

# 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

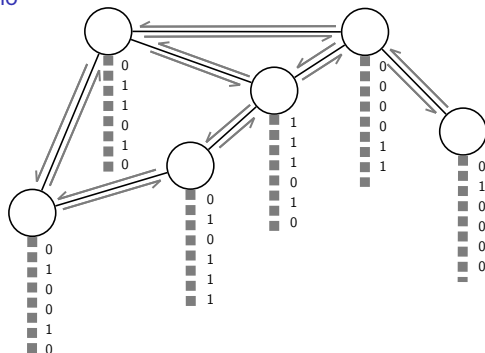
- ▶ **Threshold monitoring:** arbitrary function  $f(\cdot)$ , threshold  $T$

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

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

# System Architecture

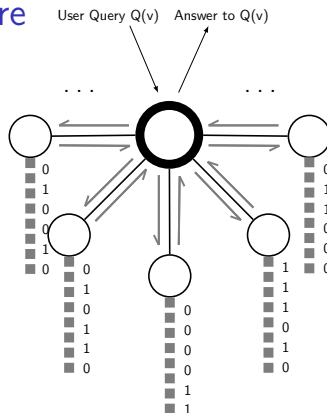
## Decentralized Scenario



**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

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

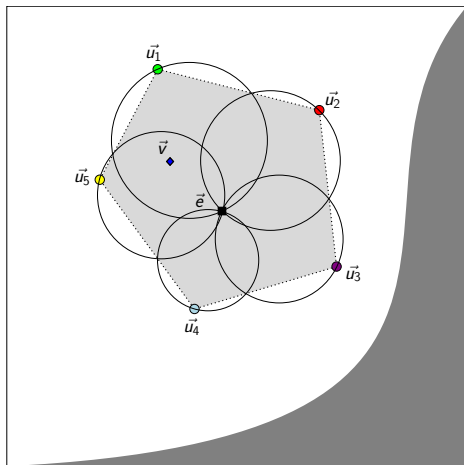
### Theorem (Sharfman et al.)

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    end
21    if new message < REQ > received then
22      Send < REP,  $\vec{v}_i(t)$ ,  $\vec{u}_i(t)$  > message to coordinator;
23      Wait for < NEW-EST, · > or < ADJ-SLK, · >
24      message from coordinator;
25    end
26    if new < NEW-EST,  $\vec{e}$  > message received then
27      Set  $\vec{e}(t) = \vec{e}$ ;
28       $\vec{v}_i' = \vec{v}_i(t)$ ;
29       $\vec{\delta}_i = \vec{0}$ ;
30    end
31    if new < ADJ-SLK,  $\Delta \vec{\delta}_i$  > message received then
32      Recompute delta vector  $\vec{\delta}_i$ ;
33    end
34  end
35 end
```

# Multi-objective Optimization



# Non-linear Constraint Optimization

## Primal Descent

# Feasible Directions

# SQP

# The Savitzky-Golay Filter

# Maximum Weight Matching

## The Primal-Dual Method

# Related Work

# Problem Formulation

# The Geometric Monitoring Framework



# The Distance-based Hierarchical Clustering

## The Idea

## The Weight Function

# The Distance-based Hierarchical Clustering

## The Algorithm

# 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



# Velocity and Acceleration Estimation via SG Filtering



# 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?