

Conclusions & Future Work

- Conclusion
- Future Work

Data Stream Systems

- ▶ **Data streams:** Continuous, high volume, size unbound, violative, probably distributed
- ▶ *Pull paradigm*
- ▶ Centralizing and/or polling → prohibitive in terms of communication overhead
- ▶ Examples: telecommunication, sensor networks

The Geometric Monitoring Method

- ▶ Threshold monitoring
- ▶ Nodes communicate when needed
 - ▶ Local constraints
 - ▶ Violation resolution (*false alarms*)
- ▶ Arbitrary function monitoring
- ▶ Tight accuracy bounds
- ▶ A promising framework for *distributed data stream monitoring*

Motivation

Problems:

- ▶ increasing node population
- ▶ data volume
- ▶ data dimensionality
- ▶ arbitrary functions
- ▶ **communication - accuracy tradeoff**

Need for:

- ▶ scalability warranties
- ▶ tight accuracy bounds
- ▶ incremental/real-time operation
- ▶ **Minimize communication while retaining accuracy bounds**

Motivation

Problems:

- ▶ increasing node population
- ▶ data volume
- ▶ data dimensionality
- ▶ arbitrary functions
- ▶ **communication - accuracy tradeoff**

Need for:

- ▶ scalability warranties
- ▶ tight accuracy bounds
- ▶ incremental/real-time operation
- ▶ Minimize communication while retaining accuracy bounds

Motivation

Problems:

- ▶ increasing node population
- ▶ data volume
- ▶ data dimensionality
- ▶ arbitrary functions
- ▶ **communication - accuracy tradeoff**

Need for:

- ▶ scalability warranties
- ▶ tight accuracy bounds
- ▶ incremental/real-time operation
- ▶ **Minimize communication while retaining accuracy bounds**

Motivation

Problems:

- ▶ increasing node population
- ▶ data volume
- ▶ data dimensionality
- ▶ arbitrary functions
- ▶ **communication - accuracy tradeoff**

Need for:

- ▶ scalability warranties
- ▶ tight accuracy bounds
- ▶ incremental/real-time operation
- ▶ **Minimize communication while retaining accuracy bounds**

Contributions

Expand the *geometric monitoring method*:

- ▶ heuristic method for violation resolution
- ▶ distance-based hierarchical node clustering
- ▶ throughout method evaluation on synthetic and real-world datasets

Contributions

Expand the *geometric monitoring method*:

- ▶ heuristic method for violation resolution
- ▶ distance-based hierarchical node clustering
- ▶ throughout method evaluation on synthetic and real-world datasets

Contributions

Expand the *geometric monitoring method*:

- ▶ heuristic method for violation resolution
- ▶ distance-based hierarchical node clustering
- ▶ throughout method evaluation on synthetic and real-world datasets

Introduction

Theoretical Background

- The Geometric Monitoring Method
- Theoretical Tools
- Related Work

Problem Statement & Implementation

- Problem Statement
- Implementation

Experimental Results

- Data & Setup
- Experiments

Conclusions & Future Work

- Conclusion
- Future Work

Geometric Threshold Monitoring

- ▶ **Threshold monitoring:** arbitrary function $f(\cdot)$, threshold T

$$f(\cdot) < T \text{ or } f(\cdot) > T$$

- ▶ **Idea:** decompose into local constraints at the nodes

System Architecture

Decentralized Scenario

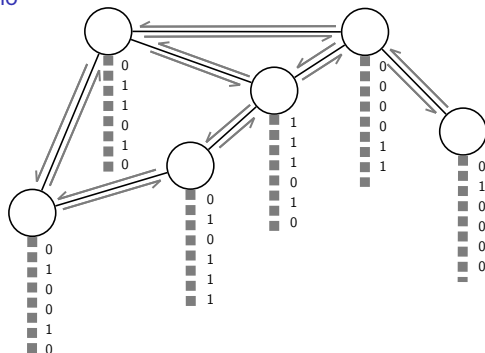


Figure: Mesh-like network topology example of the decentralized scenario. Dashed lines represent data streams and half arrows represent message exchanges.

System Architecture

Centralized Scenario

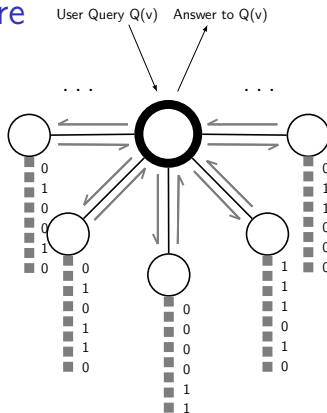


Figure: Star-like network topology example of the centralized scenario. The bold node represents the coordinator node. Dashed lines represent data streams and half arrows represent message exchanges.

Computational Model

Statistics vectors

- ▶ the *monitoring function* $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ the *threshold* $T \in \mathbb{R}$
- ▶ the *monitoring node set* : $P = \{p_1, \dots, p_n\}$
with *weights* w_1, \dots, w_n
- ▶ the *data streams* : $S = \{s_1, \dots, s_n\}$
- ▶ the *d-dimensional local statistics vectors* : $\vec{v}_1(t), \dots, \vec{v}_n(t)$
represent each node's data stream at time t

Global statistics vector

$$\vec{v}(t) = \frac{\sum_{i=1}^n w_i \vec{v}_i(t)}{\sum_{i=1}^n w_i} \quad (1)$$

Computational Model

Estimate vector

Infrequent communication between nodes/nodes-coordinator:

Estimate vector

$$\vec{e}(t) = \frac{\sum_{i=1}^n w_i \vec{v}_i'}{\sum_{i=1}^n w_i} \quad (2)$$

- ▶ the last communicated *local statistics vector* of node p_i : \vec{v}_i'
- ▶ *Local statistics divergence*: $\Delta \vec{v}_i(t) = \vec{v}_i(t) - \vec{v}_i', i = 1, \dots, n$

Decentralized drift vector

$$\vec{u}_i(t) = \vec{e}(t) + \Delta \vec{v}_i(t) \quad (3)$$

Centralized drift vector

$$\vec{u}_i(t) = \vec{e}(t) + \Delta \vec{v}_i(t) + \frac{\vec{\delta}_i}{w_i} \quad (4)$$

Computational Model

Balancing Process

Centralized scenario

Purpose: resolve possible false alarms

Balancing vector

$$\vec{b} = \frac{\sum_{p_i \in P'} w_i \vec{u}_i(t)}{\sum_{p_i \in P'} w_i} \quad (5)$$

- ▶ the *balancing set* P' : a subset of nodes
- ▶ the *slack vector* at the nodes $\vec{\delta}_i = \vec{\delta}_i' + \Delta \vec{\delta}_i$, $\sum_{p_i \in P'} \Delta \vec{\delta}_i = \vec{0}$:

$$\Delta \vec{\delta}_i = w_i \vec{b} - w_i \vec{u}_i(t) \quad \forall p_i \in P' \quad (6)$$

, readjusts the *drift vectors* (4).

Geometric Interpretation

Convexity Property

Convexity Property

$$\vec{v}(t) = \frac{\sum_{i=1}^n w_i \vec{u}_i(t)}{\sum_{i=1}^n w_i} \quad (7)$$

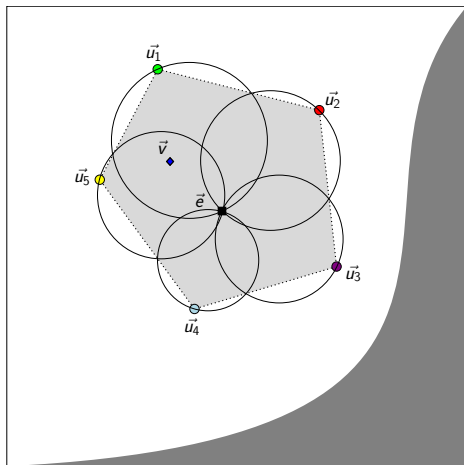
Theorem (Sharfman et al.)

Let $\vec{x}, \vec{y}_1, \dots, \vec{y}_n \in \mathbb{R}^d$ be a set of vectors in \mathbb{R}^d . Let $\text{Conv}(\vec{x}, \vec{y}_1, \dots, \vec{y}_n)$ be the convex hull of $\vec{x}, \vec{y}_1, \dots, \vec{y}_n$. Let $B(\vec{x}, \vec{y}_i)$ be a ball centered at $\frac{\vec{x} + \vec{y}_i}{2}$ and with radius of $\|\frac{\vec{x} + \vec{y}_i}{2}\|_2$ i.e., $B(\vec{x}, \vec{y}_i) = \{\vec{z} \mid \|\vec{z} - \frac{\vec{x} + \vec{y}_i}{2}\|_2 \leq \|\frac{\vec{x} + \vec{y}_i}{2}\|_2\}$, then $\text{Conv}(\vec{x}, \vec{y}_1, \dots, \vec{y}_n) \subset B(\vec{x}, \vec{y}_i)$.

Geometric Interpretation

Convexity Property & Local Constraints

Figure: Example of a convex hull (light gray) defined by the drift vectors $\vec{u}_i, i = 1, 2, 3, 4, 5$. The hull is bounded by the spheres created from the estimate vector \vec{e} and the drift vectors $\vec{u}_i, i = 1, 2, 3, 4, 5$. The global statistics vector \vec{v} is guaranteed to be contained in the convex hull of the drift vectors.



Protocol

Decentralized Algorithm

Algorithm 1: Decentralized algorithm

```
1 begin
2   foreach node  $p_i$  do                               /* Node initialization */
3       Broadcast  $\vec{v}_i(0)$ ;
4        $\vec{v}_i' = \vec{v}_i(0)$ ;
5       Wait messages from all other nodes;
6       if messages from all vectors received then
7           Compute estimate vector  $\vec{e}(t)$ ;
8       end
9   end
10  foreach node  $p_i$  do                                /* Main monitoring task */
11      foreach new  $s_i$  stream update  $\vec{v}_i(t)$  do
12          Recalculate drift vector  $\vec{u}_i(t)$ ;
13          if  $B(\vec{e}, \vec{u}_i(t))$  is not monochromatic then
14              Broadcast message  $\langle i, \vec{v}_i(t) \rangle$ ;
15              Set  $\vec{v}_i' = \vec{v}_i(t)$ ;
16          end
17          if new message  $\langle j, \vec{v}_j(t) \rangle$  received then
18              Set  $\vec{v}_j' = \vec{v}_j(t)$ ;
19              Recalculate estimate vector  $\vec{e}(t)$ ;
20              if  $B(\vec{e}, \vec{u}_i(t))$  is not monochromatic then
21                  Broadcast message  $\langle i, \vec{v}_i(t) \rangle$ ;
22                  Set  $\vec{v}_i' = \vec{v}_i(t)$ ;
23              end
24          end
25      end
26  end
27 end
```

Protocol

Centralized Algorithm

Algorithm 2: Centralized algorithm's coordinator node operation

```
1 begin
2   Wait for < INIT, · > messages from all monitoring nodes;
   /* Initialization */
3   Compute estimate vector  $\vec{e}(0)$ ;
4   if new < REP,  $\vec{v}_i(t)$ ,  $\vec{u}_i(t)$  > message received then
   /* Monitoring operation */
5      $P' = P' \cup \{ < i, \vec{v}_i(t), \vec{u}_i(t) > \}$ ;
6     Balance( $P'$ );
7   end
8 end
9 Function Balance( $P'$ ) /* Balancing Process */
10  Compute balancing vector  $\vec{b}$ ;
11  if  $B(\vec{e}, \vec{b})$  is not monochromatic then
12    if  $P - P' \neq \emptyset$  then
13      Send < REQ > message to random node in  $P - P'$  set;
14    else
15      Compute estimate vector  $\vec{e}(t)$ ;
16      Send < NEW-EST,  $\vec{e}(t)$  > message to all nodes;
17      return;
18    end
19  else
20    foreach  $p_i \in P'$  do
21      Compute slack adjustment vector  $\Delta \vec{\delta}_i$ ;
22      Send < ADJ-SLK,  $\Delta \vec{\delta}_i$  > message to node  $p_i$ ;
23      return;
24    end
25  end
26 end
```

Algorithm 3: Centralized algorithm's monitoring node operation

```
1 begin
2   foreach node  $p_i$  do /* Node initialization */
3     Send < INIT,  $\vec{v}_i(0)$  > message to coordinator;
4      $\vec{v}_i' = \vec{v}_i(0)$ ;
5      $\vec{\delta}_i = \vec{0}$ ;
6     Wait message from coordinator;
7     if < NEW-EST,  $\vec{e}$  > message received then
8       Set  $\vec{e}(t) = \vec{e}$ ;
9     end
10  end
11  end
12  foreach node  $p_i$  do /* Main monitoring task */
13    foreach new  $s_i$  stream update  $\vec{v}_i(t)$  do
14      Recalculate drift vector  $\vec{u}_i(t)$ ;
15      if  $B(\vec{e}, \vec{u}_i(t))$  is not monochromatic then
16        Send < REP,  $\vec{v}_i(t)$ ,  $\vec{u}_i(t)$  > message to coordinator;
17        Wait for < NEW-EST, · > or < ADJ-SLK, · >
18        message from coordinator;
19      end
20      if new message < REQ > received then
21        Send < REP,  $\vec{v}_i(t)$ ,  $\vec{u}_i(t)$  > message to coordinator;
22        Wait for < NEW-EST, · > or < ADJ-SLK, · >
23        message from coordinator;
24      end
25      if new < NEW-EST,  $\vec{e}$  > message received then
26        Set  $\vec{e}(t) = \vec{e}$ ;
27         $\vec{v}_i' = \vec{v}_i(t)$ ;
28         $\vec{\delta}_i = \vec{0}$ ;
29      end
30      if new < ADJ-SLK,  $\Delta \vec{\delta}_i$  > message received then
31        Recompute delta vector  $\vec{\delta}_i$ ;
32      end
33    end
34  end
35 end
```

Multi-objective Optimization

- ▶ **Multiple**, possibly **conflicting** objectives to be *simultaneously* optimized
- ▶ *Pareto optimality*(*non-dominated solutions*): optimal solutions where none of the objective functions can be optimized without the simultaneous degradation of other objective functions' values.
- ▶ Let vector of m objectives $F(x) = [F_1(x), F_2(x), \dots, F_m(x)]$:

$$\min_{x \in \mathbb{R}^n} F(x)$$

$$\text{s.t. } l \leq x \leq u$$

$$G_i = 0, i = 1, \dots, k_e$$

$$G_j \leq 0, j = k_e + 1, \dots, k$$

- ▶ Finding *Pareto optimal* solutions is generally **NP-hard**.

Non-linear Constraint Optimization

Primal Descent

Algorithm 4: Generic primal descent

```

1 begin
2   Choose initial point  $x_0 \in X$  and set  $t = 0$  ;           /* Initialization */
3   while maximum iteration limit OR convergence do       /* Search */
4      $t = t + 1$ ;
5     Determine search direction  $d_t$ ;
6     Determine step length  $s_t$ , so that  $f(x_t + s_t d_t) < f(x_t)$ ;
7     Update;
8   end
9 end

```

Feasible Directions

Usable feasible direction d_t :

- ▶ a small disposition towards direction d_t does not violate any constraint i.e.,

$$d_t^T \nabla G(x_t) \leq 0$$

- ▶ a move towards d_t reduces the objective functions value i.e.,

$$d_t^T \nabla F(x_t) < 0$$

.

SQP

- ▶ The *Lagrangian function*: $\mathcal{L}(x, \lambda) = F(x) + \sum_{i=1}^k \lambda_i G_i(x)$
- ▶ *Quadratic programming subproblems*:

$$\min_{d \in \mathbb{R}^n} \frac{1}{2} d^T H_t d + \nabla F(x_t)^T d$$

$$\nabla G_i(x_t)^T d + G_i(x_t) = 0, i = 1, \dots, k_e$$

$$\nabla G_i(x_t)^T d + G_i(x_t) \leq 0, i = k_e + 1, \dots, k$$

,where:

H_t : Hessian of the Lagrangian function at iteration t

d : search direction

The Savitzky-Golay Filter

- ▶ Low-pass smoothing filter
- ▶ *Moving window averaging* paradigm: $g_i = \sum_{n=-n_L}^{n_R} c_n f_{i+n}$
- ▶ *Least-squares fit* of polynomial $y_i(x)$ over window $n_L + n_R + 1$:

$$\sum_{j=i-n_L}^{i+n_R} (y_i(x_j) - f_j)^2 = \min$$

,where:

$$y_i(x) = a_0 + a_1 \frac{x - x_i}{\Delta x} + a_2 \left(\frac{x - x_i}{\Delta x} \right)^2 + \dots + a_M \left(\frac{x - x_i}{\Delta x} \right)^M$$

- ▶ Set g_i to the value of the fitted point x_i

Maximum Weight Matching

Let $G = (V, E)$ a graph:

- ▶ *maximum weight matching* $M \subseteq E$: a subset of edges where
 - ▶ no two edges share a common vertex
 - ▶ largest possible number of edges
 - ▶ maximizes the sum of weights

Maximum Weight Matching

The Primal-Dual Method

Primal-Dual Method

Constraints in Primal \iff Variables in Dual

Constraints in Dual \iff Variables in Primal

The primal:

$$\begin{aligned} \max \quad & \sum_{(u,v) \in E} x_{u,v} w_{u,v} \\ \text{s.t.} \quad & x_{u,v} \geq 0, (u,v) \in E \\ & \sum_{u \in e: e \in E} x_e \leq 1, u \in V \end{aligned}$$

The dual:

$$\begin{aligned} \min \quad & \sum_{u \in V} y_u \\ \text{s.t.} \quad & y_u \geq 0, u \in V \\ & y_u + y_v \geq w_{u,v}, (u,v) \in E \end{aligned}$$

Related Work

- ▶ *Safe Zones*: optimal local constraints fitted to nodes' data distributions
- ▶ Ellipsoidal bounding regions, decouplement of estimate vector from bounding ball construction
- ▶ Simple shapes as local constraints, hierarchical clustering of nodes for participation to the balancing operation
- ▶ Prediction models based on velocity and acceleration

Introduction

Theoretical Background

The Geometric Monitoring Method
Theoretical Tools
Related Work

Problem Statement & Implementation

Problem Statement
Implementation

Experimental Results

Data & Setup
Experiments

Conclusions & Future Work

Conclusion
Future Work

Problem Formulation

Reduce the communication burden of the *Geometric Monitoring* method by:

- ▶ Optimally position *drift vectors* during the *balancing process*
- ▶ Appropriate node selection for inclusion in the *balancing set*

,in order to **increase scalability** in terms of:

- ▶ node population
- ▶ stream dimensionality.

The Geometric Monitoring Framework

Assumptions

- ▶ Coordinator-based scenarion
- ▶ Instantaneous, loss-less, reliable communication
- ▶ Iterative operation based on time-steps
- ▶ System pause during violation resolution
- ▶ Coordinator node does not monitor a stream

Algorithm 5: Iterative network operation

Data: *monitoringNodes*: a list of Monitoring nodes,
coordinator: the Coordinator node

```

1 begin
2   initialization;
3   repeat
4     foreach node ∈ monitoringNodes do
5       node.DataVectorUpdate();
6       node.ComputeDriftVector();
7     end
8     foreach node ∈ monitoringNodes do
9       node.CheckForViolation();
10      if localViolation then
11        node.Report();
12        coordinator.Balance();
13      end
14    end
15  until globalViolation;
16 end
  
```

The Distance-based Hierarchical Clustering

The Idea

Node selection for *violation resolution*:

- ▶ Elevate *randomness* of the initial *geometric monitoring* method
- ▶ Hierarchical node clustering scheme
- ▶ Decouple matching from the data distribution at the nodes
- ▶ Accurately follow *global statistics vector*

The Distance-based Hierarchical Clustering

The Weight Function

Weight function

$$w_{i,j} = \sum_{t=t_0}^{t_{end}} [(f(\vec{v}_{global}(t)) - f(\frac{\vec{v}_i(t) + \vec{v}_j(t)}{2})) + (|\vec{v}_i(t) - \vec{v}_j(t)|)]$$

The Distance-based Hierarchical Clustering

The Algorithm

Algorithm 6: Recursively create Monitoring node pairs and hierarchy

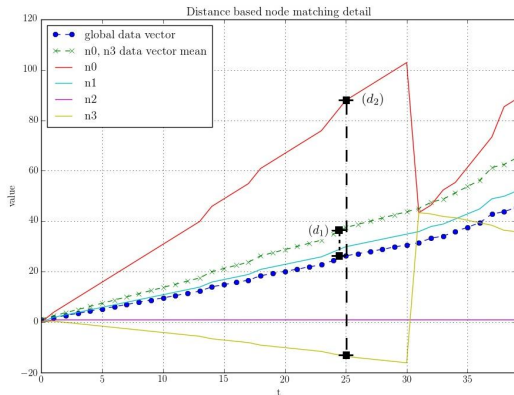
```

1 Function DistancePairer(nodes, i)
    Data: nodes =  $[(n_1, [\vec{v}_1(t_0), \dots, \vec{v}_1(t_{end})]), \dots, (n_k, [\vec{v}_k(t_0), \dots, \vec{v}_k(t_{end})])]$ : list of
        nodes with their respective data vectors, i: pair type, initial=1
    Result: nodeHierarchy: dictionary of Type-k pairs
2 if length(nodes) = 1 then                                // recursion stopping condition
3     return nodeHierarchy;
4 end
5 g = CreateCompleteGraph(nodes);                          // complete graph with nodes as
    vertices
6 foreach  $(n_i, n_j) \in g.Edges()$  do                        // assign weights to edges
7      $w_{i,j} = \sum_{t=t_0}^{t_{end}} [(f(\vec{v}_{global}(t)) - f(\frac{\vec{v}_i(t) + \vec{v}_j(t)}{2})) + (|\vec{v}_i(t) - \vec{v}_j(t)|)]$ ;
8     g.edge(ni, nj).weight = wi,j;
9 end
10 nodeHierarchy(Type-i) = g.maximalWeightMatching();      // node pairs of
    Type-i
11 DistancePairer(nodeHierarchy(Type-i), i * 2);
12 end
  
```

The Distance-based Hierarchical Clustering

Example

Figure: Distance based node matching operating on 4 nodes ($\{n_0, n_1, n_2, n_3\}$). Distance d_1 : the distance of the data vector mean of the paired nodes n_0 and n_3 from the global mean (*global data vector*), distance d_2 : denotes the in-between distance of data vectors $\vec{v}_0(t)$ and $\vec{v}_3(t)$ of the node pair. Both distances are taking part in the edge weighting process.



The Heuristic Balancing

The Idea

The Heuristic Balancing

The Optimizing Function

The Heuristic Balancing

The Function Formulation

The Heuristic Balancing

The Algorithm

An Nested Optimization Problem

Implementation Challenges

Synthetic Data

Real-world Data

Notation

RAND, DIST, DISTR Comparison

GM, HM Comparison

GM, HDM Comparison

Synthetic Data Monitoring

GM, HDM Comparison

Air Pollution Monitoring

Summary & Concluding Remarks

Future Work

The end
Questions?

Appendix

Savitzky-Golay filter matrix notation

hey you