# Scaling Geometric Monitoring Over Distributed Streams

Alexandros D. Keros

June 23, 2016

Supervised by: Prof. V.Samoladas

## Table of contents

**Introduction** Theoretical Background Problem Statement & Implementation Experimental Results Conclusions & Future Work

Overview

# Data Stream Systems

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

**Introduction** Theoretical Background Problem Statement & Implementation Experimental Results Conclusions & Future Work
○●○○ ○○○○○○○○○○○ ○ ○○ ○
 ○○○○○○ ○○○○○○○○○○○ ○○○○○ ○○
 ○

Overview

# 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*

**Introduction** Theoretical Background Problem Statement & Implementation Experimental Results Conclusions & Future Work
0000 00000000000 0 00 0
000000 00000000000 00000 00
0

Overview

# 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

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          00000000000             0                                   00                     0
              000000                  00000000000                         00000                  00
              0

Overview

# 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

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
○○●○         ○○○○○○○○○○○                 ○                                ○○                     ○
             ○○○○○○                      ○○○○○○○○○○○                        ○○○○○                  ○○
             ○

Overview

# 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

Introduction Theoretical Background Problem Statement & Implementation Experimental Results Conclusions & Future Work
0000 00000000000 0 00 0
       000000 00000000000 00000 00
       0

Overview

# 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**

**Introduction** Theoretical Background Problem Statement & Implementation Experimental Results Conclusions & Future Work
○○○● ○○○○○○○○○○○ ○ ○○ ○
 ○○○○○○ ○○○○○○○○○○○ ○○○○○ ○○
 ○

Overview

## 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  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
000●          00000000000                 0                                 00                    0
              000000                      00000000000                       00000                 00
              0

Overview

## 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  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
○○○●         ○○○○○○○○○○○              ○                                 ○○                   ○
             ○○○○○○                  ○○○○○○○○○○○                        ○○○○○                ○○
             ○

Overview

## 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**  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          ●000000000                    0                                 00                    0
              000000                        00000000000                       00000                 00
              0

The Geometric Monitoring Method

# 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

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          0●000000000             0                                   00                    0
              000000                  00000000000                        00000                 00
              0

The Geometric Monitoring Method

# 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.

Alexandros D. Keros

Scaling Geometric Monitoring Over Distributed Streams

Introduction  **Theoretical Background**  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          00●00000000                 0                                  00                    0
              000000                       00000000000                       00000                 00
              0

The Geometric Monitoring Method

# 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.

Introduction   Theoretical Background   Problem Statement & Implementation   Experimental Results   Conclusions & Future Work
0000   000●0000000   0   00   0
       000000   00000000000   00000   00
       0

The Geometric Monitoring Method

# Computational Model

Statistics vectors

- the *monitoring function* $f : \mathbb{R}^d \to \mathbb{R}$
- the *threshold* $T \in \mathbb{R}$
- the *monitoring node set* : $P = \{p_1, \ldots, p_n\}$
  with *weights* $w_1, \ldots, w_n$
- the *data streams* : $S = \{s_1, \ldots, s_n\}$
- the $d$-dimensional *local statistics vectors* : $\vec{v_1}(t), \ldots, \vec{v_n}(t)$
  represent each node's data stream at time $t$

### Global statistics vector

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

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          0000●000000              O                                   00                    O
              000000                   00000000000                         00000                 00
              O

The Geometric Monitoring Method

# Computational Model

### Estimate vector

Infrequent communication between nodes/nodes-coordinator:

---

**Estimate vector**

$$\vec{e}(t) = \frac{\sum_{i=1}^{n} w_i \vec{v_i}'}{\sum_{i=1}^{n} w_i} \tag{2}$$

---

- 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**

$$\vec{u_i}(t) = \vec{e}(t) + \Delta\vec{v_i}(t) \tag{3}$$

**Centralized drift vector**

$$\vec{u_i}(t) = \vec{e}(t) + \Delta\vec{v_i}(t) + \frac{\vec{\delta_i}}{w_i} \tag{4}$$

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          00000●00000                ○                                 ○○                   ○
              000000                     00000000000                       00000                ○○
              ○

The Geometric Monitoring Method

# Computational Model

Balancing Process

**Centralized scenario**

**Purpose**: resolve possible false alarms

### Balancing vector

$$\vec{b} = \frac{\sum_{p_i \in P'} w_i \vec{u_i}(t)}{\sum_{p_i \in P'} w_i} \tag{5}$$

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

$$\Delta\vec{\delta_i} = w_i\vec{b} - w_i\vec{u_i}(t) \ \forall \ p_i \in P' \tag{6}$$

, readjusts the *drift vectors* (4).

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          0000000●0000            0                                 00                  0
              000000                  00000000000                                           00
              0

The Geometric Monitoring Method

# Geometric Interpretation

Convexity Property

### Convexity Property

$$\vec{v}(t) = \frac{\sum_{i=1}^{n} w_i \vec{u_i}(t)}{\sum_{i=1}^{n} w_i} \tag{7}$$

### Theorem (Sharfman et al. [3])

*Let $\vec{x}, \vec{y_1}, \ldots, \vec{y_n} \in \mathbb{R}^d$ be a set of vectors in $\mathbb{R}^d$. Let $Conv(\vec{x}, \vec{y_1}, \ldots, \vec{y_n})$ be the convex hull of $\vec{x}, \vec{y_1}, \ldots, \vec{y_n}$. 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., $B(\vec{x}, \vec{y_i}) = \{\vec{z} \mid \|\vec{z} - \frac{\vec{x}+\vec{y_i}}{2}\|_2 \le \|\frac{\vec{x}+\vec{y_i}}{2}\|_2\}$, then $Conv(vecx, \vec{y_1}, \ldots, \vec{y_n}) \subset B(\vec{x}, \vec{y_i})$.*

Introduction   Theoretical Background   Problem Statement & Implementation   Experimental Results   Conclusions & Future Work
0000           0000000●000            O                                     OO                    O
               000000                 00000000000                           00000                 OO
               O

The Geometric Monitoring Method

# Geometric Interpretation

Convexity Property

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.

Introduction    Theoretical Background    Problem Statement & Implementation    Experimental Results    Conclusions & Future Work
0000            00000000●00                 0                                    00                      0
                000000                      00000000000                          00000                   00
                0

The Geometric Monitoring Method

# Geometric Interpretation

Local Constraints

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000          000000000●0                 0                                  00                    0
              000000                      00000000000                        00000                00
              0

The Geometric Monitoring Method

# Protocol
Decentralized Algorithm

Introduction    Theoretical Background    Problem Statement & Implementation    Experimental Results    Conclusions & Future Work
0000            0000000000●              ○                                    ○○                      ○
                000000                   00000000000                          00000                   ○○
                ○

The Geometric Monitoring Method

# Protocol

## Centralized Algorithm

Theoretical Tools

# Multi-objective Optimization

Theoretical Tools

# Non-linear Constraint Optimization

Primal Descent

Introduction    Theoretical Background    Problem Statement & Implementation    Experimental Results    Conclusions & Future Work
0000            00000000000               0                                     00                      0
                000●000                   00000000000                           00000                   00
                0

Theoretical Tools

# Feasible Directions

Theoretical Tools

# SQP

Theoretical Tools

# The Savitzky-Golay Filter

# Maximum Weight Matching

The Primal-Dual Method

Introduction   Theoretical Background   Problem Statement & Implementation   Experimental Results   Conclusions & Future Work
0000           00000000000              0                                     00                     0
               000000                   00000000000                           00000                  00
               ●

Related Work

# Related Work

Introduction  Theoretical Background  **Problem Statement & Implementation**  Experimental Results  Conclusions & Future Work
0000          00000000000            ●                                     00                    0
              000000                 00000000000                           00000                 00
              0

Problem Statement

# Problem Formulation

Introduction    Theoretical Background    Problem Statement & Implementation    Experimental Results    Conclusions & Future Work
0000            00000000000               ○                                     ○○                      ○
                000000                    ●000000000                            00000                   ○○
                ○

Implementation

# The Geometric Monitoring Framework

Introduction  Theoretical Background  **Problem Statement & Implementation**  Experimental Results  Conclusions & Future Work
0000          00000000000              0                                        00                     0
              000000                   0●000000000                              00000                  00
              0

Implementation

# The Distance-based Hierarchical Clustering

The Idea

Introduction  Theoretical Background  **Problem Statement & Implementation**  Experimental Results  Conclusions & Future Work
0000          00000000000                 0                                     00                    0
              000000                      00●00000000                           00000                 00
              0

Implementation

# The Distance-based Hierarchical Clustering

The Weight Function

Implementation

# The Distance-based Hierarchical Clustering

The Algorithm

Introduction    Theoretical Background    Problem Statement & Implementation    Experimental Results    Conclusions & Future Work
0000             00000000000                 0                                    00                       0
                 000000                      00000●000000                         00000                    00
                 0

Implementation

# The Heuristic Balancing

## The Idea

Introduction    Theoretical Background    **Problem Statement & Implementation**    Experimental Results    Conclusions & Future Work
0000    00000000000    ○    00    ○
         000000         00000●00000    00000    00
         ○

Implementation

# The Heuristic Balancing

The Optimizing Function

Implementation

# The Heuristic Balancing

The Function Formulation

Implementation

# The Heuristic Balancing

The Algorithm

Introduction
0000

Theoretical Background
00000000000
000000
0

Problem Statement & Implementation
0
00000000●00

Experimental Results
00
00000

Conclusions & Future Work
0
00

Implementation

# An Nested Optimization Problem

Introduction   Theoretical Background   **Problem Statement & Implementation**   Experimental Results   Conclusions & Future Work
0000          00000000000              0000000000●0                                  00                     0
              000000                   ○                                             00000                  00
              ○

Implementation

# Velocity and Acceleration Estimation via SG Filtering

Introduction    Theoretical Background    **Problem Statement & Implementation**    Experimental Results    Conclusions & Future Work
0000            00000000000               O                                        OO                     O
                000000                    0000000000●                              00000                  OO
                O

Implementation

## Implementation Challenges

Data & Setup

# Synthetic Data

Introduction    Theoretical Background    Problem Statement & Implementation    **Experimental Results**    Conclusions & Future Work
0000            00000000000               0                                     0●                          0
                000000                    00000000000                           00000                       00
                0

Data & Setup

# Real-world Data

Introduction    Theoretical Background    Problem Statement & Implementation    **Experimental Results**    Conclusions & Future Work
0000            00000000000               O                                    OO                         O
                000000                    00000000000                          ●0000                      OO
                O

Experiments

# Notation

Introduction     Theoretical Background     Problem Statement & Implementation     **Experimental Results**     Conclusions & Future Work
0000            00000000000                 ○                                        ○○                  ○
                000000                      00000000000                              0●000               ○○
                ○

Experiments

# RAND, DIST, DISTR Comparison

Introduction    Theoretical Background    Problem Statement & Implementation    **Experimental Results**    Conclusions & Future Work
0000           00000000000                ○                                     ○○                        ○
               000000                     00000000000                           00●00                     ○○
               ○

Experiments

# GM, HM Comparison

Introduction    Theoretical Background    Problem Statement & Implementation    **Experimental Results**    Conclusions & Future Work
0000            00000000000               ○                                    ○○                        ○
                000000                    00000000000                           000●0                     ○○
                ○

Experiments

# GM, HDM Comparison

Synthetic Data Monitoring

Introduction Theoretical Background Problem Statement & Implementation **Experimental Results** Conclusions & Future Work
0000          00000000000                 000000000000                          00                    0
              000000                       0                                     0000●                  00
              0

Experiments

# GM, HDM Comparison

Air Pollution Monitoring

Introduction  Theoretical Background  Problem Statement & Implementation  Experimental Results  Conclusions & Future Work
0000         00000000000            0                                00                     ●
             000000                 00000000000                      00000                  00
             0

Conclusion

# Summary & Concluding Remarks

Introduction   Theoretical Background   Problem Statement & Implementation   Experimental Results   Conclusions & Future Work
0000           00000000000              0                                    00                     0
               000000                   00000000000                          00000                  ●0
               0

Future Work

# Future Work

The end
Questions?