# Data Science in Biomedicine: Risk Estimators and Pattern Mining

Panos M. Pardalos<sup>1,2</sup>

<sup>1</sup>Center for Applied Optimization Department of Industrial and Systems Engineering University of Florida, USA

http://www.ise.ufl.edu/cao/

<sup>2</sup>Laboratory of Algorithms and Technologies for Networks Analysis (LATNA) National Research University Higher School of Economics, Russia

https://nnov.hse.ru/en/latna/

# BIOMAT 2017 (Moscow, Russia) October 29 - November 3, 2017

PM Pardalos Data Science in Biomedicine

# This talk is dedicated to the memory of my friend and colleague Chris Floudas



PM Pardalos Data Science in Biomedicine





# 2 Acute Kidney Injury and Sepsis

- sCr and 90-Day Mortality After Surgery
- Risk Prediction Models for Sepsis and AKI
- Multivariate Time Series Analysis for AKI Prediction

- The proliferation of massive datasets brings with it a series of special computational challenges.
- This data avalanche arises in a wide range of scientific and commercial applications.
- With advances in computer and information technologies, many of these challenges are beginning to be addressed

#### The 5 V's that define Big Data



#### http://bigdata.black/featured/what-is-big-data

PM Pardalos Data Science in Biomedicine

# **Data Growth**



PM Pardalos Data Science in Biomedicine

Acute Kidney Injury and Sepsis

# Historical Cost of Computer Memory and Storage



# Approximately 10 times cheaper every 5 years

http://www.jcmit.net/diskprice.
htm

Acute Kidney Injury and Sepsis

# **Memory Capacity**



# Approximately 10 times more every 5 years

http://www.jcmit.net

Acute Kidney Injury and Sepsis

# **Processing Power**



Microprocessor Transistor Counts 1971-2011 & Moore's Law

https://en.wikipedia.org/ wiki/Transistor\_count

every

# Some Examples of Massive Data Sources

- WWW, Internet
- Weather
- Telecommunications
- Media
- Demographics Data
- Financial Data
- Social Networks
- Biological Networks
- Health Care

# AT&T Call Graph

# • Call graph:

- vetrices are phone numbers
- two are connected if a phone call was made from one to another
- One-day call graph had 53,767,087 vertices and 170 million edges [1].
- [1] Abello, J., Pardalos, P.M, and Resende,M.G.C. On maximum clique problems in very large graphs. DIMACS series 50, 119-130 (1999).

Acute Kidney Injury and Sepsis

### **Human Genome Project**



[1] Eisenstein, M. Big data: The power of petabytes. Nature 527, S2S4 (2015).

# Human Brain

- Enormous size and complexity: **85 billion** neurons and **100 trillion** synapses [1].
- Small volumes of brain tissue can be fully reconstructed [2].
- 175 exabytes are projected to be required to store human brain connectome [3].
- Azevedo, F.A. *et al.* Equal numbers of neuronal and nonneuronal cells make the human brain an isometrically scaled-up primate brain. J Comp Neurol.513(5), 532-541 (2009).
- [2] Helmstaedter, M. Cellular-resolution connectomics: challenges of dense neural circuit reconstruction. Nature Methods 10, 501-507 (2013).
- [3] Mikula, S. Progress Towards Mammalian Whole-Brain Cellular Connectomics. Frontiers in Neuroanatomy 10, 62 (2016).

# Electronic Medical Records (EMRs)

 EMRs are collections of partially-structured records of medical data.

# • Data come in different formats:

- demographic data
- time series
- imaging
- plain text
- etc.
- Examples:
  - MIMIC-III: deidentified health data associated with
    - $\sim$  40,000 critical care patients https://mimic.physionet.org/
  - University of Florida Health Integrated Data Repository

#### **Selected Books**



Abello, James, Panos M. Pardalos, and Mauricio GC Resende, eds. *Handbook of massive data sets.* Vol. 4. Springer, 2013.



Pardalos, Panos M., Vladimir L. Boginski, and Vazacopoulos Alkis, eds. *Data mining in biomedicine*.Vol. 7. Springer Science & Business Media, 2008.



Pardalos, Panos M., Petros Xanthopoulos, and Michalis Zervakis, eds. *Data Mining for Biomarker Discovery.* Vol. 65. Springer Science & Business Media, 2012.









# Acute Kidney Injury and Sepsis

- sCr and 90-Day Mortality After Surgery
- Risk Prediction Models for Sepsis and AKI
- Multivariate Time Series Analysis for AKI Prediction



- Acute Kidney Injury (AKI) is one of common complications in post-operative patients.
- The in-hospital mortality rate for patients with AKI may be as high as 60%.
- An elevated serum creatinine (sCr) level in blood is commonly recognized as an indicator of AKI.
- The degree to which patterns of sCr change are associated with in-hospital mortality is unknown

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### Causes of AKI



Prerenal. Sudden and severe drop in blood pressure (shock) or interruption of blood flow to the kidneys from severe injury or illness.

- Intrarenal. Direct damage to the kidneys by inflammation, toxins, drugs, infection, or reduced blood supply.
- Postrenal. sudden obstruction of urine flow due to enlarged prostate, kidney stones, bladder tumor, or injury.

# Sepsis

- Sepsis is a severe and dysregulated inflammatory response to infection characterized by end-organ dysfunction distant from the primary site of infection.
- Development of (AKI) during sepsis:
  - increases patient morbidity,
  - predicts higher mortality,
  - has a significant effect on multiple organ functions,
  - is associated with an increased length of stay in the intensive care unit,
  - and hence consumes considerable healthcare resources
- [1] Zarjou, A., & Agarwal, A. (2011). Sepsis and acute kidney injury. Journal of the American Society of Nephrology, 22(6), 999-1006.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction



# 2 Acute Kidney Injury and Sepsis

- sCr and 90-Day Mortality After Surgery
- Risk Prediction Models for Sepsis and AKI
- Multivariate Time Series Analysis for AKI Prediction

# **Study Objective**

Objectives:

- To develop a comprehensive model for assessing the mortality risk in post-operative patients
- To establish a quantitative relationship between sCr pattern and mortality risk

Results:

- A probabilistic model was designed to assess mortality risk in post-operative patients
- The model provided high discriminative capacity and accuracy
- A quantitative association between sCr time series and mortality risk was established
- Simple and informative sCr risk factors were derived



- We performed a retrospective study involving patients admitted to Shands Hospital (Gainesville, FL, USA) between January 1, 2000 and November 30, 2010.
- For each patient who underwent a surgery, detailed clinical and outcome data were collected.
- The analysis included 46,299 patients.

# **Other Covariates**

- emergent surgery status
- type of surgical procedure
- age
- race
- gender
- intensive care unit (ICU) admission
- Charlson-Deyo comorbidity index (CCI)
- renal replacement therapy (RRT)
- time of surgery (days elapsed between admission and the surgical operation)
- operating surgeons' unique identifier

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

# **Two Cohorts of Patients**

- A subcohort of **7,766** patients with either
  - history of chronic kidney disease (CKD) prior to admission or
  - with the baseline estimated glomerular filtration rate (eGFR) of less than 60 ml/min/1.73 m2.

Since these patients may have **different sCr kinetics**, we performed a separate analysis for this cohort.

The remaining cohort consists of 38,533 patients.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

# **Two Cohorts Comparison**

#### Table 1

Baseline patients' characteristics stratified by ninety-day mortality.

	Baseline eGFR<60 (n=7,	0 ml/min/1.73 m <sup>2</sup> 766)	Baseline eGFR≥60 ml/min/1.73 m <sup>2</sup> (n=38,533)		
	Non-Survivors (n=848)	Survivors (n=6,918)	Non-Survivors (n= 1,597)	Survivors (n= 36,936)	
Patients' Characteristics					
Age in years, median (25th, 75th)	73 (63, 80)	67 (55, 75) <sup>a</sup>	64 (52, 74)	54 (41, 65) <sup>a</sup>	
Female gender, n (%)	360 (42)	3135 (45)	700 (44)	18,509 (50) <sup>a</sup>	

# Feature Extraction From sCr Time Series

- Model0: preoperative data only (no sCr time series).
   Used as a baseline.
- Model1: preoperative data and
  - the minimum of the sCr values available within six months of admission including the admission day value (BaseCr),
  - the absolute values for the maximum sCr (MaxCr), and
  - the last measured sCr during index hospitalization (LastCr).
- Model2: preoperative data and
  - BaseCr,
  - MaxCr BaseCr,
  - LastCr BaseCr.
- Model3: preoperative data and
  - BaseCr,
  - MaxCr/BaseCr,
  - LastCr/BaseCr.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

# **Mortality Risk**

### Probabilistic score

$$\log \frac{P(C=1|X=x)}{P(C=0|X=x)} = \alpha + \sum_{i=1}^{m} f_i(x_i),$$
(1)

 $C \in \{0, 1\}$  is an outcome,  $X = (X_1, \ldots, X_m)$  - the risk factors,  $x = (x_1, \ldots, x_m)$  - the values of these factors,  $f_i$  - is a nonlinear risk function associated with risk factor *i*, and  $\alpha$  is a free term.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### **Nonlinear Risk Factor Estimation**

From the probabilistic score equation we can derive

$$\alpha + \sum_{i=1}^{m} f_i(x_i) = \log \frac{P(X = x | C = 1)}{P(X = x | C = 0)} + \log \frac{P(C = 1)}{P(C = 0)}.$$
 (2)

Under assumption of risk factors independence:

$$\alpha + \sum_{i=1}^{m} f_i(x_i) = \sum_{i=1}^{m} \log \frac{P(X_i = x_i | C = 1)}{P(X_i = x_i | C = 0)} + \log \frac{P(C = 1)}{P(C = 0)}.$$
 (3)

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### **Nonlinear Risk Factor Estimation**

$$P(X_i|C=c) = \frac{\#\{j: C^j = c, x_i^j = x\}}{\#\{j: C^j = c\}}$$
(4)

For continuous factors,  $P(X_i = x_i | C)$  implies  $P(X_i \in [x_i - \epsilon_i, x_i + \epsilon_i])$ , where  $\epsilon_i$ 's are chosen appropriately.

- Nonlinear risk functions f<sub>i</sub> were estimated with cubic splines via a local scoring algorithm.
- This back fitting algorithm uses an iterative method which approximates the risk function, where the change in residuals is less than a prespecified threshold.
- The degrees of freedom for each spline were estimated by maximizing a restricted likelihood function.
- From the initial data set the predetermined set of variables was used for backward-stepwise selection based on AIC criteria.

### Degrees of freedom for f<sub>i</sub>'s after backward-stepwise selection

	Baseline eGFR≥60 ml/min/1.73 m <sup>2</sup>			Baseline eGFR<60 ml/min/1.73 m <sup>2</sup>				
Variables	Model 0 <sup>a</sup>	Model 1 <sup>b</sup>	Model 2 <sup>c</sup>	Model 3 <sup>d</sup>	Model 0 <sup>a</sup>	Model 1 <sup>b</sup>	Model 2 <sup>c</sup>	Model 3 <sup>d</sup>
Categorical or Nominal								
Type of surgical procedure	8.4	8.6	8.6	8.5	8.2	1.6	1.6	1.2
Emergent surgery	1	1	1	1	1	1	1	1
Gender	1	1	1	1	1	1	1	1
Race	1	1	1	1	1	1	1	1
Operating surgeon	1.8	1.4	1.5	1.7	1	1.5	1.9	1.6
Intensive care unit admission	1	1	1	1	1	1	1	1
Charlson comorbidity index	5.2	4.9	5.1	5.2	4.5	5.1	1	1
Time of surgery	х	Х	Х	Х	Х	Х	Х	х
Renal Replacement therapy	NA	NA	NA	NA	NA	NA	NA	NA
Continuous								
Age	3.8	3.9	3.9	3.9	3.6	4.1	4.0	3.75
BaseCr		1.7	2.4	Х		3.4	1	х
MaxCr		6				4.3		
LastCr		7.1				3.8		
MaxCr - BaseCr			5				4.2	
LastCr - BaseCr			7				4.1	
MaxCr/BaseCr				5.4				2.5
LastCr/BaseCr				7				5.9

The cohort excludes patients requiring renal replacement therapy. All values represent degree of freedom except for variables with

the symbol "X" that were discarded during the feature selection procedure.

Risk factors with estimated degrees of freedom close to 1 were not smoothed in a final model, instead the original values of risk factors were used. Therefore, the final model has the following form:

$$\log \frac{P(C = 1 | X = x)}{P(C = 0 | X = x)} = \alpha + \sum_{i \in I} w_i x_i + \sum_{i \notin I} f_i(x_i), \quad (5)$$

where *I* is a set of risk factors with estimated degrees of freedom close to 1.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

# **Classification Accuracy Evaluation**

- 70/30 cross-validation analysis has been performed
- Random split into 70% training subset and 30% testing subset
- The average over 100 runs results are reported

Acute Kidney Injury and Sepsis

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### **Resulting Non-linear Risk Functions**



Figure: Nonlinear functions for the association between maxCr/baseCr, llastCr/baseCr, age, and CCI and ninety-day mortality for the patients with chronic kidney disease or baseline estimated glomerular filtration rate < 60 ml/min/1.73 m2 after excluding patients on renal replacement therapy. Panels (A) and (B) show unadjusted and adjusted nonlinear functions, respectively. The Y axis represents risk probability for ninety-days mortality ranging from 0 to 1. The shaded areas represent 95% prediction intervals for the function values.

Acute Kidney Injury and Sepsis

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### **Resulting Non-linear Risk Functions**



**Figure:** Nonlinear functions for the association between maxCr/baseCr, lastCr/baseCr, age, and CCI and ninety-day mortality for the patients with baseline EGFT  $\geq 60$  ml/min/1.73 m2 after excluding patients on renal replacement therapy. Panels (A) and (B) show unadjusted and adjusted nonlinear functions, respectively. The Y axis represents risk probability for ninety-days mortality ranging from 0 to 1. The shaded areas represent 95% prediction intervals for the function values.

Introduction Acute Kidney Injury and Sepsis Multivariate Time Series Analysis for AKI Predi

# Association Between pattern in sCr change and 90-day mortality



Figure: Association between pattern in sCr change and 90-day mortality adjusted for preoperative clinical factors among (Å) patients with baseline estimated glomerular filtration rate  $\geq 60$  ml/min/1.73 m2 and (B) patients with chronic kidney disease or baseline estimated glomerular filtration rate < 60 ml/min/1.73 m2. Left panel: The log odds ratios for 90-day mortality based on adjusted nonlinear functions for maxCr/baseCr and lastCr/baseCr with respect to a pattern with no change in sCr. Right panel: Distribution of maxCr/baseCr and lastCr/baseCr values for non-survivors and survivors. \*There were no patients in the dataset with the combination of the maxCr/baseCr and lastCr/baseCr and lastCr/baseCr
sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### References

- Korenkevych, D., Ozrazgat-Baslanti, T., Thottakkara, P., Hobson, C. E., Pardalos, P., Momcilovic, P., & Bihorac, A. (2016). The pattern of longitudinal change in serum creatinine and 90-day mortality after major surgery. Annals of Surgery, 263(6), 1219-1227.
- [2] Rosenthal, M., Korenkevych, D., Baslanti, T. O., Webel, A., Glerum, A., Momcilovic, P., ... & Bihorac, A. (2013). 703: Algorithms can accurately predict risk for major postoperative complications using preoperative data. Critical Care Medicine, 41(12), A173.
- [3] Bihorac, A., Korenkevych, D., Baslanti, T. O., Momcilovic, P., Pardalos, P., Segal, M., & Moore, F. (2013). 714: Database Communication Enables Machine Learning Classifiers to Predict Postoperative AKI in ICU. Critical Care Medicine, 41(12), A176.
- [4] Baslanti, T. O., Korenkevych, D., Momcilovic, P., Pardalos, P., Hobson, C., & Bihorac, A. (2012). Mathematical modeling of the association between the pattern of change in postoperative serum creatinine and hospital mortality. In CRITICAL CARE MEDICINE (Vol. 40, No. 12, pp. U131-U131). 530 WALNUT ST, PHILADELPHIA, PA 19106-3621 USA: LIPPINCOTT WILLIAMS & WILKINS.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction



## 2 Acute Kidney Injury and Sepsis

- sCr and 90-Day Mortality After Surgery
- Risk Prediction Models for Sepsis and AKI
- Multivariate Time Series Analysis for AKI Prediction

## **Study Outline**

## Objective:

 To compare performance of risk prediction models for forecasting postoperative sepsis and acute kidney injury.

## Design:

- Retrospective single center cohort study of adult surgical patients admitted between 2000 and 2010.
- 50,318 adult patients undergoing major surgery.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Study Outline**

# Algorithms:

- Logistic Regression (LR)
- Generalized Additive Models (GAMs)
- Naïve Bayes (NB) and
- Support Vector Machines (SVMs)

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Development Flow from Raw Data to Model Building**



#### PM Pardalos Data Science in Biomedicine

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### **Model Comparison**

Model		Acute Kidney Injury		Severe Sepsis						
	Accuracy (95% CI)	AUC (95% CI)	PPV (95% CI)	Accuracy (95% Cl)	AUC (95% CI)	PPV (95% CI)				
Logistic Regression Model	0.752 (0.746,0.758)	0.824 (0.818,0.828) <sup>b</sup>	0.725 (0.714,0.737)	0.773 (0.762,0.781)	0.851 (0.840,0.8560)	0.811 (0.785,0.833)				
GAMs	0.756 (0.751,0.761)	0.827 (0.821,0.832) <sup>a</sup>	0.719 (0.706,0.729)	0.775 (0.766,0.783)	0.852 (0.840,0.863)	0.806 (0.779,0.832)				
Naïve Bayes Model	0.744 (0.738,0.749)	0.797 (0.791,0.803) <sup>a,b</sup>	0.545 (0.534,0.558)	0.805 (0.798,0.811)	0.83 (0.819,0.841) <sup>a,b</sup>	0.689 (0.659,0.716)				
SVM	0.767 (0.757,0.774)	0.819 (0.811,0.828) <sup>a,b</sup>	0.662 (0.648,0.676)	0.71 (0.689,0.731)	0.762 (0.733,0.782) <sup>a,b</sup>	0.677 (0.619,0.722)				
After feature selection with LASSO										
Logistic Regression Model	0.753 (0.747,0.757)	0.824 (0.818,0.830) <sup>b</sup>	0.726 (0.714,0.738)	0.772 (0.760,0.780)	0.85 (0.838,0.863) <sup>b</sup>	0.812 (0.781,0.838)				
GAMs	0.757 (0.752,0.762)	0.828 (0.822,0.833) <sup>a</sup>	0.72 (0.706,0.732)	0.774 (0.766,0.780)	0.851 (0.842,0.862)	0.806 (0.783,0.831)				
Naïve Bayes Model	0.744 (0.737,0.750)	0.797 (0.789,0.804) <sup>a,b</sup>	0.545 (0.533,0.556)	0.806 (0.800,0.813)	0.831 (0.817,0.841) <sup>a,b</sup>	0.69 (0.659,0.711)				
SVM	0.767 (0.759,0.774)	0.82 (0.812,0.829) <sup>a,b</sup>	0.665 (0.646,0.685)	0.697 (0.684,0.713)	0.757 (0.736,0.779) <sup>a,b</sup>	0.689 (0.652,0.732)				
After feature extraction	with 5 principal com	oonents								
Logistic Regression Model	0.774 (0.769,0.781)	0.853 (0.849,0.859) <sup>a,b</sup>	0.758 (0.746,0.767)	0.818 (0.809,0.824)	0.904 (0.895,0.913) <sup>a,b</sup>	0.854 (0.841,0.880)				
GAMs	0.773 (0.768,0.777)	0.858 (0.853,0.862) <sup>a,b</sup>	0.784 (0.771,0.793)	0.826 (0.819,0.833)	0.909 (0.902,0.917) <sup>a,b</sup>	0.86 (0.843,0.878)				
Naïve Bayes Model	0.741 (0.735,0.747)	0.819 (0.814,0.826) <sup>a,b</sup>	0.666 (0.651,0.677)	0.805 (0.797,0.815)	0.882 (0.874,0.890) <sup>a,b</sup>	0.839 (0.822,0.866)				
SVM	0.777 (0.767,0.782)	0.857 (0.850,0.862) <sup>a,b</sup>	0.735 (0.725,0.750)	0.85 (0.737,0.897)	0.877 (0.828,0.904) <sup>a,b</sup>	0.751 (0.667,0.850)				

<sup>a</sup> p<0.05 for AUC comparison with respect to logistic regression model without any data reduction.

<sup>b</sup> p<0.05 for AUC comparison with respect to GAMs model without any data reduction.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### **GAMs: Predicted Risk Functions**



Figure: Predicted risk functions for the association between (A) acute kidney injury and (B) severe sepsis and continuous variables. Risk functions were generated from multivariate generalized additive models and logistic regression models. GAM, generalized additive model; DoF, degree of freedom; GFR, glomerular filtration rate.

## Results

- AUC ROC for different models ranged between 0.797 and 0.858 for AKI and between 0.757 and 0.909 for severe sepsis.
- LR, GAM, and SVM had better performance compared to NB model
- GAMs additionally accounted for non-linearity of continuous clinical variables as depicted in their risk patterns plots
- Reducing the input feature space with LASSO had minimal effect on prediction performance, while feature extraction using principal component analysis improved performance of the models.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### References

[1] Thottakkara, P., Ozrazgat-Baslanti, T., Hupf, B. B., Rashidi, P., Pardalos, P., Momcilovic, P., & Bihorac, A. (2016). Application of machine learning techniques to high-dimensional clinical data to forecast postoperative complications. PloS one, 11(5), e0155705.



## Acute Kidney Injury and Sepsis

- sCr and 90-Day Mortality After Surgery
- Risk Prediction Models for Sepsis and AKI
- Multivariate Time Series Analysis for AKI Prediction

## **Objectives**

## Objectives:

- Utilize information hidden in multivariate time series (heart rate, blood pressure, etc.) of physiological recordings collected during a surgical procedure.
- A faster algorithm for pattern mining form multivariate time series for on-line applications.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Augmented AKI Data Set**

- For 5202 patients, intra-operative recordings were available
- We decided to concentrate on two time series:
  - Blood Pressure
  - Heart Rate

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## Intersection of two Problems

- Time Series Classification
  - Data is a collection of records, where each record consists of a number of time series and an outcome associated with it
  - Time series are sampled unevenly in time
  - Time series can be both numerical and categorical
- Pattern Mining as Class-Specific Feature Extraction

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Mining Frequent Patterns**



sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## From Time Series to Temporal Abstractions



From Temporal Abstractions to Multivariate State Sequences



## **Multivariate State Sequences: Definitions**

- S = (F, V) is a state where is F is a variable label and V ∈ Σ is an abstraction value.
- *E* = (*F*, *V*, *s*, *e*) is a **state interval** where (*F*, *V*) is a state, and *s* and *e* are the start and end times of the state interval.
- Z = ⟨E<sub>1</sub>,..., E<sub>l</sub>⟩, E<sub>i</sub>.s ≤ E<sub>i+1</sub>.s, 1 ≤ i ≤ l − 1, is a Multivariate State Sequence (MSS) consisting of *l* state intervals which are arranged in the non-decreasing order of their start times.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### **Temporal Patterns**



**Fig. 4** A temporal pattern with states  $\langle (F_1, B), (F_3, A), (F_2, C), (F_3, B) \rangle$  and temporal relations  $R_{1,2} = c$ ,  $R_{1,3} = b$ ,  $R_{1,4} = b$ ,  $R_{2,3} = c$ ,  $R_{2,4} = b$  and  $R_{3,4} = c$ 

## **Temporal Pattern Definition**

- $P = (\langle S_1, \ldots, S_k \rangle, R)$  is a **Temporal Pattern** (TP) of length  $k \ (|P| = k)$  with states  $S_1, \ldots, S_k$ , where R is a (upper-triangular) matrix describing pair-wise temporal relationships between the states:  $R_{i,j} \in \{b, c\}, \ 1 \le i < j \le k$ .
- Given an MSS Z = ⟨E<sub>1</sub>, E<sub>2</sub>,..., E<sub>I</sub>⟩ and a temporal pattern P = (⟨S<sub>1</sub>,..., S<sub>k</sub>⟩, R), we say that Z contains P, denoted as P ∈ Z, if there is an injective mapping π : {1,..., k} → {1,..., I} (k ≤ I) that matches every state S<sub>i</sub> in P to a state interval E<sub>π(i)</sub> in Z such that:

**●** 
$$S_i.F = E_{\pi(i)}.F$$
 and  $S_i.V = E_{\pi(i)}.V$ ,  $1 \le i \le k$ ,  
**●**  $\pi(i) < \pi(j), i < j$ ,  
**●**  $R(E_{\pi(i)}, E_{\pi(j)}) = R_{i,j}, i < j$ .

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### **Temporal Patterns in MSSs**







sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### **Frequent Temporal Patterns**



PM Pardalos Data Science in Biomedicine

## Frequent Temporal Pattern Mining (FTPM)

```
Input: D, 1-FTPs
Output: FTPs
FTPs \leftarrow 1-FTPs:
k \leftarrow 1;
while |k-FTPs| > 0 and no other stopping criteria are met do
   (k+1)-FTPs \leftarrow \emptyset:
   (k+1)-candidates \leftarrow CreateCandidates(k-FTPs, 1-FTPs);
   forall the TP \in (k+1)-candidates do
       if TP is FTP in D then
        (k+1)-FTPs \leftarrow (k+1)-FTPs \cup \{TP\};
       end
   end
   FTPs \leftarrow FTPs \cup (k+1) - FTPs;
   k \leftarrow k+1:
end
```

 $\label{eq:algorithm:thegeneral} \textbf{Algorithm:thegeneral framework}.$ 

### References

- Batal, I. et al. An efficient pattern mining approach for event detection in multivariate temporal data. Knowl Inf Syst 46, 115–150 (2016).
- [2] Batal, I. et al. A Pattern Mining Approach for Classifying Multivariate Temporal Data. Proceedings (IEEE Int Conf Bioinformatics Biomed), 358–365 (2011).
- [3] Batal, I. et al. Multivariate time series classification with temporal abstractions. Proceedings of the Twenty-Second International FLAIRS Conference, 344-349 (2009).

Introduction Acute Kidney Injury and Sepsis Multivariate Time Series Analysis for AKI Prediction

## **Apriory Property**

## Apriory Property:

$$P.p\_ids = \bigcap_{\widetilde{P} \in sub(P)} \widetilde{P}.ids = \bigcap_{\widetilde{P} \in sub_k(P)} \widetilde{P}.ids$$



PM Pardalos Data Science in Biomedicine

Frequent Temporal Pattern Mining with Extended Lists (FTPMwEL)

Main Idea:

 Store locations where the first state of the pattern appears inside each record

Definitions:

- Let *P.pos*[*i*] and *P.p\_pos*[*i*] denote all **positions** and **potential positions** at which *P* appears and may appear inside *Z<sub>i</sub>*, respectfully.
- Thus, *P.ppos*[*i*] can be computed as the following:

$$P.p\_pos[i] = \bigcap_{\widetilde{P} \in X} \widetilde{P}.pos[i],$$

where  $X = \operatorname{sub}_k(P) \setminus \operatorname{parent}(P)$ .

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction





sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction



sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction



sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Extended Lists: Example**



PM Pardalos Data Science in Biomedicine

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction



sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### **Extended Lists: Example**



<0,3>

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### **Extended Lists**



PM Pardalos Data Science in Biomedicine

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Speeding-up Search**



PM Pardalos Data Science in Biomedicine

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### **FTPMwEL Algorithm**

Input: D. 1-FTPs Output: FTPs  $FTPs \leftarrow 1$ -FTPs:  $k \leftarrow 1$ : while |k-FTPs| > 0 and no other stopping criteria are met do (k + 1)-FTPs  $\leftarrow \emptyset$ ; (k + 1)-candidates  $\leftarrow$  CreateCoherentCandidates(k-FTPs, 1-FTPs); forall the  $P \in (k + 1)$ -candidates do if  $\exists \tilde{P} \in sub_k(P)$  such that  $\tilde{P}$  isn't frequent then continue end  $P.p_ids \leftarrow \cap_{\overline{P}Gsub_k(P)} \overline{P}.ids;$ if not PotentiallyFrequent(P) then 1 continue end for all the  $id \in P.p.ids$  do  $P.p.pos[id] = \bigcap_{\overline{P} \in subs_i(P) \setminus parent(P)} \widetilde{P}.pos[id];$ if  $P.p.pos[id] = \emptyset$  then  $P.p_ids \leftarrow P.p_ids \setminus id$ end else  $P.p\_ex\_list[id] \leftarrow CreateLinks(P.p\_pos[id], P.parent.pos[id]);$ if  $P.p\_ex\_list[id] = \emptyset$  then  $P.p_ids \leftarrow P.p_ids \setminus id$ end end end if not PotentiallyFrequent(P) then 1 continue end  $P.ids \leftarrow P.p\_ids$ ; for all the  $id \in P.ids$  do  $P.ex\_list[id] \leftarrow FindPositionsAndLinks(P.p.ex\_list[id]);$ if  $P.ex_{ist}[id] = \emptyset$  then  $P.ids \leftarrow P.ids \setminus id$ end end if P is FTP in D then (k+1)-FTPs  $\leftarrow$  (k+1)-FTPs  $\cup$   $\{P\}$ ; end end  $FTPs \leftarrow FTPs \cup (k + 1) - FTPs$ ;  $k \leftarrow k+1$ ; end Algorithm 2: Frequent Temporal Pattern Mining with Extended Lists

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Mining FTPs on AKI dataset**

$\theta$	max	found	FTPM	FTPM,	Speed-	FTPM	FTPM,	Memory
	k	k	wEL, sec	sec	up	wEL,	MB	Ratio
					Ratio	MB		
0.5	5	5	5013.54	22768.01	4.54	14300.29	398.83	112.66
0.6	5	5	2097.22	11146.29	5.31	5105.14	126.93	40.22
0.7	Inf	7*(13)	8156.93	> 86400	> 10.59	179.8	NA	NA
0.8	Inf	10	569.16	4404.06	7.74	17.24	1.26	13.7
0.9	Inf	6	36.64	82.05	2.24	11.42	1.68	6.78

Table 1: Computational time and memory usage comparison of FTPMwEL and FTPM on the AKI dataset for a varying threshold level.

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

## **Mining FTPs on AKI dataset**


Introduction

Acute Kidney Injury and Sepsis

sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

### Mining FTPs on datasets from UCR archive

	Dataset	FTPM	FTPM.	Speed-	FTPM	FTPM.	Memory
		wEL sec	800	up Ratio	wEL.	MB	Ratio
			000	up runo	MB		. and
_							
1	BeetleFlv	114.13	120.11	1.05	209.57	14.88	13.08
2	BirdChicken	38.38	53.37	1.39	105.48	5.89	16.91
3	Coffee	42.09	46.84	1.11	86,98	1.17	73.47
4	Computers	299.14	672.23	2.25	175.04	0.52	333.41
5	DistalPhalanx	13.16	24.34	1.85	39.91	0.52	75.24
	OutlineCorrect						
6	Earthquakes	1661.45	666.89	-2.49	993.54	7.2	137.01
7	ECG200	58.06	115.58	1.99	128.51	3.46	36.09
8	ECGFiveDays	13.44	15.87	1.18	33.62	0.53	62.76
9	FordA	26578.2	21167.66	-1.26	23044.7	401.63	56.38
10	FordB	15227.69	12381.59	-1.23	13714.19	250.28	53.8
11	Gun_Point	1.25	2.28	1.82	0.07	0	NA
12	Ham	1001.84	1104.86	1.1	1318.66	44.4	28.7
13	HandOutlines	42.47	38.55	-1.1	116.25	4.2	26.71
14	Herring	92.62	164.44	1.78	286.96	11.68	23.56
15	ItalyPowerDemand	1.28	2.33	1.82	1.51	0	NA
16	Lighting2	737.92	2329.12	3.16	608.77	13.39	44.46
17	MiddlePhalanx	101.23	120.23	1.19	304.71	9.5	31.07
	OutlineCorrect						
18	MoteStrain 2.48	3.61	1.46	2.28		0	NA
19	PhalangesOutlines	95.14	159.27	1.67	405.23	14.16	27.62
	Correct						
20	ProximalPhalanx	42.34	63.23	1.49	174.46	4.72	35.97
	OutlineCorrect						
21	ShapeletSim	1584.8	9980.11	6.3	878.67	18.49	46.52
22	SonyAIBORobot Surface	49.41	74.31	1.5	145.55	10.27	13.18
23	SonyAIBORobot	99.64	135.33	1.36	236.55	14.56	15.24
94	Strawborry	01.41	127.24	1.8	221.09	10.90	21.18
24	TooSormontation1	108.99	979.45	1.3	240.24	15.05	20.22
20	ToeSegmentation?	173.06	235.05	1.37	391.07	17.63	17.97
20	TwoLeadECG	3.95	4 19	1.00	3.88	0	NA
28	wafer	215.3	494.89	2.3	881.46	35.67	23.71
20	Wine	31.67	32.01	1.04	65.16	3.26	18.98
30	WormsTwoClass	854.77	2412.22	2.82	930.53	30.74	20.20
31	Voga	146.48	210.16	1.43	445 57	14.67	29.38
- 01	1080	110/110	210.10	1110	440101	4,4101	20100

PM Pardalos

Introduction Acute Kidney Injury and Sepsis sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### Other results

- The best pattern has accuracy of 76.8% alone
- FTPs are almost the same after sampling of records

## Project participants:

- Anton Kocheturov
- Petar Momcilovic
- Azra Bihorac
- Dmytro Korenkevych
- Panos M Pardalos

## Acknowledgments:

AK was supported by the grant by UF Informatics Institute. AB, PP and PM were supported by grant R01 GM-110240 by the National Institute of General Medical Sciences - National Institutes of Health. Introduction Acute Kidney Injury and Sepsis sCr and 90-Day Mortality After Surgery Risk Prediction Models for Sepsis and AKI Multivariate Time Series Analysis for AKI Prediction

#### The End

# Thank you!

PM Pardalos Data Science in Biomedicine