Real-Time Data Mining

Data Mining for the XXI Century

João Gama jgama@fep.up.pt

INESC TEC, University of Porto, Portugal ISEP 2020





э

・ロト ・四ト ・ヨト ・ヨト

Motivation

Case Study

Clustering Time Series Growing the Structure Adapting to Change Properties of ODAC

Final Comments



Outline

Motivation

Case Study

Clustering Time Series

Growing the Structure Adapting to Change Properties of ODAC

Final Comments



Context





ヘロト ヘ週ト ヘヨト ヘヨト

Industry 4.0

We have machines that collect, process, and send information to other machines





э

We have machines that collect, process, and send information to other machines





Internet of Things

HOT ANALYTICS			Q3/2016 Insights			that empower you to understand IoT markets				
IoT Segment			Global share of IoT projects ¹			Details				
							Americas	Europe	APAC	Trend ²
1	$\underline{\mathbf{A}}$	Connected Industry				22%	43%	30%	20%	$\langle \! $
2	ьđ	Smart City				20%	31%	47%	15%	1
3	۲	Smart Energy			13%		49%	24%	25%	\bigtriangledown
4	-	Connected Car			13%		43%	33%	17%	$\langle \! \rangle$
5	<	Other	89	%			46%	33%	13%	$\overline{\langle}$
6	97	Smart Agriculture	6%				48%	31%	17%	٨
7	Ē	Connected Building ³	5%			35	48%	33%	12%	1
8	Ô	Connected Health	5%		- 🍾 🔨		61%	30%	6%	\bigcirc
ٞ۞)	Smart Retail	4%		K N = 640 global, pub announced loT proj	licly	52%	30%	13%	۲
10		Smart Supply Chain	4%		Americas Europe APAC	MEA N/A	57%	35%	4%	Ø
1. Based on 640+ publicly known enterprise IoT projects. (Not including consumer IoT projects e.g., Weanables, Smart Home) 2. Trend based on IoT Analytics's 02/2016 IoT Employment Statistics Tracker 3. Not including Consumer Smart Home Solutions Source: IoT Analytics 2016 Global overview of 640 enterprise IoT we cares (August 2016)										



◆□▶ ◆□▶ ◆目▶ ◆目▶ ─ 目

The Growth of Digital Data...

The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020



Source: IDC's Digital Universe Study, sponsored by EMC, December 2012

Memory unit	Size	Binary size
kilobyte (kB/KB)	10 ³	2 ¹⁰
megabyte (MB)	10 ⁶	2 ²⁰
gigabyte (GB)	10 ⁹	2 ³⁰
terabyte (TB)	10 ¹²	2 ⁴⁰
petabyte (PB)	10 ¹⁵	2 ⁵⁰
exabyte (EB)	10 ¹⁸	2 ⁶⁰
zettabyte (ZB)	10 ²¹	2 ⁷⁰
yottabyte (YB)	10 ²⁴	2 ⁸⁰

Tools: time ago ...

Tools seemed quite powerful





Tools

Problems

・ロト ・ 同ト ・ ヨト ・ ヨト

э

Tools: nowadays ...

Last few years







The Value of Information ...



Elapsed Time



<ロト <回ト < 三ト < 三ト = 三

A World in Movement

- The new characteristics of data:
 - Time and space: The objects of analysis exist in time and space. Often they are able to move.
 - Dynamic environment: The objects exist in a dynamic and evolving environment.
 - Information processing capability: The objects have limited information processing capabilities
 - Locality: The objects know only their local spatio-temporal environment;
 - Distributed Environment: Objects will be able to exchange information with other objects.
- Main Goal:
 - Real-Time Analysis: decision models have to evolve in correspondence with the evolving environment.



These characteristics imply:

- Switch from one-shot learning to continuously learning dynamic models that evolve over time.
- In this context, finite training sets, static models, and stationary distributions will have to be completely thought anew.
- Computational resources are finite. Algorithms will have to use *limited computational resources* (in terms of computations, memory, space and time, communications).



Data Stream Computational Model

- 1. One-pass algorithms: random access to data has high cost
- 2. Limited computational resources:

time, memory, bandwidth

3. Anytime prediction



(日)

э

Outline

Motivation

Case Study

Clustering Time Series

Growing the Structure Adapting to Change Properties of ODAC

Final Comments



Scenario



Electrical power Network: Sensors all around network monitor measurements of interest.



Sensors produce continuous flow of data at high speed:

- Send information at different time scales;
- Act in adversary conditions: they are prone to noise, weather conditions, battery conditions, etc;
- Huge number of Sensors, variable along time
- Geographic distribution:
 - The topology of the network and the position of the sensors are known.



Illustrative Learning Tasks:

Cluster Analysis

Identification of Profiles: Urban, Rural, Industrial, etc.

Predictive Analysis

- Predict the value measured by each sensor for different time horizons.
- Prediction of peaks on the demand.
- Monitoring Evolution
 - Change Detection
 - Detect changes in the behavior of sensors;
 - Detect Failures and Abnormal Activities;
 - Extreme Values, Anomalies and Outliers Detection
 - Identification of critical points in load evolution;



This problem has been addressed time ago:

Strategy

- Select a finite sample
- Generate a static model (cluster structure, neural nets, Kalman filters, Wavelets, etc)

- Very good performance in next month!
- Six months later: Retrain everything!

This problem has been addressed time ago:

Strategy

- Select a finite sample
- Generate a static model (cluster structure, neural nets, Kalman filters, Wavelets, etc)
- Very good performance in next month!
- Six months later: Retrain everything!

What is the Problem?

The world is not static! Things change over time.



The Data Stream Phenomenon

Highly detailed, automatic, rapid data feeds.

- Internet: traffic logs, user queries, email, financial,
- Telecommunications: phone calls, sms,
- Astronomical surveys: optical, radio,.
- Sensor networks: many more observation points ...
- Most of these data will never be seen by a human!
- Need for near-real time analysis of data feeds.
- Monitoring, intrusion, anomalous activity Classification, Prediction, Complex correlations, Detect outliers, extreme events, etc



Continuous flow of data generated at **high-speed** in **Dynamic**, **Time-changing** environments.

The usual approaches for *querying*, *clustering* and *prediction* use **batch procedures** cannot cope with this streaming setting. Machine Learning algorithms assume:

Instances are independent and generated at random according to some probability distribution D.

► It is required that *D* is stationary

Practice: *finite* training sets, *static* models.

We need to maintain **Decision models** in **real time**. Decision Models must be capable of:

- incorporating new information at the speed data arrives;
- detecting changes and adapting the decision models to the most recent information.
- forgetting outdated information;

Unbounded training sets, dynamic models.



Outline

Motivation

Case Study

Clustering Time Series

Growing the Structure Adapting to Change Properties of ODAC

Final Comments



Goal: Continuously maintain a clustering structure from evolving time series data streams.

- Ability to Incorporate new Information;
- Process new Information at the rate it is available.
- Ability to Detect and React to *changes* in the Cluster's Structure.

Clustering of *variables* (sensors) not examples! The standard technique of transposing the working-matrix does not work: transpose is a blocking operator!

Online Divisive-Agglomerative Clustering, Rodrigues & Gama, 2008 **Goal:** Continuously maintain a hierarchical cluster's structure from evolving time series data streams.

- Performs hierarchical clustering
- Continuously monitor the evolution of clusters' diameters
- Two Operators:
 - Splitting: expand the structure more data, more detailed clusters
 - Merge: contract the structure reacting to changes.
- Split and merge criteria are supported by a confidence level given by the Hoeffding bounds.



Feeding ODAC

Each example is processed once.

Only sufficient statistics at leaves are updated.

Sufficient Statistics: a triangular matrix of the correlations between variables in a leaf.

Released when a leaf expands to a node.



 $C_1 = \{ \, x_2 \, , \, x_3 \} \, , \, C_2 = \{ \, x_4 \, , \dots \, , \, x_{m-1} \} \, , \, C_3 = \{ \, x_1 \, , \, x_m \}$



◆日 > < 同 > < 国 > < 国 >

Distance between time Series: $rnomc(a, b) = \sqrt{\frac{1-corr(a,b)}{2}}$ where corr(a, b) is the Pearson Correlation coefficient: $corr(a, b) = \frac{P - \frac{AB}{n}}{\sqrt{A_2 - \frac{A^2}{n}}\sqrt{B_2 - \frac{B^2}{n}}}$

The *sufficient statistics* needed to compute the correlation are easily updated at each time step:

$$A = \sum a_i, \ B = \sum b_i, \ A_2 = \sum a_i^2, \ B_2 = \sum b_i^2, \ P = \sum a_i b_i$$





A 日 > A 四 > A 回 > A 回 >

The Splitting Operator: Expanding a Leaf





The base Idea

A small sample can often be enough to choose a near optimal decision

(*Mining High-Speed Data Streams*, P. Domingos, G. Hulten; KDD00)

- Collect sufficient statistics from a small set of examples
- Estimate the merit of each alternative

How large should be the sample?

- The wrong idea: Fixed sized, defined apriori without looking for the data;
- The right idea: Choose the sample size that allow to differentiate between the alternatives.



Expanding a leaf: How large should be the sample? Let

- $d_1 = d(a, b)$ the farthest distance
- ► *d*₂ the second farthest distance

Question:

Is d_1 a stable option?

what if we observe more examples?

Hoeffding bound:

Split if $d_1 - d_2 > \epsilon$ with $\epsilon = \sqrt{\frac{R^2 ln(1/\delta)}{2n}}$ where R is the range of the random variable; δ is a user confidence level, and n is the number of observed data points.



- Suppose we have made *n* independent observations of a random variable *r* whose range is *R*.
- The Hoeffding bound states that:
 - With probability 1δ
 - The true mean of r is in the range $\overline{r} \pm \epsilon$ where $\epsilon = \sqrt{\frac{R^2 ln(1/\delta)}{2n}}$
- Independent of the probability distribution generating the examples.



The Expand Operator: Expanding a Leaf





A multi-window system: each node (and leaves) receive examples from different time-windows.





The Merge Operator: Change Detection

Time Series Concept Drift:

- Time evolving time-series
- Changes in the distribution generating the observations.
- Clustering Concept Drift
 - Changes in the way time series correlate with each other
 - Change in the cluster Structure.





The Splitting Criteria guarantees that cluster's diameters monotonically decrease.

- Assume Clusters: c_j with descendants c_k and c_s .
- If diameter(c_k) − diameter(c_j) > ε OR diameter(c_s) − diameter(c_j) > ε
 - Change in the correlation structure!
 - Merge clusters c_k and c_s into c_j .





- 日本 - 4 国本 - 4 国本

The Electrical Load Demand Problem





< □ > < @ > < E > < E >

The Electrical Load Demand Problem





(ㅁ▶ ◀륨▶ ◀불▶ ◀불▶ ...

Properties of ODAC

- For stationary data the cluster's diameters monotonically decrease.
- Constant update time/memory consumption with respect to the number of examples!
- Every time a split is reported
 - the time to process the next example decreases, and
 - the space used by the new leaves is less than that used by the parent.



・ロト ・ 同ト ・ ヨト ・ ヨト

Evolution of Processing Speed





Hoeffding Algorithms

Classification:

Mining high-speed data streams, P. Domingos, G. Hulten, KDD, 2000

Regression:

Learning model trees from evolving data streams; Ikonomovska, Gama, Dzeroski; Data Min. Knowl. Discov. 2011

Decision Rules: Learning Decision Rules from Data Streams, J. Gama, P. Kosina; IJCAI 2011

Regression Rules E. Almeida, C. Ferreira, J. Gama: Adaptive Model Rules from Data Streams. ECML/PKDD 2013

Clustering: Hierarchical Clustering of Time-Series Data Streams. Rodrigues, Gama, IEEE TKDE 20(5): 615-627 (2008)

Multiple Models:

. . .

Ensembles of Restricted Hoeffding Trees. Bifet, Frank, Holmes, Pfahringer; ACM TIST: 2012

J. Duarte, J. Gama, Ensembles of Adaptive Model Rules from High-Speed Data Streams. BigMine 2014.



The number of examples required to expand a node only depends on the Hoeffding bound.

- Low variance models: Stable decisions with statistical support.
- Low overfiting:

Examples are processed only once.

- No need for pruning; Decisions with statistical support;
- Convergence: Hoeffding Algorithms becomes asymptotically close to that of a batch learner. The expected disagreement is δ/p; where p is the probability that an example fall into a leaf.



Outline

Motivation

Case Study

Clustering Time Series

Growing the Structure Adapting to Change Properties of ODAC

Final Comments



Massive Online Analysis

Config	gure Eval	JatePrequ	uential -l t	rees.Hoe	ffdingTree -s	generators.WaveformGenerato	r		Run	
command			stat	status		time elapsed	current activity	% complete		
valuate	Prequential -	l (trees.H	o	r	unning	10m11s	Evaluating learner	21,22		
valuate	Prequential -	l trees.Ho			unning	11m13s	Evaluating learner	12,25		
						Pause Resume Ca	Delete			
					Preview (1	1m13s) Refresh Auto re	fresh: every second 👻			
4982E	-7,890000	0.0,84	.6,76.9	036145	8431755,8	900000.0,-11239.0,3211	.0,1606.0,1606.0,24.0,0.	0,0.0,-Infinity		
4488E	-7,900000	0.0,83	.8,75.6	566322	904254,90	000000.0,-11330.0,3237.	0,1619.0,1619.0,25.0,0.0	,0.0,-Infinity		
7947E	-7,910000	0.0,86	.0,78.9	278489	5482129,9	9100000.0,-11505.0,3287	.0,1644.0,1644.0,25.0,0.	0,0.0,-Infinity		
7032E	-7,920000	0.0,86	.1,79.1	478522	2877956,9	200000.0,-11589.0,3311	.0,1656.0,1656.0,25.0,0.	0,0.0,-Infinity		
1544E	-7,930000	0.0,85	.399999	999999	99,78.113	360239611082,9300000.0,	-11757.0,3359.0,1680.0,1	680.0,25.0,0.0,0.0,-Infi	nity	
2432E	-7,940000	0.0,85	.0,77.4	774436	5982531,9	9400000.0,-11841.0,3383	.0,1692.0,1692.0,25.0,0.	0,0.0,-Infinity		
8526E	-7,950000	0.0,85	.3,77.8	983927	4706439,9	9500000.0,-11960.0,3417	.0,1709.0,1709.0,25.0,0.	0,0.0,-Infinity		
2083E	-7,960000	0.0,84	. 899999	999999	99,77.348	327846916366,9600000.0,	-12100.0,3457.0,1729.0,1	729.0,25.0,0.0,0.0,-infi	nity	
001E-	7,9700000	.0,85.	1,77.62	6/8//9	233454,91	100000.0,-12219.0,3491.	0,1746.0,1746.0,25.0,0.0	,0.0,-infinity		
`	6					Evport as tyt fil	0		,	
Evalua	tion					Export up reach				
Value	,					Plot				
	Measure Cu		Current Mean		an	Zoom in Y Zoom out Y				
۲	Accuracy	84,90	83,30	85,06	81,43				Juck	
0	Карра	77,30	74,91	77,57	72,13	00.00			*	
0	Ram-Hours	0.00	0.00	0.00	0.00	and a part of the	some margan a hula		=	
	Tara	670.60	401.07	220 70	227.17	Mode a weat out	which the proper a terreture	n. IulaMir		
0	nme	670,60	401,87	556,76	237,17	. MANNAN	and a merilian ware	with a		
O	Memory	0,01	0,02	0,01	0,01	Y¥' ''				
									-	
						Ч				

New Tools Emerge





(日)

A Generic Model for Adaptive Learning Algorithms



A generic schema for an online adaptive learning algorithm.

(A survey on concept drift adaptation, J.Gama et al, ACM-CSUR 2014)



э

Learning from data streams:

- Learning is not one-shot: is an evolving process;
- We need to monitor the learning process;
- Opens the possibility to reasoning about the learning

・ロト ・ 四ト ・ ヨト ・ ヨト ・ ヨ

- What changed in the decision structure last week?
- Which patterns disappeared/ appeared last week?
- Which patterns are growing/shrinking this month?

・ロト ・ 四ト ・ ヨト ・ ヨト ・ ヨ

Mine the evolution of decision structures.

Intelligent systems must:

- be able to adapt continuously to changing environmental conditions and evolving user habits and needs.
- be capable of predictive self-diagnosis.

The development of such self-configuring, self-optimizing, and self-repairing systems is a major scientific and engineering challenge.



Real-time learning: An existential pleasure!

Thank you!

