# Synthesizing Scenario-based Dataset for User Behavior Pattern Mining

**University of Ontario INSTITUTE OF TECHNOLOGY**

Weina Ma
Kamran Sartipi

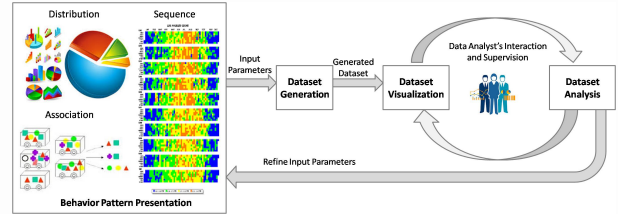{Weina.Ma, Kamran.Sartipi}@uoit.ca

## Motivations

- Identifying user behavior patterns from audit logs is valuable for system security of monitoring authorized users.
- Due to restricted access to production event-logs, security and privacy issues, and high costs of real datasets, synthetic event-log datasets are crucial in designing and evaluating data analytics approaches.
- A controlled event-log simulation environment provides the data analysts various synthetic dataset containing embedded interesting patterns and features. The produced testing datasets reduce the algorithm evaluation time.

## Behavior Pattern

```
{
    {
        "id": "P-00001",
        "description": "typical doctor's workflow in radiology department",
        "context": {
            "role": "doctor",
            "department": "radiology"
        },
        "sequence": {
            "action": "create an order",
            "action": "read historical exams",
            "action": "create an exam",
            "action": "create a report"
        },
        "interval": "within 3 days",
        "support": "10%"
    }
    {
        "id": "P-00002",
        "description": "daily nursing ward-round",
        "context": {
            "role": "nurse"
        },
        "sequence": {
            "location": "ward-A",
            "location": "ward-B",
            "location": "ward-C"
        },
        "interval": "daily",
        "support": "15%"
    }
}
```

## Architecture



We developed an interactive data exploration environment to such a design-generate-visualize-analyze-optimize process.

- *Design*: statistical characteristics (distribution), association pattern, sequence pattern
- *Generate*: produce a dataset that contains predefined attributes and patterns
- *Visualize*: extract simplified workable information from generated dataset
- *Analyze*: verify the differences between generated dataset and expected dataset
- *Optimize*: refine input parameters

## Proposed Approach

We proposed a synthetic event-log generator that effectively assists data analysts in designing scenario-driven event-logs with embedded user behavior patterns, and visually analyzing the quality of the generated datasets. The toolkit includes three layers:

- *Behavior pattern representation layer*: allows data analysts to design interesting features and patterns that will be injected into the dataset.
- *Dataset generation layer*: creates datasets that are controlled by data size, data distribution, and the designed behavior patterns.
- *Dataset visualization layer and analysis layer*: provides an interactive exploration environment for visual analysis of the quality of generated datasets.

## Generator Algorithm

**Algorithm:** dataset-generator
**Input:** $A$, $V$, $B$, $avg$, $D_{All}$, $U_{All}$
**Output:** $E$

```
1   I = Ø, S = Ø , E = Ø , U = Ø
2   for B_j in B do begin
3       association pattern i_j = B_j-c
4       sequence pattern s_j=B_j-s
5       apply time constraint B_j-t to s_j
6       randomly select B_j-sup users U-B_j
7   end for
8   for a_i in A do begin
9       build biased random value select function F_k
10  end for
11  for each user u_u in U_All do begin
12      for each day d_p in D_All do begin
13          x = randomly selected an integer around avg
14          generate an empty event sequence u_u-d_p-s_x
15          for U-B_j in U do begin
16              if u_u exists in U-B_j then
17                  insert constraint-based sequence pattern s_j to
                    u_u-d_p-s_x
18                  insert association pattern i_j to u_u-d_p-s_x
19              end if
20          end for
21          for each event e_i in u_u-d_p-s_x do begin
22              for each empty attribute a_i in e_i do begin
23                  call function F_k to assign random value v_ik to e_i
24              end for
25          end for
26      end for
27  end for
```

## Dataset Generator Output



| U-1 30 days (U-B_1, U-B_2, U-B_3) | U-2 30 days | U-3 30 days | U-4 30 days (U-B_1, U-B_2, U-B_3) |
|---|---|---|---|
| 1  1  U-1 D-1 T-7 R-5 L-6 A-1 P-190 | 2  1  U-1 D-2 T-2 R-5 L-3 A-6 P-94 | Random Values based on F_k | Random Values based on F_k | 91  1  U-4 D-1 T-2 R-9 L-12 A-6 P-252 |

A slice of generated event dataset for 4 users where the average events per user per day is 20. Users U-1 and U-4 are selected for insertion of 3 behavior patterns B1, B2 and B3, which are highlighted with different colors.

# Implementation

## Dataset Design

- Design an event dataset to simulate user-system interactions in distributed medical imaging systems.
- Each event has 6 attributes, where *Event=<User, Location, Action, Patient, Date, Time>*.
- Table II defined attribute distributions.
- Table III defined 9 typical user behavior patterns that constitute ordering, timing, and sequencing.
- Produced 30,000 events with randomly selected attribute values but following predefined distribution; predefined behavior patterns are inserted into the events.

**TABLE II.   ATTRIBUTE DISTRIBUTION DEFINITION**

| Attribute | Rep | Domain | Type | Mu | Sigma |
|---|---|---|---|---|---|
| User | U- | 100 | Random | | |
| Location | L- | 15 | Normal | 8 | 4 |
| Action | A- | 16 | Normal | 8 | 3 |
| Patient | P- | 300 | Normal | 150 | 50 |
| Date | D- | 30 | Normal | 15 | 5 |
| Time | T- | 24 | Normal | 11 | 4 |

**TABLE III.   ANALYST DEFINED BEHAVIOR PATTERNS**

| Pattern Id | Sequence | Support |
|---|---|---|
| P-00001 | office-1-Juravinski-Hamilton, office-3-Juravinski-Hamilton, office-4-Juravinski-Hamilton | 30% |
| P-00002 | office-3-Lakeridge-Oshawa, office-4-Lakeridge-Oshawa, office-5-Lakeridge-Oshawa | 25% |
| P-00003 | office-2-McMaster-Hamilton, office-1-McMaster-Hamilton, office-3-McMaster-Hamilton | 20% |
| P-00004 | read exam, read report, update report | 30% |
| P-00005 | read exam, read order, create exam | 25% |
| P-00006 | create profile, read profile, update profile | 20% |
| P-00007 | 11:00, 12:00, 14:00 | 30% |
| P-00008 | 10:00, 11:00, 12:00 | 25% |
| P-00009 | 14:00, 15:00, 16:00 | 20% |

## Visually Analysis of Generated Dataset

We developed a toolkit that can produce the following visual graphs for analyzing the dataset:

- *Sequence overview*: sequence is an ordered list of events performed by one person per day.
- *Frequent sequential patterns*: are subsequences that appear frequently among all user sequences.
- *Clustering based on sequence similarity*: divides the frequent sequential patterns into a number of clusters
- *Clustering representatives*: explores the representative patterns of each cluster