



DETC2013-13155

A Privacy Preserving Data Mining Methodology for Dynamically Predicting Emerging Human Threats

Tuesday, August 6th, 2013

Gautam Manohar & Conrad S. Tucker

**{gautam.atulya@gmail.com,
ctucker4@psu.edu, }**



Presentation Overview

- Research Motivation and Background
- Methodology
 - The Knowledge Discovery process
 - Data Acquisition and Storage
 - Data Mining Predictive Model Construction
 - Result Interpretation and Output
- Application Case Study
- Results and Discussion
- Conclusion and Path Forward





RESEARCH MOTIVATION



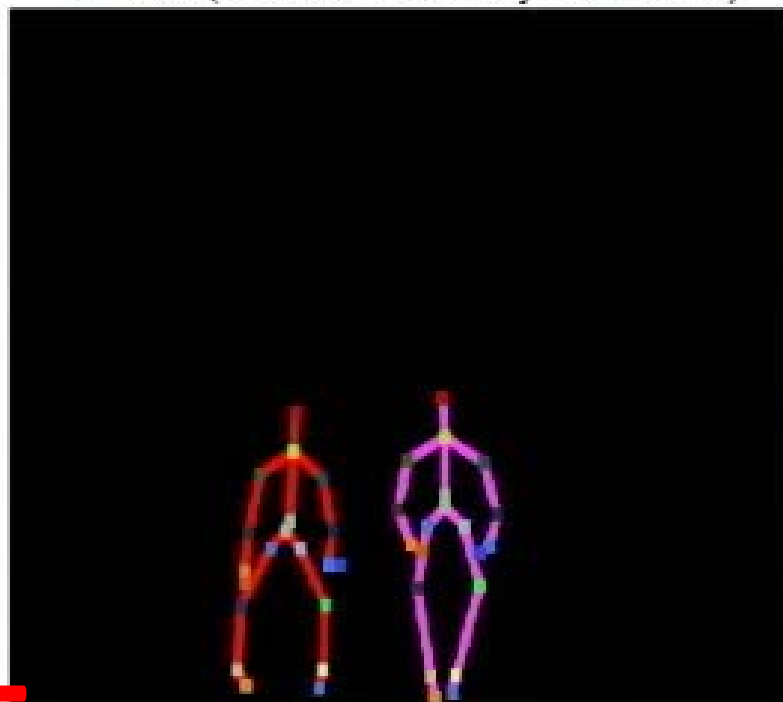
Motivation



Depth Stream



Skeleton (rendered if full body fits in frame)



Capturing Emergence

30 fps



Windows Explorer sidebar showing a file tree:

- Explorer
- bin\ Kinect (5 projects)
- MetaCore
- MetaPlayer
- MetaRecorder
- MetaViewer**
- Properties
- References
- app.config
- app.xaml
- MainWindow.xaml
- ReadMe.txt
- SkeletalViewer.ico
- VisualKinecl
- Properties
- References
- EventArgs
- EventData
- Timeline
- VirtualNLI
- Clics
- Views
- AssemblyName

System tray: 3:54 PM, 2/25/2011

Motivation and Background

- Existing systems are passive and more useful for post-incident analysis.
- Privacy issues with most existing systems become a hindrance in public use (I.e. the need to preserve **Personally Identifiable Information (PII)**)

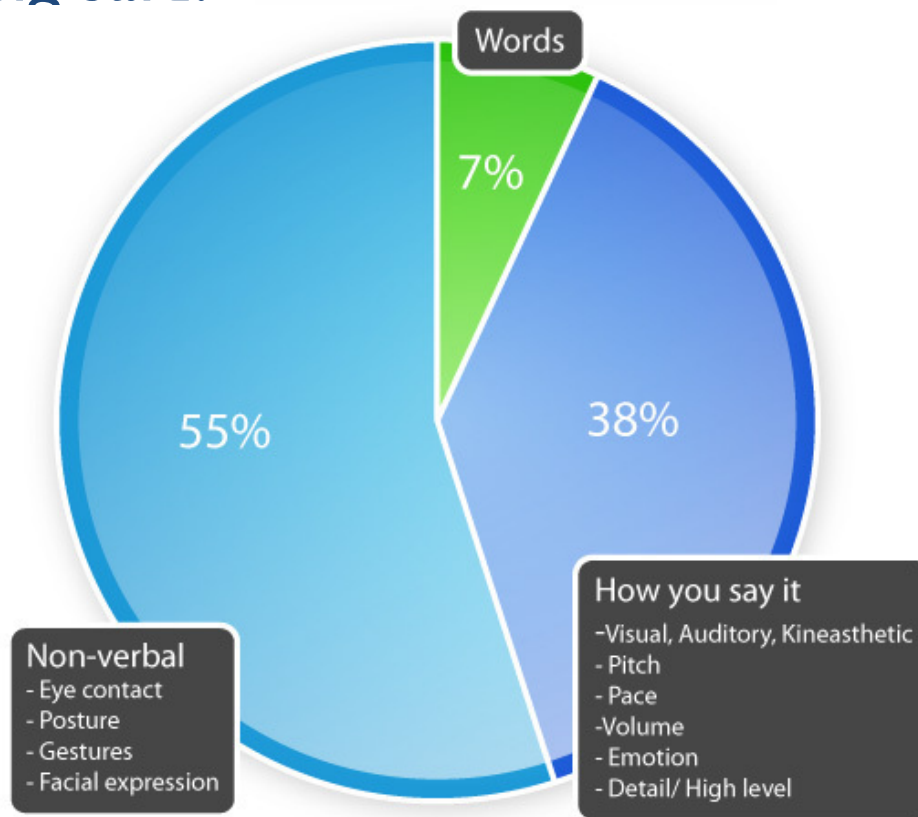




Why Individual Body Movement Data?

“The most important thing in communication is to hear what isn't being said.”

Peter F. Drucker





RESEARCH METHODOLOGY

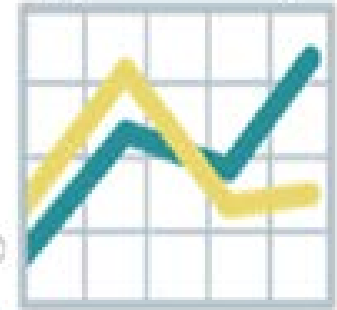
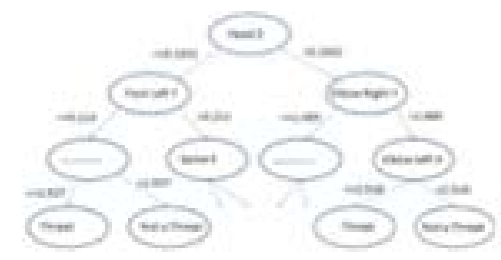
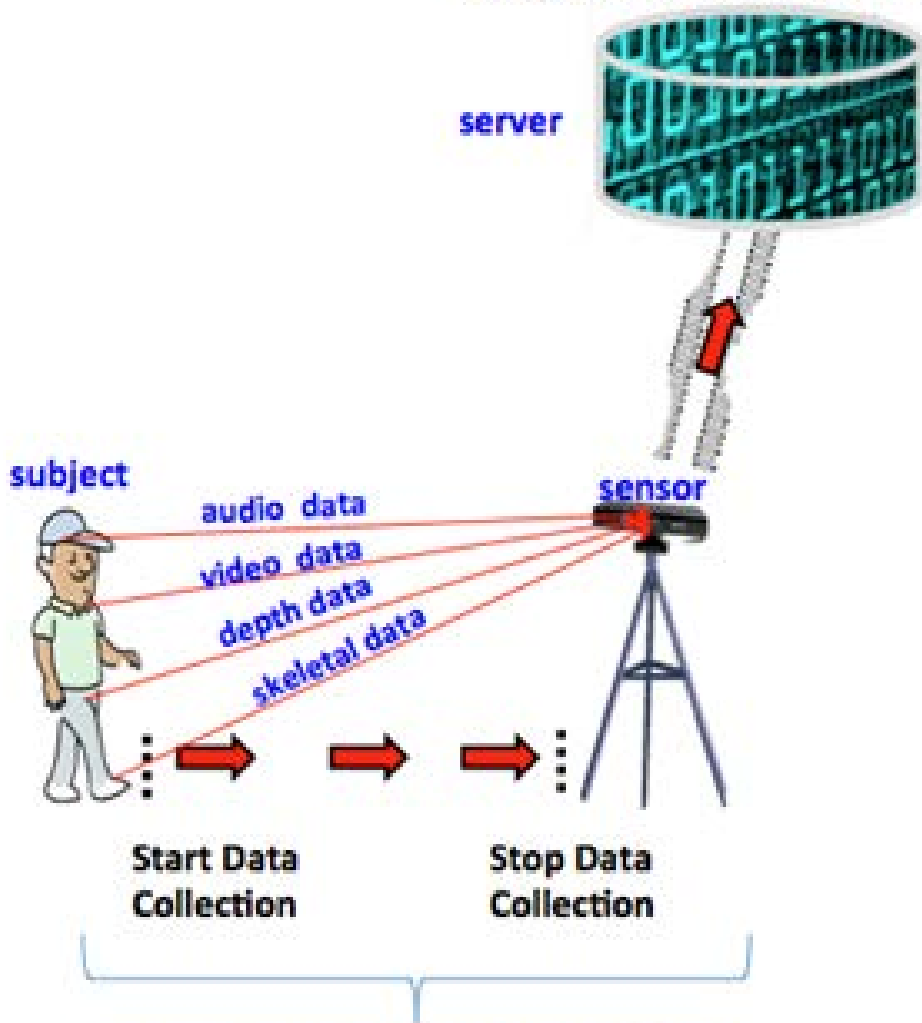


Proposed Methodology

Step 2: Data Transfer and Storage



Step 3: Data Mining Knowledge Discovery



Step 4: Decision support GUI

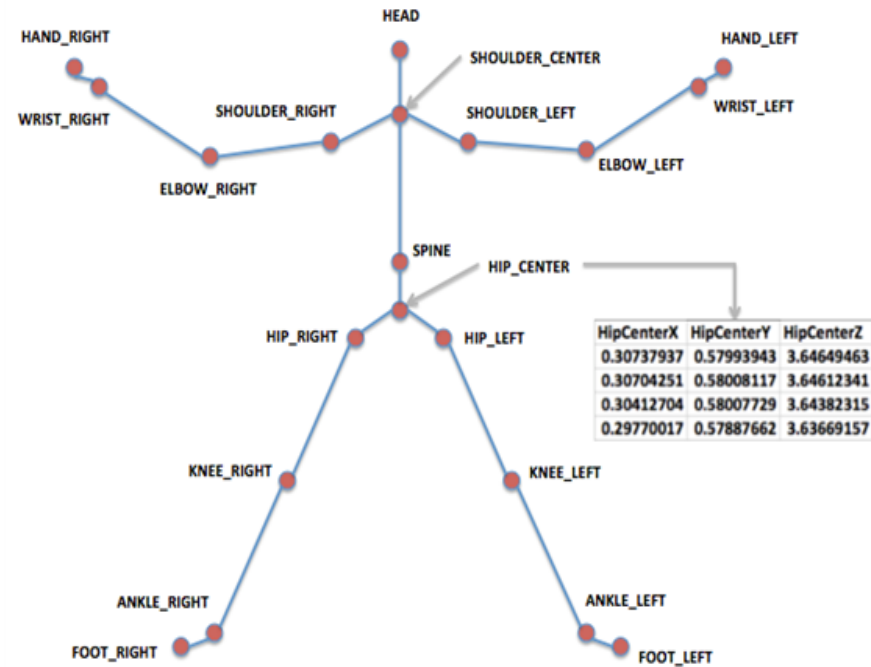


Step 1: Sensor Data Acquisition



Step 1: Data Acquisition

- Data acquisition hardware setup consists of a sensor system with:
 - an RGB video camera, and
 - an infrared depth sensor
- Output from sensors is used to create a virtual skeleton of the subject with 20 nodes as shown
- Each nodes collects data pertaining to:
 - 3D Spatial Coordinates (X,Y,Z)
 - Timestamp



High Fidelity Data, Privacy Preserving

Velocities of each node



Large Scale Data Base

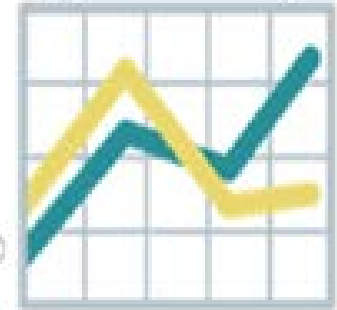
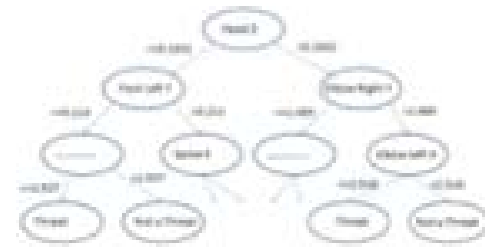
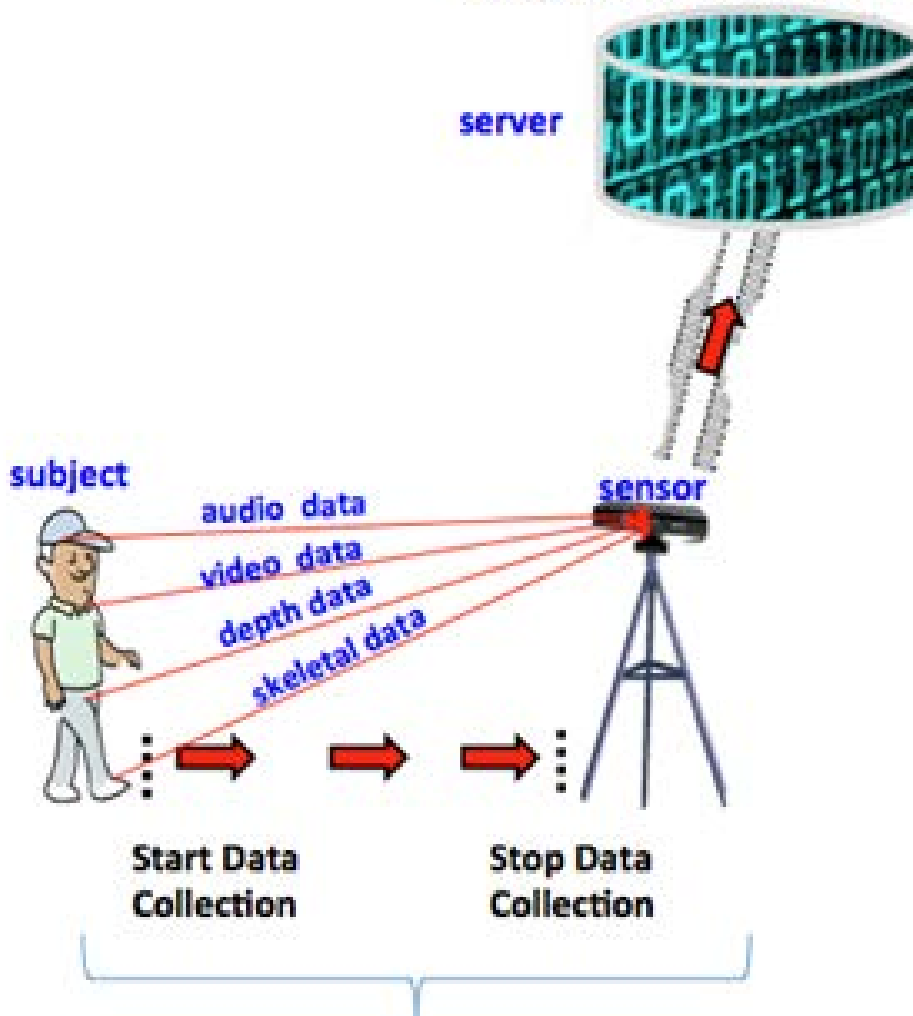
| Head X-Coordinate | Head Y-Coordinate | Head Z-Coordinate | Spine X-Coordinate | Spine Y-Coordinate | Spine Z-Coordinate | Shoulder X-Coordinate | Shoulder Y-Coordinate | Shoulder Z-Coordinate |
|-------------------|-------------------|-------------------|--------------------|--------------------|--------------------|-----------------------|-----------------------|-----------------------|
| 0.6400 | 0.7900 | 0.3500 | 0.7000 | 0.2400 | 0.5000 | 0.2600 | 0.1200 | 0.4300 |
| 0.8200 | 0.3800 | 0.7200 | 0.0600 | 0.4600 | 0.0700 | 0.2500 | 0.7700 | 0.6900 |
| 0.2300 | 0.9500 | 0.9200 | 0.9100 | 0.9700 | 0.3800 | 0.8100 | 0.6400 | 0.3000 |
| 0.9700 | 1.0000 | 0.4700 | 0.8600 | 0.7100 | 0.2400 | 0.1300 | 0.4600 | 0.2600 |
| 0.9700 | 0.2800 | 0.7100 | 0.1600 | 0.0700 | 0.0900 | 0.0700 | 0.5600 | 0.9200 |
| 0.4800 | 0.9900 | 0.9800 | 0.4300 | 0.6800 | 0.0400 | 0.4100 | 0.4900 | 0.7700 |
| 0.5000 | 0.7300 | 0.6700 | 0.6600 | 0.1500 | 0.0900 | 0.1700 | 0.6100 | 0.7600 |
| 0.2500 | 0.8500 | 0.6000 | 0.3900 | 0.9300 | 1.0000 | 0.8200 | 0.5100 | 0.4100 |
| 0.5000 | 0.6300 | 0.7900 | 0.0500 | 0.8100 | 0.8000 | 0.9300 | 1.0000 | 0.0500 |
| 0.6000 | 0.5400 | 0.5800 | 0.2300 | 0.2200 | 0.1200 | 0.1200 | 0.8500 | 0.8200 |
| 0.8300 | 0.7900 | 0.5400 | 0.5900 | 0.3200 | 0.9900 | 0.4100 | 0.5600 | 0.1400 |
| 1.0000 | 0.6900 | 0.6300 | 0.4700 | 0.6100 | 0.2400 | 0.1200 | 0.5700 | 0.7800 |
| 0.4400 | 0.0700 | 0.1100 | 0.7200 | 0.9000 | 0.4600 | 0.5100 | 0.8800 | 0.2600 |
| 0.1600 | 0.4700 | 0.7400 | 0.7000 | 0.5800 | 0.7100 | 1.0000 | 0.9100 | 0.2400 |
| 0.4300 | 0.6500 | 0.0000 | 0.2400 | 0.0200 | 0.6500 | 0.8300 | 0.8100 | 0.0000 |
| 0.0700 | 0.5700 | 0.3300 | 0.2900 | 0.3300 | 0.4000 | 0.7500 | 0.6500 | 0.6100 |
| 0.4900 | 0.9200 | 0.7100 | 0.3200 | 0.6300 | 0.4400 | 0.8000 | 0.9100 | 0.4200 |
| 0.6300 | 0.1300 | 0.4600 | 0.2100 | 0.2300 | 0.3100 | 0.3200 | 0.5400 | 0.7700 |
| 0.9500 | 0.1100 | 0.8100 | 0.1400 | 0.2300 | 0.8800 | 0.7200 | 0.1100 | 0.7200 |
| 0.7000 | 0.4700 | 0.1100 | 0.2800 | 0.0600 | 0.9700 | 0.2800 | 0.1700 | 0.6700 |
| 0.1300 | 0.4200 | 0.5300 | 0.1400 | 0.7300 | 0.7100 | 0.7200 | 0.3900 | 0.4400 |
| 0.5800 | 0.8800 | 0.7600 | 0.1600 | 0.6000 | 0.8300 | 0.1800 | 1.0000 | 0.4400 |
| 0.9800 | 0.9300 | 0.9900 | 0.4400 | 0.0600 | 0.0800 | 0.6100 | 0.9500 | 0.6100 |
| 0.0700 | 0.7200 | 0.7200 | 0.4200 | 0.8400 | 0.1700 | 0.0200 | 0.7600 | 0.7700 |
| 0.8800 | 0.6900 | 0.4400 | 0.4800 | 0.8900 | 0.3700 | 0.7500 | 0.8500 | 0.1700 |



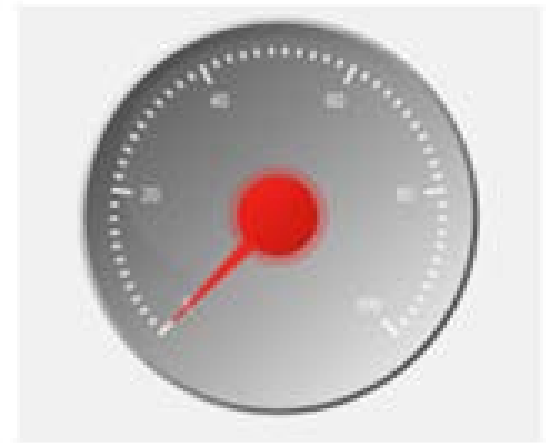
Proposed Methodology

Step 2: Data Transfer and Storage

Step 3: Data Mining
Knowledge Discovery



Step 4: Decision support GUI



Step 1: Sensor Data Acquisition





Step 2: Data Transfer and Storage

- The data is stored in a structured Relational Database with fields for the following measures:
 - Timestamp
 - Euclidean Coordinates
 - Velocities of each node
 - Boolean “Threat Class” defining whether the data collected during training was for a threat action or not.





Step 2: Data Transfer and Storage

- The data is stored in a structured Relational Database with fields for the following measures:

| Timestamp | Node X Coordinate | Node Y <u>Coord.</u> | Node Z <u>Coord.</u> | Velocity of Node X <u>Coord.</u> | | Threat Class |
|-----------|-------------------|----------------------|----------------------|----------------------------------|-------|--------------|
| 1346879 | 3.04356 | 2.98750 | 1.25673 | 0.25677 | | FALSE |
| 1346902 | 3.25831 | 5.63178 | 4.77721 | 6.78103 | | TRUE |
| ⋮ | | | | | | |

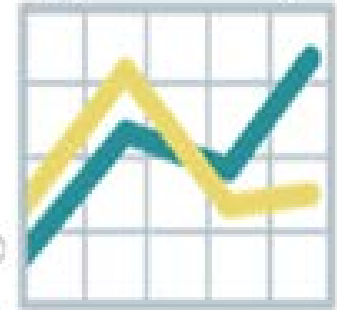
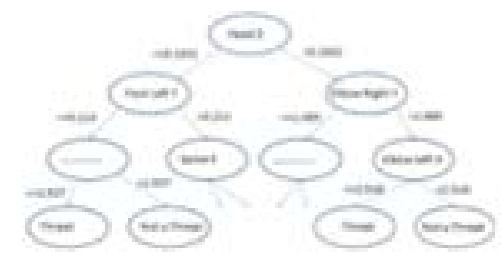
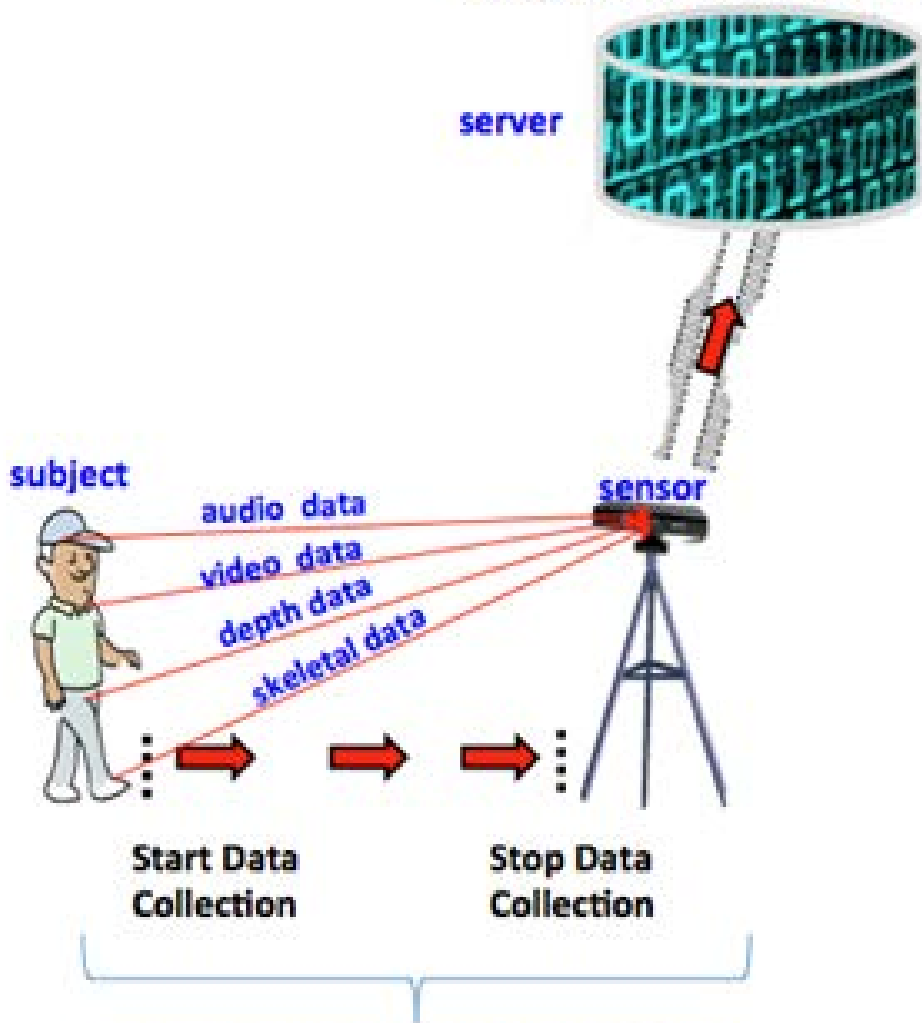


Proposed Methodology

Step 2: Data Transfer and Storage



Step 3: Data Mining
Knowledge Discovery



Step 4: Decision support GUI

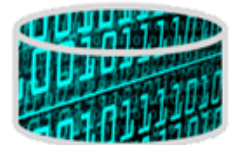
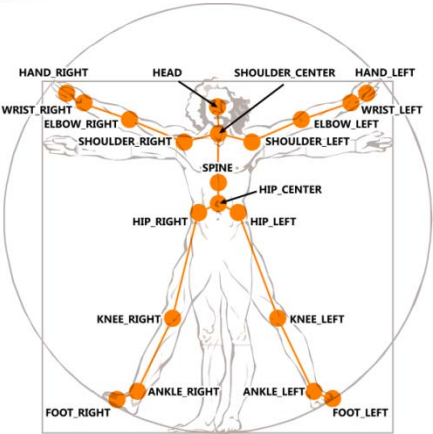


Step 1: Sensor Data Acquisition





Step 3: Data Mining/Knowledge Discovery

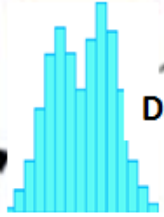


Data

Data Selection and Cleaning



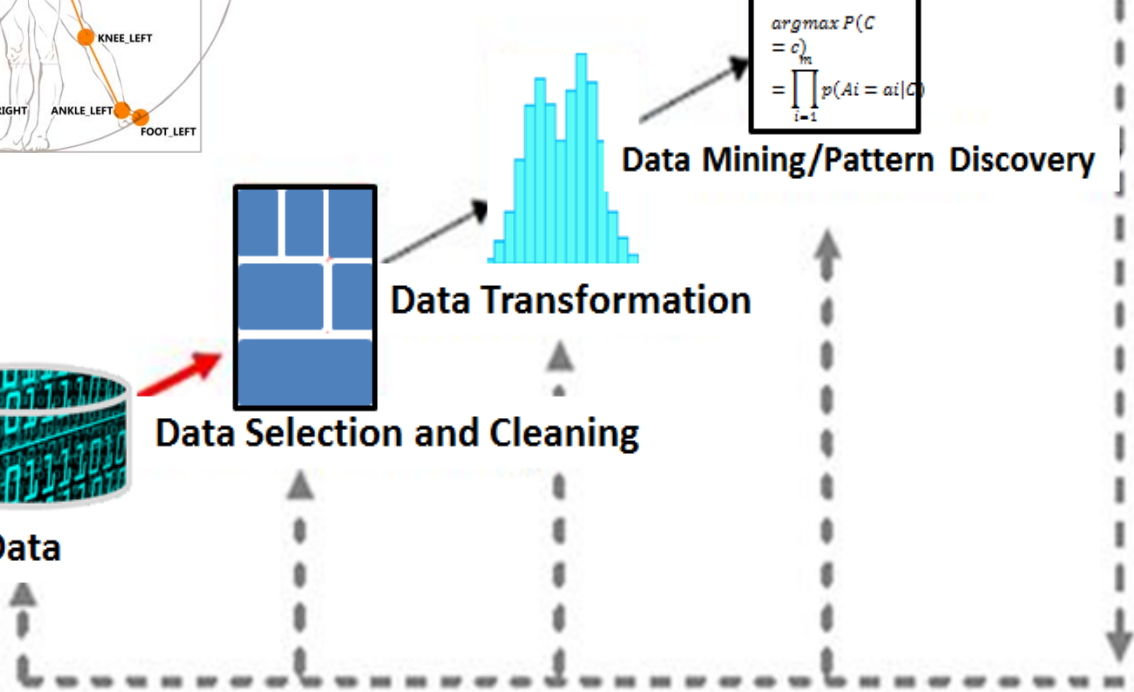
Data Transformation



Data Mining/Pattern Discovery

$$\begin{aligned}
 & \text{info}_m(A) \\
 & - \sum_{i=1}^m f_i \log_2 f_i \\
 & \text{argmax } P(C) \\
 & = c_m \\
 & = \prod_{i=1}^m p(A_i = a_i | C)
 \end{aligned}$$

Interpretation/Evaluation





Knowledge Discovery in Data Bases

$$\begin{aligned}
 & \text{info}(A) \\
 & - \sum_{i=1}^m f_i \log_2 f_i \\
 & \text{argmax } P(C \\
 & = c) \\
 & = \prod_{i=1}^m p(A_i = a_i | C)
 \end{aligned}$$

Data Mining/Pattern Discovery

Supervised Learning

Unsupervised Learning



Supervised VS Unsupervised Learning

Supervised

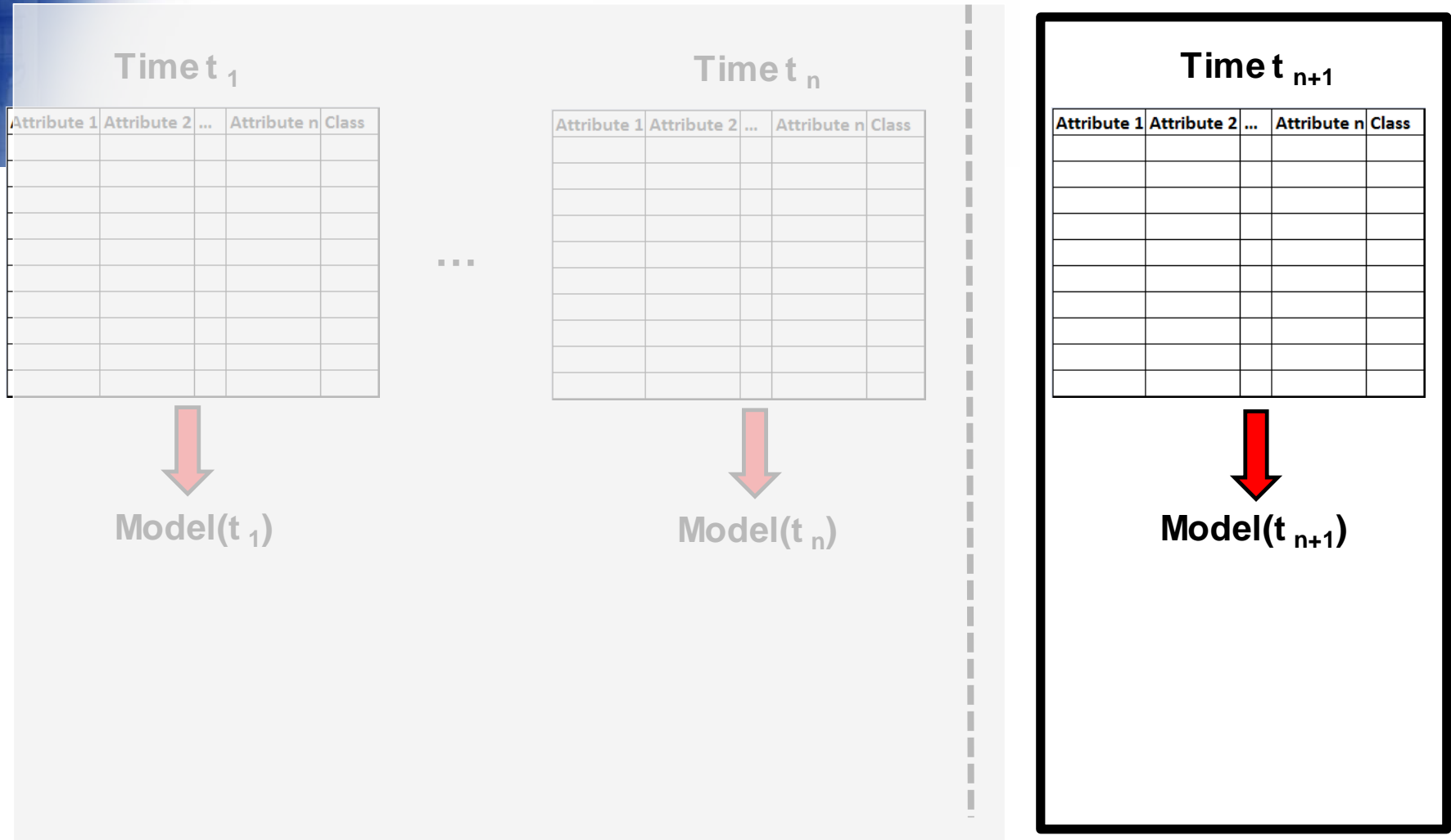
- $y=F(x)$: true function
- D : labeled training set
- $D: \{x_i, F(x_i)\}$
- Learn:
 $G(x)$: model trained to predict labels D
- Goal:
 $E[(F(x)-G(x))^2] \approx 0$
- Well defined criteria:
Accuracy, RMSE, ...

Unsupervised

- Generator: true model
- D : unlabeled data sample
- $D: \{x_i\}$
- Learn
Underlying data structure
- Goal:
Find natural patterns
- Well defined criteria:
varies



Capturing Threat Emergence



Data Mining Decision Tree Induction

Given a time stamped Data Set (t),

| Feature 1 | Feature 2 | ... | Feature N | Class |
|------------------|------------------|-----|------------------|------------------|
| A _{1,1} | A _{2,1} | | A _{N,1} | C _{j,1} |
| . | . | | . | . |
| . | . | | . | . |
| . | . | | . | . |
| A _{1,M} | A _{2,M} | | A _{N,M} | C _{j,M} |

$$Entropy(T) = - \sum_j p(C_j | T) \log_2 p(C_j | T)$$

$$GAIN(X) = Entropy(T) - \left(\sum_{i=1}^k \frac{T_i}{T} Entropy_X(T_i) \right)$$

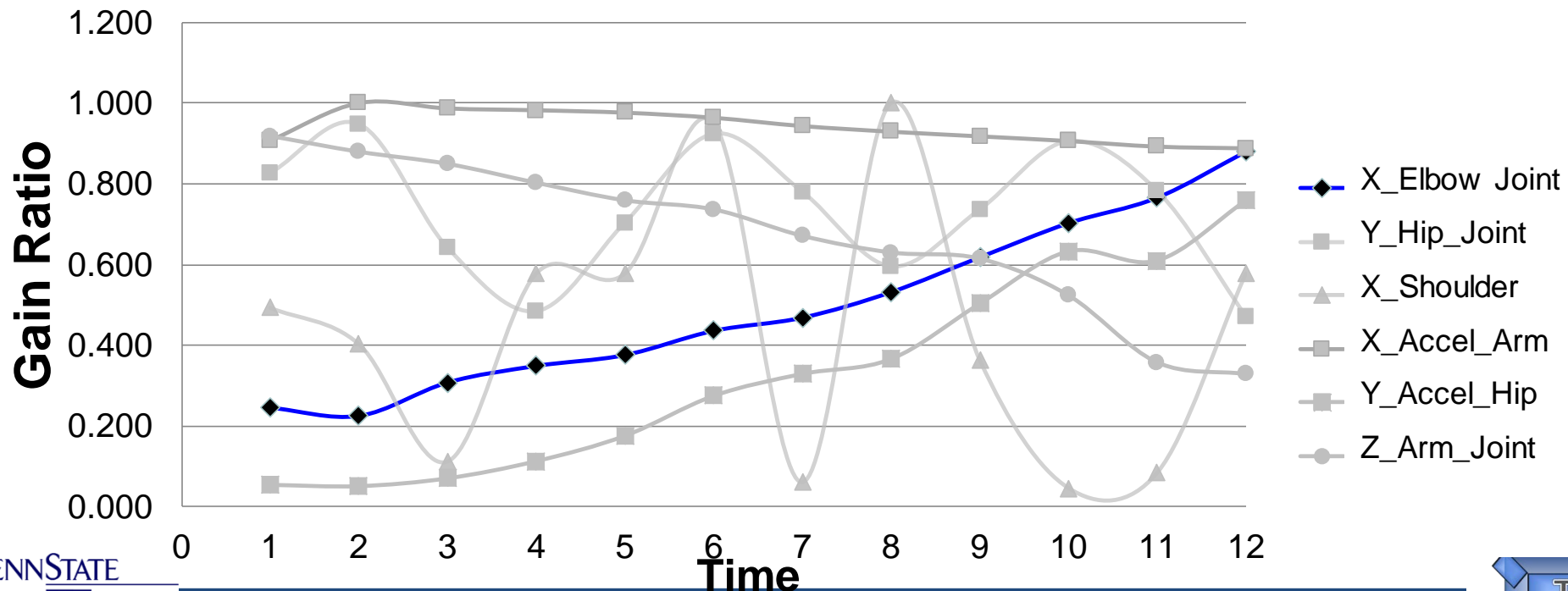
$$Gain\ ratio(X) = \frac{Gain(X)}{- \sum_{i=1}^k \frac{|T_i|}{|T|} \cdot \log_2 \frac{|T_i|}{|T|}}$$

Tucker C., H.M. Kim, "Trend Mining for Predictive Product Design", Transactions of ASME: Journal of Mechanical Design, Vol. 133, No. 11, 2011.



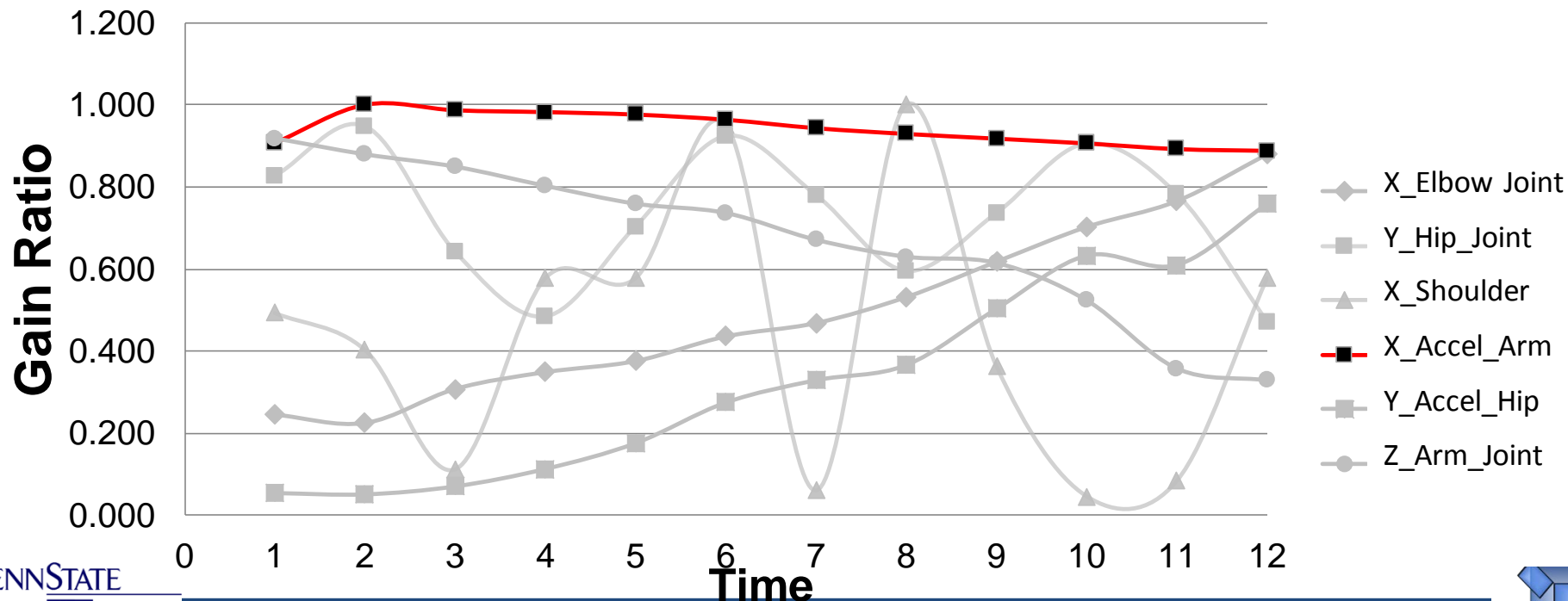
| Features | Time Series Gain Ratio | | | | | | | | | | | | Predict |
|---------------|------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------------|
| | t1 | t2 | t3 | t4 | t5 | t6 | t7 | t8 | t9 | t10 | t11 | t12 | t13_predict |
| X_Elbow_Joint | 0.245 | 0.225 | 0.308 | 0.349 | 0.376 | 0.436 | 0.468 | 0.532 | 0.618 | 0.702 | 0.765 | 0.879 | 0.919 |
| Y_Hip_Joint | 0.827 | 0.948 | 0.642 | 0.485 | 0.704 | 0.924 | 0.780 | 0.596 | 0.737 | 0.906 | 0.782 | 0.472 | 0.789 |
| X_Shoulder | 0.493 | 0.403 | 0.112 | 0.578 | 0.578 | 0.951 | 0.061 | 1.000 | 0.363 | 0.046 | 0.084 | 0.578 | 0.541 |
| X_Accel_Arm | 0.907 | 1.000 | 0.987 | 0.982 | 0.976 | 0.963 | 0.943 | 0.929 | 0.917 | 0.906 | 0.892 | 0.888 | 0.877 |
| Y_Accel_Hip | 0.054 | 0.051 | 0.070 | 0.113 | 0.176 | 0.275 | 0.329 | 0.366 | 0.503 | 0.633 | 0.610 | 0.759 | 0.842 |
| Z_Arm_Joint | 0.918 | 0.879 | 0.849 | 0.803 | 0.759 | 0.737 | 0.671 | 0.630 | 0.615 | 0.524 | 0.358 | 0.329 | 0.270 |

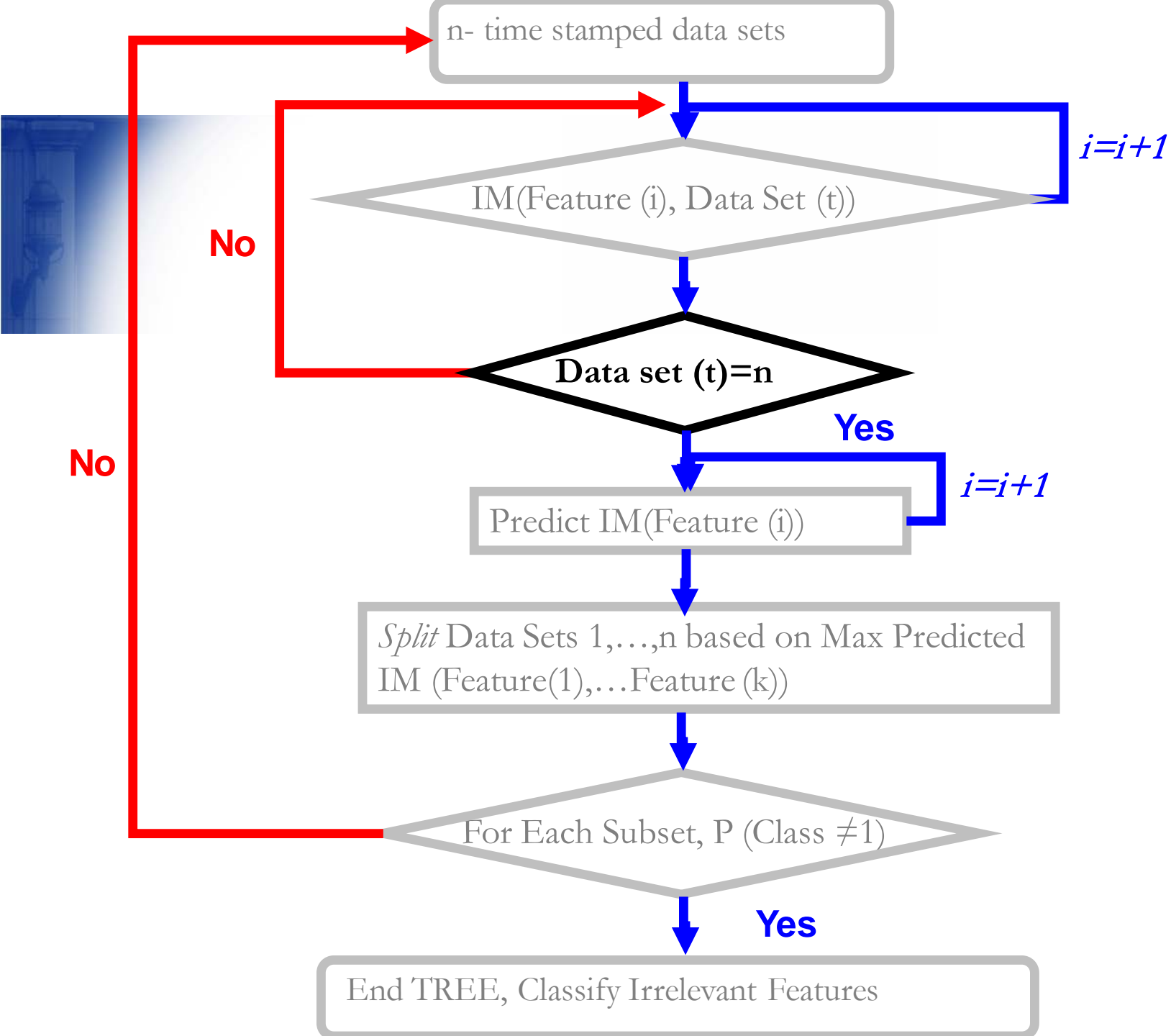
Feature Gain Ratio Plot Over Time

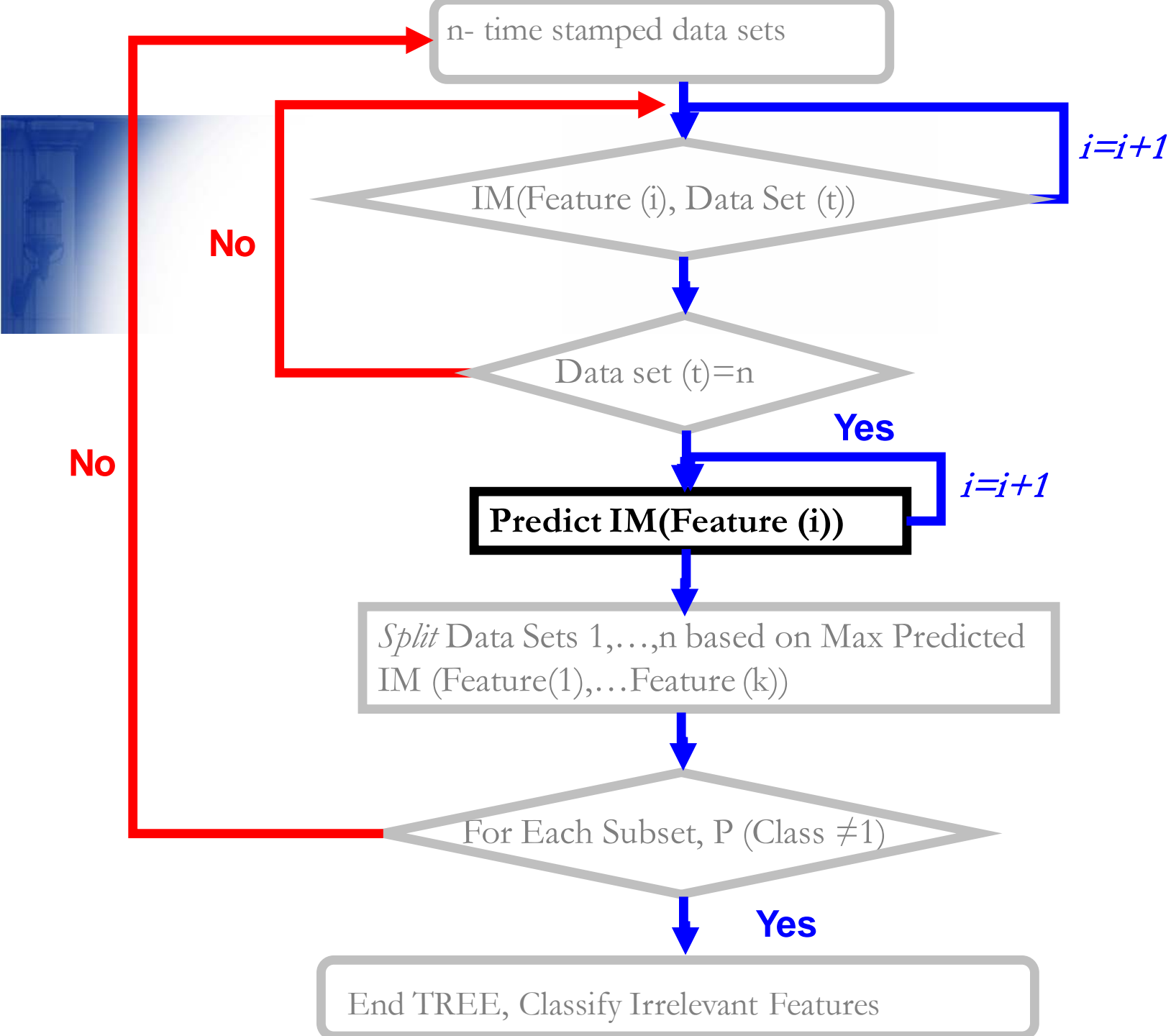


| Features | Time Series Gain Ratio | | | | | | | | | | | | Predict |
|---------------|------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------------|
| | t1 | t2 | t3 | t4 | t5 | t6 | t7 | t8 | t9 | t10 | t11 | t12 | t13_predict |
| X_Elbow_Joint | 0.245 | 0.225 | 0.308 | 0.349 | 0.376 | 0.436 | 0.468 | 0.532 | 0.618 | 0.702 | 0.765 | 0.879 | 0.919 |
| Y_Hip_Joint | 0.827 | 0.948 | 0.642 | 0.485 | 0.704 | 0.924 | 0.780 | 0.596 | 0.737 | 0.906 | 0.782 | 0.472 | 0.789 |
| X_Shoulder | 0.493 | 0.403 | 0.112 | 0.578 | 0.578 | 0.951 | 0.061 | 1.000 | 0.363 | 0.046 | 0.084 | 0.578 | 0.541 |
| X_Accel_Arm | 0.907 | 1.000 | 0.987 | 0.982 | 0.976 | 0.963 | 0.943 | 0.929 | 0.917 | 0.906 | 0.892 | 0.888 | 0.877 |
| Y_Accel_Hip | 0.054 | 0.051 | 0.070 | 0.113 | 0.176 | 0.275 | 0.329 | 0.366 | 0.503 | 0.633 | 0.610 | 0.759 | 0.842 |
| Z_Arm_Joint | 0.918 | 0.879 | 0.849 | 0.803 | 0.759 | 0.737 | 0.671 | 0.630 | 0.615 | 0.524 | 0.358 | 0.329 | 0.270 |

Feature Gain Ratio Plot Over Time









Holt-Winters Forecasting

The (k) step-ahead forecasting model is defined as:

$$\hat{y}_t(k) = L_t + kT_t + I_{t-s+k}$$

Where:

Level L_t (the level component):

$$L_t = \alpha(y_t - I_{t-s}) + (1 - \alpha)(L_{t-1} + T_{t-1})$$

Trend T_t (the slope component):

$$T_t = \gamma(L_t - L_{t-1}) + (1 - \gamma)T_{t-1}$$

Season I_t (the seasonal component):

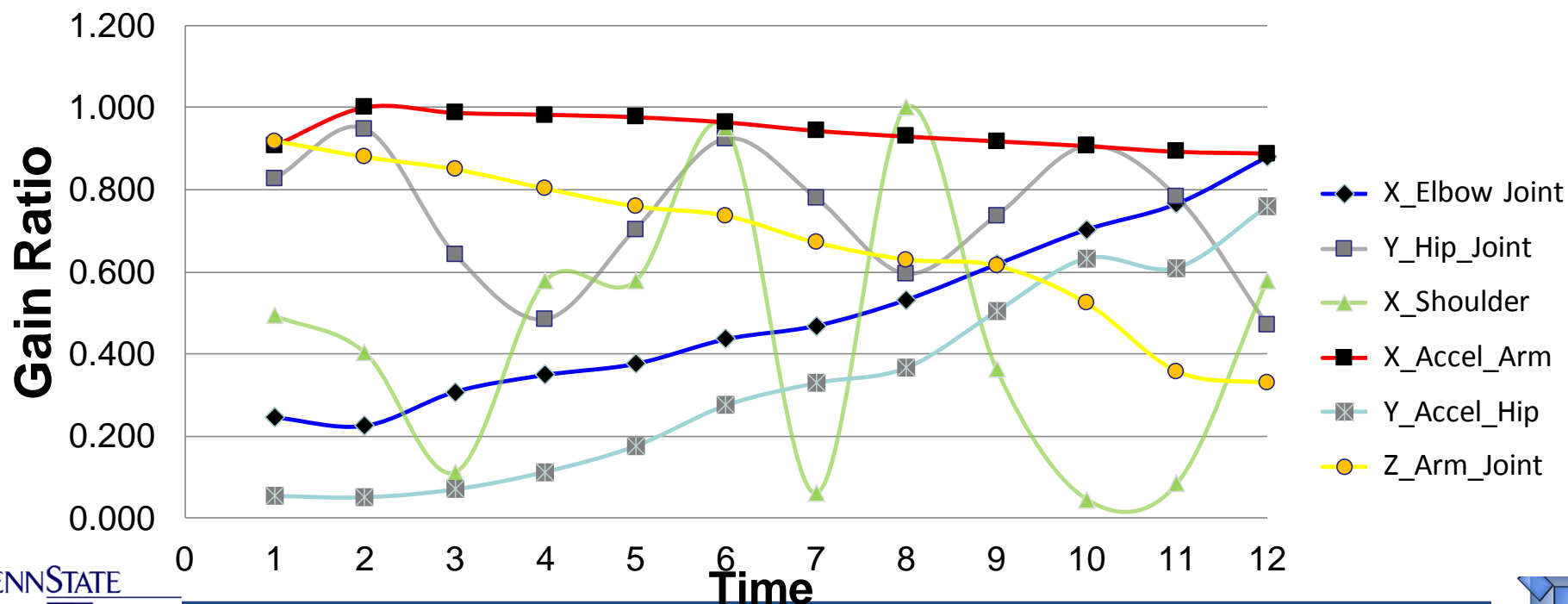
$$I_t = \delta(y_t - L_t) + (1 - \delta)I_{t-s}$$

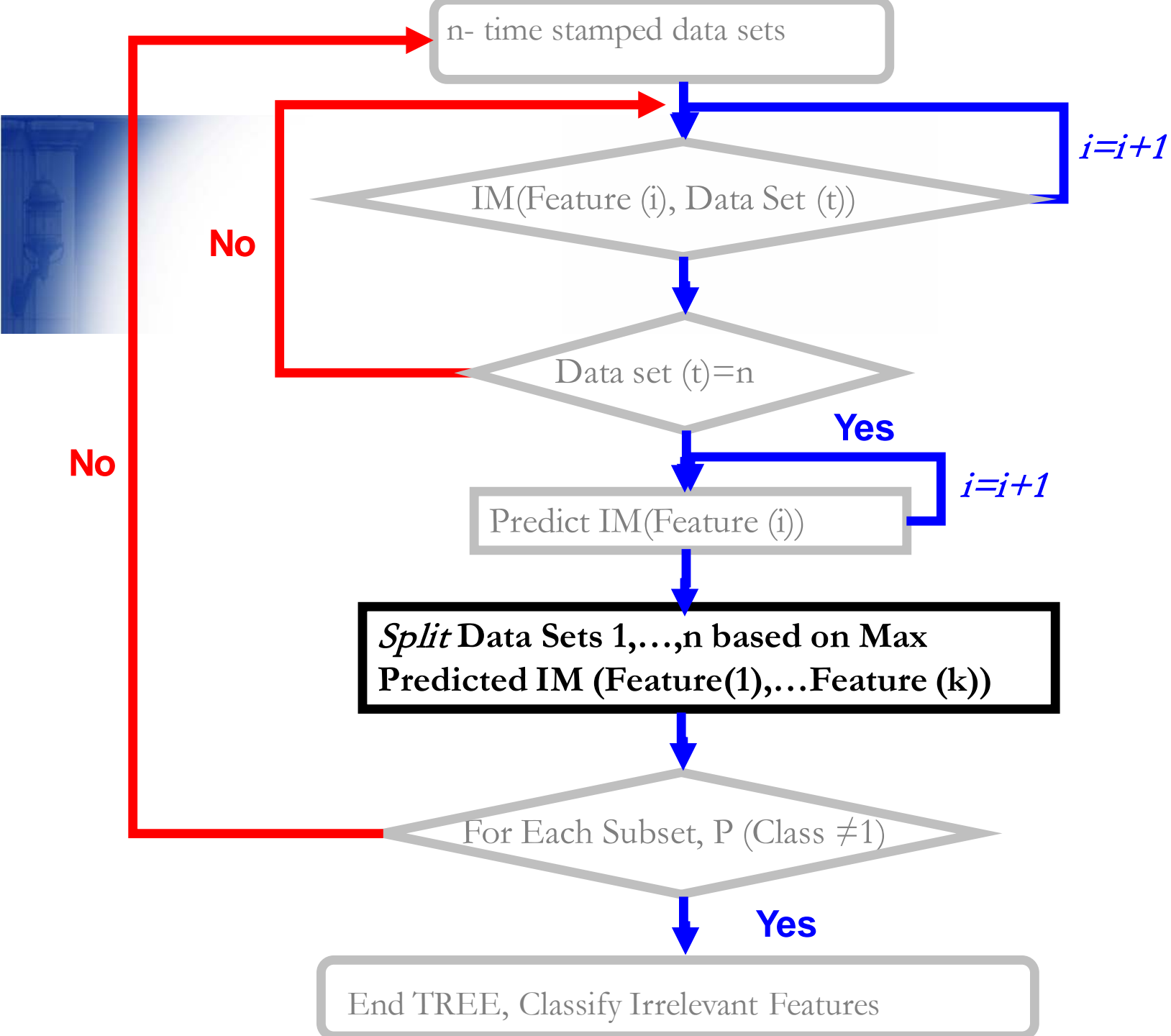
The smoothing parameters α, γ, δ , are in the range $\{0,1\}$



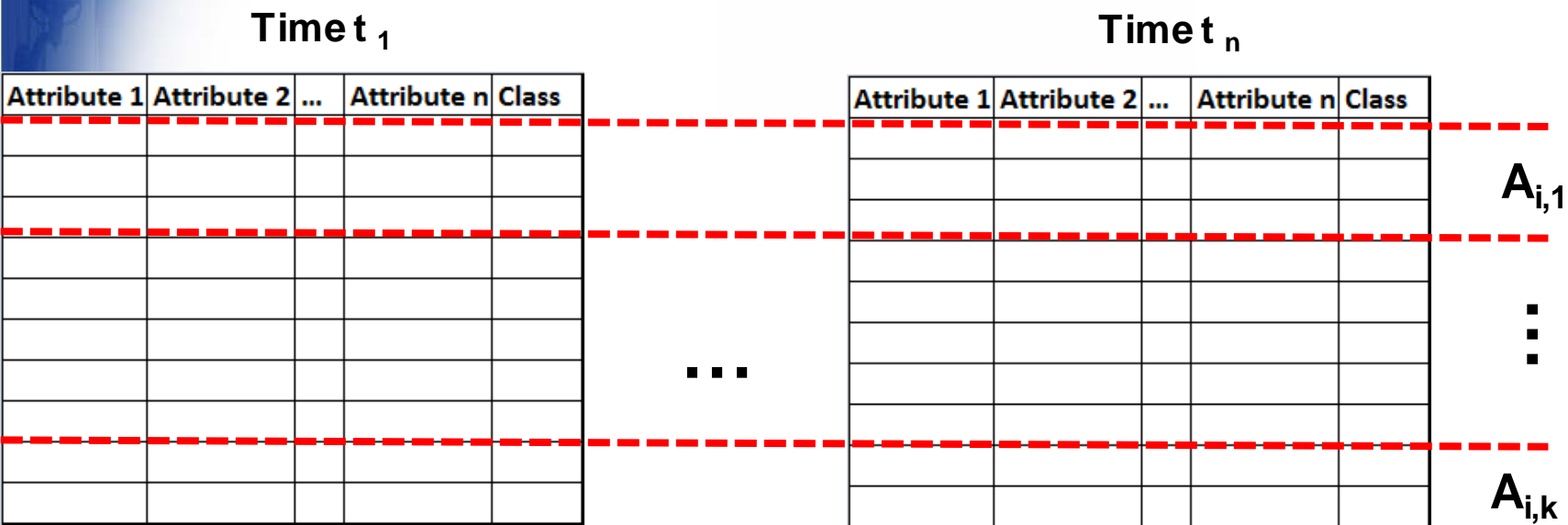
| PENNSTATE | Time Series Gain Ratio | | | | | | | | | | | | Predict |
|---------------|------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------------|
| Features | t1 | t2 | t3 | t4 | t5 | t6 | t7 | t8 | t9 | t10 | t11 | t12 | t13_predict |
| X_Elbow_Joint | 0.245 | 0.225 | 0.308 | 0.349 | 0.376 | 0.436 | 0.468 | 0.532 | 0.618 | 0.702 | 0.765 | 0.879 | 0.919 |
| Y_Hip_Joint | 0.827 | 0.948 | 0.642 | 0.485 | 0.704 | 0.924 | 0.780 | 0.596 | 0.737 | 0.906 | 0.782 | 0.472 | 0.789 |
| X_Shoulder | 0.493 | 0.403 | 0.112 | 0.578 | 0.578 | 0.951 | 0.061 | 1.000 | 0.363 | 0.046 | 0.084 | 0.578 | 0.541 |
| X_Accel_Arm | 0.907 | 1.000 | 0.987 | 0.982 | 0.976 | 0.963 | 0.943 | 0.929 | 0.917 | 0.906 | 0.892 | 0.888 | 0.877 |
| Y_Accel_Hip | 0.054 | 0.051 | 0.070 | 0.113 | 0.176 | 0.275 | 0.329 | 0.366 | 0.503 | 0.633 | 0.610 | 0.759 | 0.842 |
| Z_Arm_Joint | 0.918 | 0.879 | 0.849 | 0.803 | 0.759 | 0.737 | 0.671 | 0.630 | 0.615 | 0.524 | 0.358 | 0.329 | 0.270 |

Feature Gain Ratio Plot Over Time

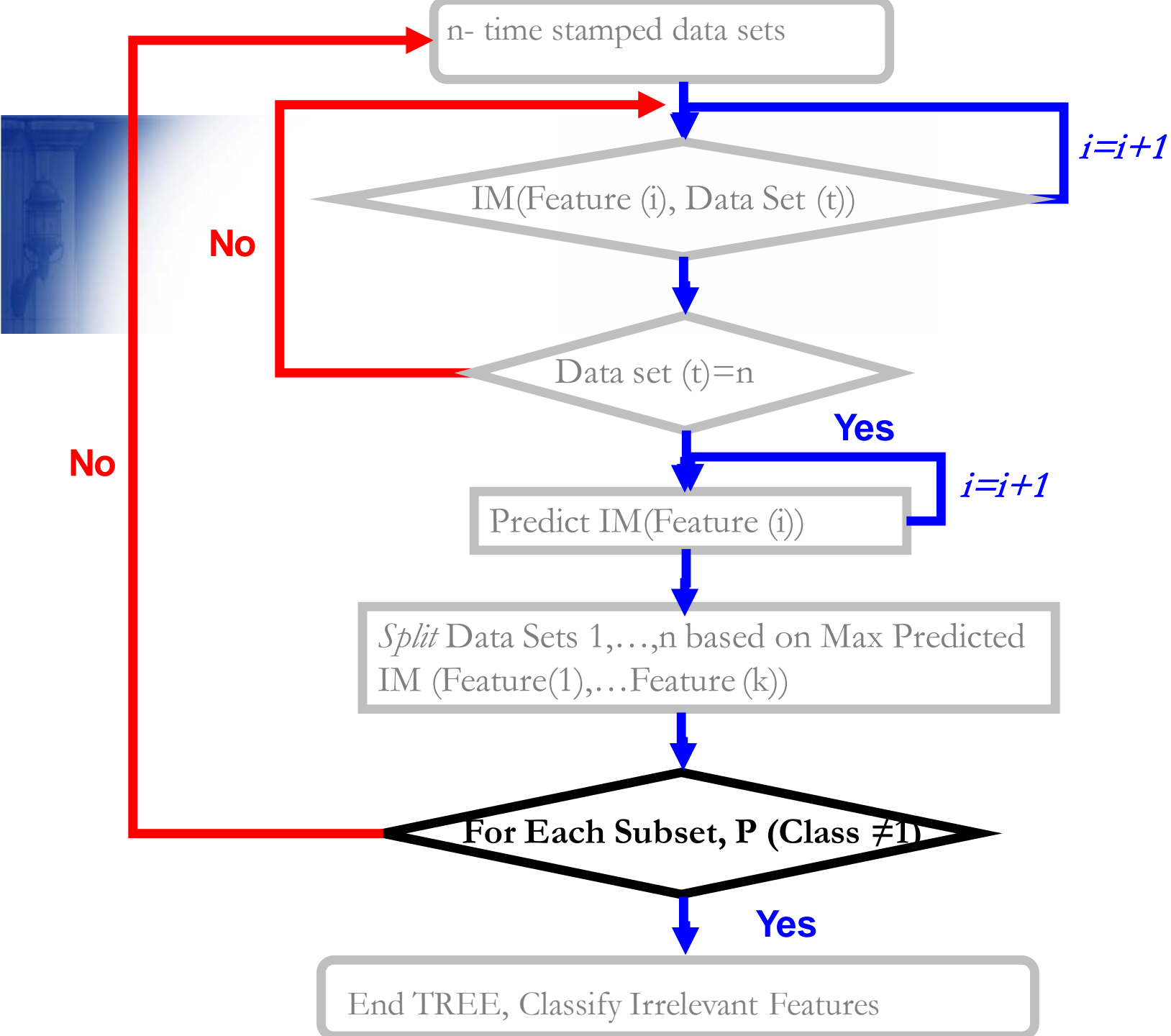


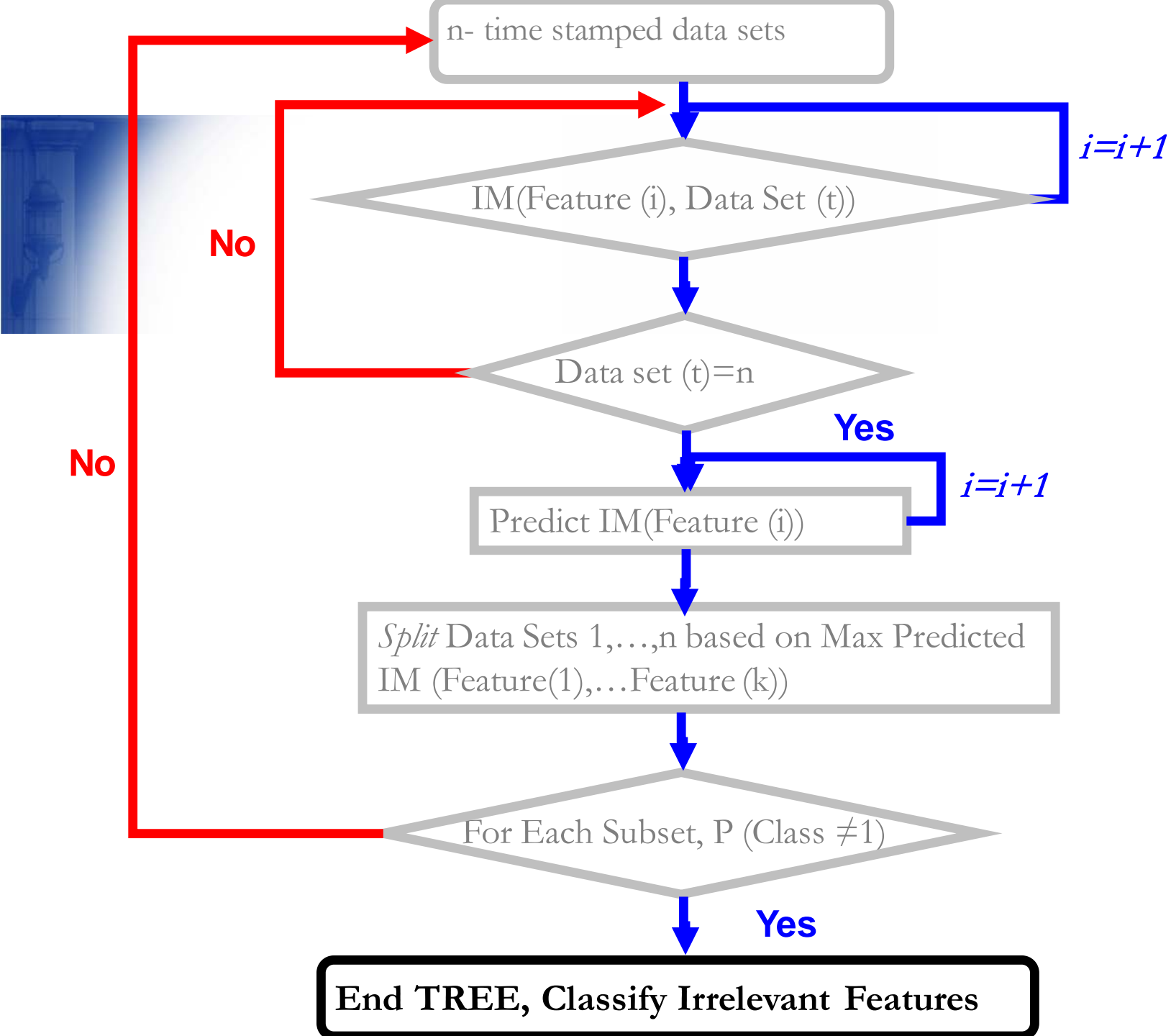


Split Data Sets (t_1, \dots, t_n) : Max IM



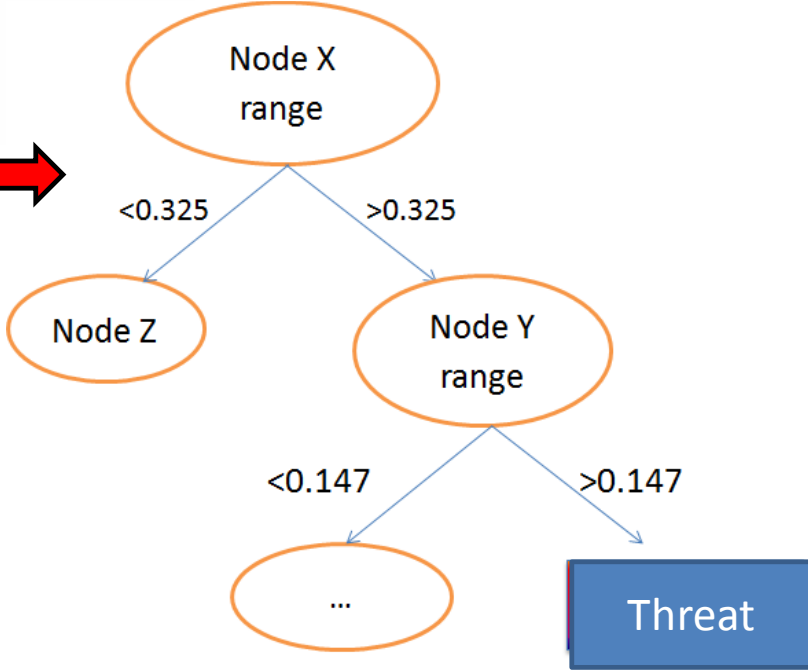
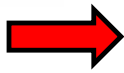
Split Data Sets (1,..,n) based on k mutually exclusive Feature values of Feature A_i





Data Mining Predictive Model

| Time t_1 | | | | | Time t_n | | | | |
|-------------|-------------|-----|-------------|-------|-------------|-------------|-----|-------------|-------|
| Attribute 1 | Attribute 2 | ... | Attribute n | Class | Attribute 1 | Attribute 2 | ... | Attribute n | Class |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |

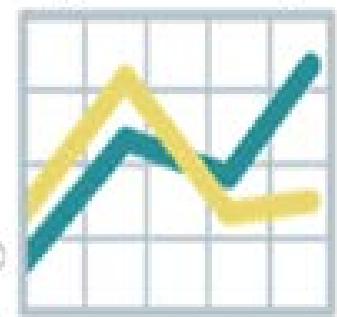
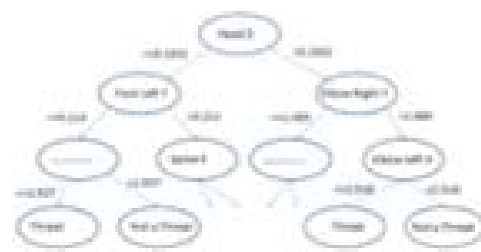
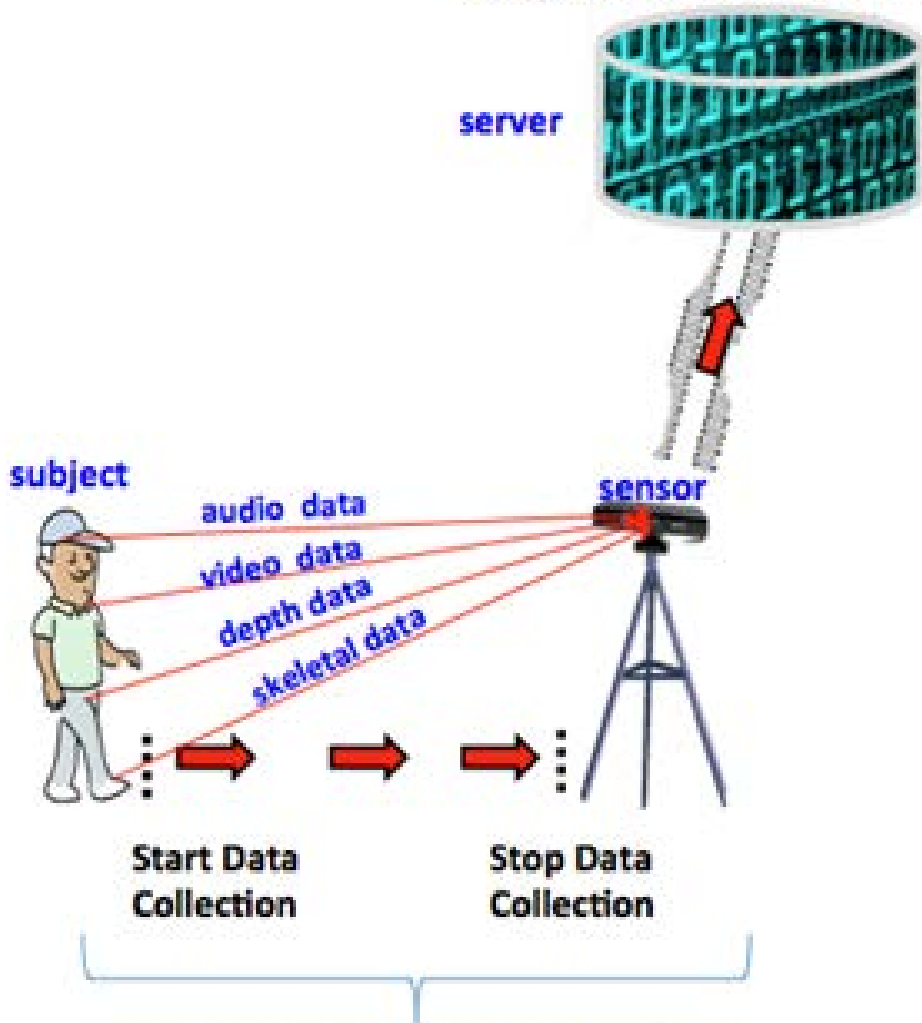


Proposed Methodology

Step 2: Data Transfer and Storage



Step 3: Data Mining
Knowledge Discovery



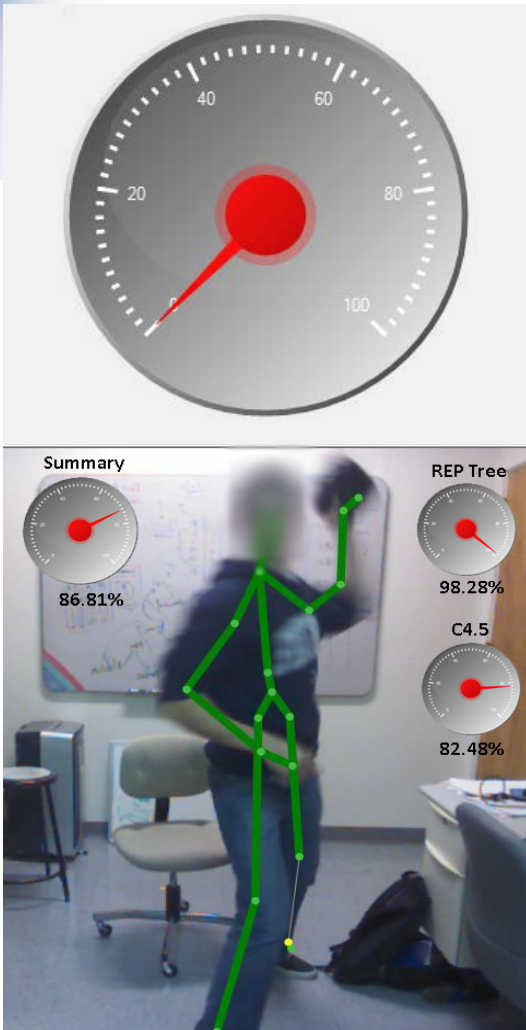
Step 4: Decision
support GUI



Step 1: Sensor Data Acquisition



Step 4: Decision Support



- Early Warning System (EWS) is a graphical user interface (GUI) that display the “percentage probability of threat/violent action being committed”.



APPLICATION CASE STUDY



Possible Threat Scenario



BBC UK (2008)

CASE STUDY: TEST DATA

- Voluntary participants from the University community were invited to enact the threat and non-threat actions
- Recreated in an indoor space, similar to a high profile speech
- The data collected is then used to train the predictive models
- The study was approved by the IRB and the ORP at the Pennsylvania State University, University Park campus, under the title “A Dynamic Pattern Recognition Framework for Mining and Predicting Emerging Threats” and is filed as IRB # 40258.
- **Study: 24 Subjects spanning 2 months**



THREAT PREDICTION RESULTS

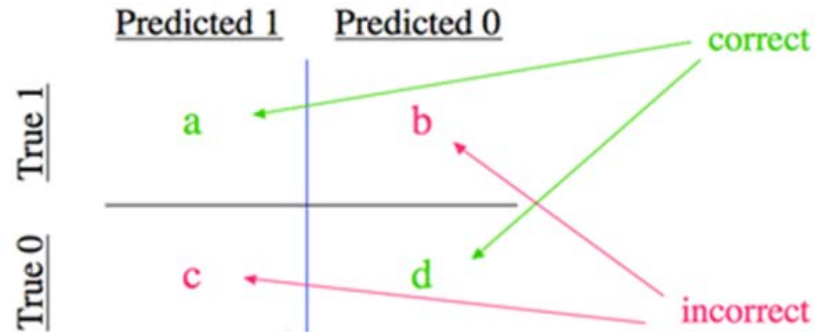
Low level threat prediction



High level threat prediction



RESULTS



$$\text{Accuracy} = \frac{a+d}{a+b+c+d}$$

$$\text{Precision} = \frac{a}{a+c}$$

$$\text{Recall} = \frac{a}{a+b}$$

Accuracy of Ensemble Methods: 86.8%



CONCLUSION AND FUTURE WORK



Conclusion and Future Work

- The most common surveillance systems today are reactive in nature and are not capable of actively predicting the emergence of a threat by analyzing past data collected.
- Privacy preserving data mining methodology
- This methodology takes the first step towards addressing these issues while providing promising results
- Expand the definition of “threat”





ACKNOWLEDGEMENTS AND REFERENCES

Contributors:

- Dr. Conrad S. Tucker, D.A.T.A. Lab members, Research Participants from PSU.

References:

1. Quinlan, J. R. C4.5: Programs for Machine Learning. Morgan Kaufmann Publishers, 1993.
2. Joshi, Karuna Pande. "Analysis of data mining algorithms." *University of Minnesota*. Retrieved July 25 (1997): 2005.
3. J. Han, M. Kamber, J. Pei, Data Mining: Concepts and Techniques, Third edition, 2011.
4. Data-Driven Decision Tree Classification for Product Portfolio Design Optimization, Conrad S. Tucker and Harrison M. Kim, J. Comput. Inf. Sci. Eng. 9, 041004 (2009), DOI:10.1115/1.3243634.
5. J. L. Raheja, A. Chaudhary, K. Singal, Tracking of fingertips and centers of palm using KINECT, International Conference on Computational Intelligence, Modeling & Simulation, 2011, 248-252.
6. Ya-Li Hou and Grantham K.H. Pang, Human detection in crowded scenes, IEEE international conference on image processing, 2010, 721-724.

