

# Distributed Pattern Discovery in Multiple Streams

Jimeng Sun<sup>†</sup> Spiros Papadimitriou<sup>§</sup> Christos Faloutsos\*<sup>†</sup>

<sup>†</sup>Carnegie Mellon University    <sup>§</sup>IBM Watson Research Center  
{jimeng,christos}@cs.cmu.edu    spapadim@us.ibm.com

**Abstract.** Given  $m$  groups of streams which consist of  $n_1, \dots, n_m$  co-evolving streams in each group, we want to: (i) incrementally find local patterns within a single group, (ii) efficiently obtain global patterns across groups, and more importantly, (iii) efficiently do that in real time while limiting shared information across groups. In this paper, we present a distributed, hierarchical algorithm addressing these problems. Our experimental case study confirms that the proposed method can perform hierarchical correlation detection efficiently and effectively.

## 1 Introduction

Streams are often inherently correlated and it is possible to reduce hundreds of numerical streams into just a handful of *patterns* that compactly describe the key trends and dramatically reduce the complexity of further data processing. Multiple co-evolving streams often arise in a large distributed system, such as computer networks and sensor networks. Centralized approaches usually will not work in this setting. The reasons are: **(i) Communication constraint**; it is too expensive to transfer all data to a central node for processing and mining. **(ii) Power consumption**; in a wireless sensor network, minimizing information exchange is crucial because many sensors have very limited power. **(iii) Robustness concerns**; centralized approaches always suffer from single point of failure. **(iv) Privacy concerns**; in any network connecting multiple autonomous systems (e.g., multiple companies forming a collaborative network), no system is willing to share all the information, while they all want to know the global patterns. To sum up, a **distributed online algorithm** is highly needed to address all the above concerns.

To address this problem, we propose a hierarchical framework that intuitively works as follows: 1) Each autonomous system first finds its local patterns and shares them with other groups (details in section 4). 2) Global patterns are discovered based on the shared local patterns (details in section 5). 3) From the global patterns, each autonomous system further refines/verifies their local patterns.

\* This material is based upon work supported by the National Science Foundation under Grants No. IIS-0209107 IIS-0205224 INT-0318547 SENSOR-0329549 IIS-0326322 This work is also supported in part by the Pennsylvania Infrastructure Technology Alliance (PITA), a partnership of Carnegie Mellon, Lehigh University and the Commonwealth of Pennsylvania's Department of Community and Economic Development (DCED). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation, or other funding parties.

## 2 Problem formalization

Given  $m$  groups of streams which consist of  $\{n_1, \dots, n_m\}$  co-evolving numeric streams, respectively, we want to solve the following two problems: (i) incrementally find patterns within a single group (*local pattern monitoring*), and (ii) efficiently obtain global patterns from all the local patterns (*global pattern detection*).

More specifically, *local pattern monitoring* can be modelled as a function,

$$F_L : (S_i(t+1, :), G(t, :)) \rightarrow L_i(t+1, :), \quad (1)$$

where the inputs are 1) the new input point  $S_i(t+1, :)$  at time  $t+1$  and the current global pattern  $G(t, :)$  and the output is the local pattern  $L_i(t+1, :)$  at time  $t+1$ . Details on constructing such a function will be explained in section 4. Likewise, *global pattern detection* is modelled as another function,

$$F_G : (L_1(t+1, :), \dots, L_m(t+1, :)) \rightarrow G(t+1, :), \quad (2)$$

where the inputs are local patterns  $L_i(t+1, :)$  from all groups at time  $t+1$  and the output is the new global pattern  $G(t+1, :)$ .

## 3 Distributed mining framework

In this section, we introduce the general framework for distributed mining. More specifically, we present the meta-algorithm to show the overall flow, using  $F_L$  (*local patterns monitoring*) and  $F_G$  (*global patterns detection*) as black boxes.

Intuitively, it is natural that global patterns are computed based on all local patterns from  $m$  groups. On the other hand, it might be a surprise that the local patterns of group  $i$  take as input both the stream measurements of group  $i$  and the global patterns. Stream measurements are a natural set of inputs, since local patterns are their summary. However, we also need global patterns as another input so that local patterns can be represented consistently across all groups. This is important at the next stage, when constructing global patterns out of the local patterns; we elaborate on this later. The meta-algorithm is the following:

Algorithm DISTRIBUTEDMINING

---

0. (*Initialization*) At  $t = 0$ , set  $G(t, :) \leftarrow \text{null}$

1. For all  $t > 1$

(*Update local patterns*) For  $i \leftarrow 1$  to  $m$ , set  $L_i(t, :) := F_L(S_i(t, :), G(t-1, :))$

(*update global patterns*) Set  $G(t, :) := F_G(L_1, \dots, L_m)$

---

## 4 Local pattern monitoring

In this section we present the method for discovering patterns within a stream group. More specifically, we explain the details of function  $F_L$  (Equation 1). We first describe the intuition behind the algorithm and then present the algorithm formally. Finally we discuss how to determine the number of local patterns  $k_i$ .

The goal of  $F_L$  is to find the low dimensional projection  $L_i(t, :)$  and the participation weights  $W_{i,t}$  so as to guarantee that the reconstruction error  $\|S_i(t, :) - \hat{S}_i(t, :)\|^2$  over time is predictably small.

**Tracking local patterns:** The first step is, for a given  $k_i$ , to incrementally update the  $k \times n_i$  participation weight matrix  $W_{i,t}$ , which serves as a basis of the low-dimensional projection for  $S_i(t, :)$ . Later in this section, we describe the method for choosing  $k_i$ . For the moment, assume that the number of patterns  $k_i$  is given.

The main idea behind the algorithm is to read the new values  $S_i(t+1, :) \equiv [S_i(t+1, 1), \dots, S_i(t+1, n_i)]$  from the  $n_i$  streams of group  $i$  at time  $t+1$ , and perform three steps: **(1) Compute** the low dimensional projection  $y_j, 1 \leq j \leq k_i$ , based on the *current* weights  $W_{i,t}$ , by projecting  $S_i(t+1, :)$  onto these. **(2) Estimate** the reconstruction error ( $e_j$  below) and the energy. **(3) Compute**  $W_{i,t+1}$  and output the *actual* local pattern  $L_i(t+1, :)$ .

The term  $\lambda$  is an exponential forgetting factor between 0 and 1, which helps adapt to more recent behavior. For instance,  $\lambda = 1$  means putting equal weights on all historical data, while smaller  $\lambda$  means putting higher weight on more recent data. We chose 0.96 throughout all experiments. As long as the values of  $S_i$  do not vary wildly, the exact value of  $\lambda$  is not crucial.

Algorithm  $F_L$

---

**Input:** new vector  $S_i(t+1, :)$ , old global patterns  $G(t, :)$

**Output:** local patterns ( $k_i$ -dimensional projection)  $L_i(t+1, :)$

1. Initialize  $\mathbf{x}_1 := S_i(t+1, :)$ .

2. For  $1 \leq j \leq k$ , we perform the following in order:

$$y_j := \mathbf{x}_j W_{i,t}(j, :)^T \quad (y_j = \text{projection onto } W_{i,t}(j, :))$$

If  $G(t, :) = \text{null}$ , then  $G(t, j) := y_j$  (handling boundary case)

$$d_j \leftarrow \lambda d_j + y_j^2 \quad (\text{local energy, determining update magnitude})$$

$$\mathbf{e} := \mathbf{x}_j - G(t, j) W_{i,t}(j, :)^T \quad (\text{error, } \mathbf{e} \perp W_{i,t}(j, :))$$

$W_{i,t+1}(j, :) \leftarrow W_{i,t}(j, :) + \frac{1}{d_j} G(t, j) \mathbf{e}$  (update participation weight)

$$\mathbf{x}_{j+1} := \mathbf{x}_j - G(t, j) W_{i,t+1}(j, :)^T \quad (\text{repeat with remainder of } \mathbf{x}).$$

3. Compute the new projection  $L_i(t+1, :) := S_i(t+1, :) W_{i,t+1}^T$

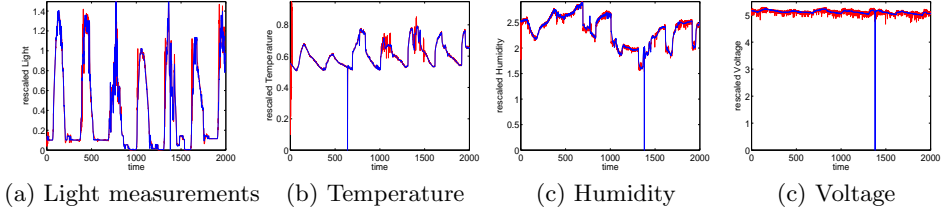
---

**Detecting the number of local patterns:** In practice, we do not know the number  $k_i$  of local patterns. We propose to estimate  $k_i$  on the fly, so that we maintain a high percentage  $f_{i,E}$  of the *energy*  $E_{i,t}$ . For each group, we have a low-energy and a high-energy threshold,  $f_{i,E}$  and  $F_{i,E}$ , respectively. We keep enough local patterns  $k_i$ , so the retained energy is within the range  $[f_{i,E} \cdot E_{i,t}, F_{i,E} \cdot E_{i,t}]$ .

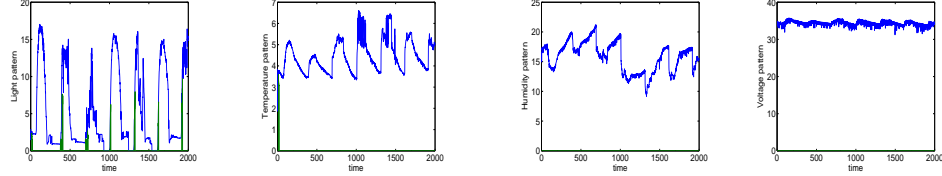
## 5 Global pattern detection

In this section we present the method for obtaining global patterns over all groups. More specifically, we explain the details of function  $F_G$  (Equation 2).

First of all, what is a global pattern? Similar to local pattern, global pattern is low dimensional projections of the streams from all groups. Loosely speaking, assume only one global group exists which consists of all streams, the global



**Fig. 1** original measurements (blue) and reconstruction (red) are very close



**Fig. 2.** Local patterns

patterns are the local patterns obtained by applying  $F_L$  on the global group—this is essentially the centralized approach. In other words, we want to obtain the result of the centralized approach without centralized computation.

The algorithm exactly follows the lemma above. The  $j$ -th global pattern is the sum of all the  $j$ -th local patterns from  $m$  groups.

Algorithm  $F_G$

**Input:** all local patterns  $L_1(t, :), \dots, L_m(t, :)$

**Output:** global patterns  $G(t, :)$

0. Set  $k := \max(k_i)$  for  $1 \leq i \leq m$

1. For  $1 \leq j \leq k$ , set  $G(t, j) := \sum_{i=1}^m L_i(t, j)$  (if  $j > k_i$  then  $L_i(t, j) \equiv 0$ )

## 6 Experimental case study

The **Motes** dataset consists of 4 groups of sensor measurements (i.e., light intensity, humidity, temperature, battery voltages) collected using 48 Berkeley Mote sensors at different locations in a lab, over a period of a month.

The main characteristics (see the blue curves in Figure 1) are: (1) Light measurements exhibit a clear global periodic pattern (daily cycle) with occasional big spikes from some sensors (outliers), (2) Temperature shows a weak daily cycle and a lot of bursts. (3) Humidity does not have any regular pattern. (4) Voltage is almost flat with a small downward trend.

The reconstruction is very good (see the red curves in Figure 1(a)), with relative error below 6%. Furthermore, the local patterns from different groups are correlated well with the original measurements (see Figure 2). The global patterns (in Figure 3) are combinations of different patterns from all groups and reveal the overall behavior of all the groups.

In terms of accuracy, everything boils down to the quality of the summary provided by the local/global patterns. To this end, we use the relative reconstruction error as the evaluation metric. The best performance is obtained when accurate global patterns are known to all groups. But this requires exchanging up-to-date local/global patterns at every timestamp among all groups, which is

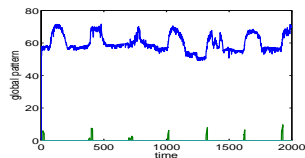


Fig. 3. Global patterns

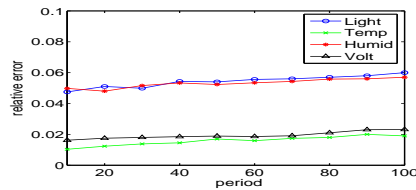


Fig. 4. Error increases slowly

prohibitively expensive. One efficient way to deal with this problem is to increase the communication period, which is the number of timestamps between successive local/global pattern transmissions. Overall, the relative error rate increases very slowly as the communication period increases. This implies that we can dramatically reduce communication with minimal sacrifice of accuracy.

## 7 Related work

**Distributed data mining** Most of works on distributed data mining focus on extending classic (centralized) data mining algorithms into distributed environment, such as association rules mining [4], frequent item sets [6]. Web is a popular distributed environment. Several techniques are proposed specifically for that, for example, distributed top-k query [2] and Bayes-net mining on web [3]. But our focus are on finding numeric patterns, which is different.

**Privacy preserving data mining** The most related discussion here is on how much privacy can be protected using subspace projection method [1,5]. Liu et al. [5] discuss the subspace projection method and propose a possible method to breach the protection using Independent component analysis(ICA). All the method provides a good insight on the issues on privacy protection. Our method focuses more on incremental online computation of subspace projection.

## 8 Conclusion

We focus on finding patterns in a large number of distributed streams. More specifically, we first find local patterns within each group, where the number of local patterns is automatically determined based on reconstruction error. Next, global patterns are identified, based on the local patterns from all groups. We evaluated our method on several datasets, where it indeed discovered the patterns. We gain significant communication savings, with small accuracy loss.

## References

1. C. Agrawal and P. Yu. A condensation approach to privacy preserving data mining. In *EDBT*, 2004.
2. B. Babcock and C. Olston. Distributed Top-K Monitoring. In *SIGMOD*, 2003.
3. R. Chen, S. Krishnamoorthy, and H. Kargupta. Distributed Web Mining using Bayesian Networks from Multiple Data Streams. In *ICDM*, pages 281–288, 2001.
4. D. W. Cheung, V. T. Ng, A. W. Fu, and Y. Fu. Efficient Mining of Association Rules in Distributed Databases. *TKDE*, 8:911–922, 1996.
5. K. Liu, H. Kargupta, and J. Ryan. Multiplicative noise, random projection, and privacy preserving data mining from distributed multi-party data. In *TKDE*, 2005.
6. K. K. Loo, I. Tong, B. Kao, and D. Cheung. Online Algorithms for Mining Inter-Stream Associations From Large Sensor Networks. In *PAKDD*, 2005.