

# A GEO-stack for Big Data

Driving Spatial Analysis Beyond the Limits of Traditional Storage

Joana Simões<sup>1</sup>

<sup>1</sup>BDigital, CASA, CICS.NOVA

March 8, 2015





Big Data Word Cloud; source:<http://olap.com/big-data/>

# Table of Contents

- 1 The Value of Data
- 2 The Big Data Revolution
- 3 An Use Case
- 4 Final Remarks

# Warning



\\* This presentation may contain tech talk, such as: databases, clusters, NoSQL.

If you are susceptible to these concepts, you may want to leave the room **now**.

From this point beyond, you are at your own risk! \\*

# The Value of Data

## Not a New Idea!

- Matthew Fontaine Maury  
(1806-1873).



# The Value of Data

## Not a New Idea!

- Matthew Fontaine Maury (1806-1873).
- Foresaw the hidden value on captain's ships logs, when analysed collectively.
- Used time series data to carry out analysis that would enable him to recommend optimal shipping routes.



# The Value of Data

## Data Mining, Open-source and Crowd-sourcing

Date.	Mile.	Portion	Course steered.	WEATHER.		Temperature.	Depth of the water in fathoms.	Depth of the water in fathoms.	Rate.	Guns.
				Direction.	Fog.		Wind.	Sea.	Wind.	Sea.
<i>A.M.</i>										
1	4	85 <sup>6</sup>	S.E.	2	280.50	90.90	88.80	88.80	Rank 0	3
2	5	-	N.W.	2	-	85.95	85.95	-	0	+
3	5	-	West	1	-	-	-	-	0	+
4	5	-	-	1	-	-	-	-	0	+
5	5	-	S.W.	1	-	88.80	88.80	-	0	+
6	2	-	S.W.	1	280.50	85.90	85.90	-	0	+
7	5	-	S.E.	1	180.00	80.90	80.90	-	0	+
8	5	-	-	1	-	85.95	85.95	-	0	+
9	4	85 <sup>6</sup>	ESE	2	280.50	90.90	88.80	88.80	Rank 0	3
10	5	-	S.E.	1	280.50	90.90	88.80	88.80	0	+
11	5	"East"	Calm	1	280.50	90.90	88.80	88.80	0	+
Nom.	4	85 <sup>6</sup>	ESE	-	-	85.95	85.95	-	0	+

# The Value of Data

## Data Mining, Open-source and Crowd-sourcing

- In 1848, Captain Jackson was the first person to try the route recommended by Maury, and as a result he was able to save 17 days on his outbound trip.

Bore. No.	Polaris.	Course steered.	Wind.	Barometer.			Temperature.			Port of the Cape Horn by squadron.	Port of the Clouds by squadron.	Date of the Clouds by squadron.	Breadth of the Clouds by squadron.	Rate.	Guns.
				Direction.	Ferm.	Leverage.	Height in Tenths of an Inch.	At Sea in the Clouds.	Water at Sea in the Clouds.						
<i>A. M.</i>															
1.	2° 4' E. 85° 2'	S. E.	2.	28.0	51.94	98.85	95.225	Rank 0	3	Steam alone.					
2.	2° 5' S.	" "	2.	"	-	-	95.95	95.95	"	"	0	"			
3.	2° 5' S.	" "	Best.	1.	"	"	"	"	"	"	"	"	"		
4.	2° 5' S.	" "	" "	1.	"	"	"	"	"	CC	"	"	"		
5.	2° 5' S.	" "	S. E.	1.	"	"	98.86	94.91	"	"	0	"			
6.	2° 2' S.	" "	S. E.	1.	28.0	51.95	98.85	"	"	"	0	"			
7.	2° 5' S.	" "	S. E.	1.	28.0	51.95	98.85	"	"	"	0	"			
8.	2° 5' S.	" "	" "	1.	"	"	95.95	94.91	"	"	0	"			
9.	2° 4' E. 85° 2'	" "	" "	2.	28.0	51.94	98.85	95.225	"	"	2250 fath. distance.				
10.	2° 8' S.	" "	S. E.	1.	28.0	51.95	99.95	95.225	CC	"	0	"			
11.	2° 8' S.	" "	Clouds.	1.	28.0	51.95	99.95	95.225	CC	"	0	"			
Nom.	2° 4' E. 85° 2'	" "	" "	6.	"	"	95.95	94.91	"	"	0	"			

# The Value of Data

## Data Mining, Open-source and Crowd-sourcing

- In 1848, Captain Jackson was the first person to try the route recommended by Maury, and as a result he was able to save 17 days on his outbound trip.
- Apart from collecting existing logs, Maury encouraged the collection of more regular and systematic time series, by creating a template.

Date, Month	Latitude	Course steered.	Wind.	VELOCITY.			Rate,	Guns,
				Direction.	Ferm.	Leverage.		
1. 6. 85.2	S.E.	2	28.0. 50. 44. 08. 05. 00.000	W.	0	0	Rank 0	3
2. 7. 85.2	W.	2	- - - 45. 05. 05. 05. - - -	N.	0	0	-	-
3. 7. 85.2	W.	1	- - - - - - - - - - -	N.	0	0	-	-
4. 7. 85.2	W.	1	- - - - - - - - - - -	E.C.	0	0	-	-
5. 7. 85.2	S.E.	1	- - - 45. 05. 05. 05. - - -	N.	0	0	-	-
6. 7. 85.2	S.E.	1	28.0. 50. 44. 08. 05. - - -	N.	0	0	-	-
7. 7. 85.2	S.E.	1	28.0. 50. 44. 08. 05. - - -	N.	0	0	-	-
8. 7. 85.2	W.	1	- - - 45. 05. 05. 05. - - -	N.	0	0	-	-
9. 7. 85.2	E.S.E.	2	28.0. 50. 44. 05. 05. 00.000	N.	0	0	225000000	-
10. 7. 85.2	S.E.	1	28.0. 50. 44. 05. 05. 00.000	N.	0	0	-	-
11. 7. 85.2	Calms	1	28.0. 50. 44. 05. 05. 00.000	N.	0	0	-	-
Nom.	7. 85.2							

# The Value of Data

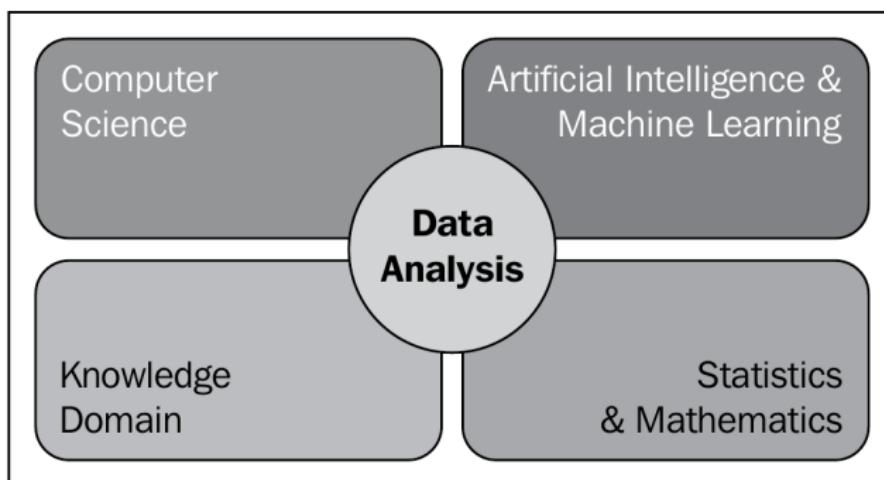
## Data Mining, Open-source and Crowd-sourcing

- In 1848, Captain Jackson was the first person to try the route recommended by Maury, and as a result he was able to save 17 days on his outbound trip.
- Apart from collecting existing logs, Maury encouraged the collection of more regular and systematic time series, by creating a template.
- Collected Data: longitude, latitude, currents, magnetic variation, air and water temperature, general wind direction, etc.

Bore. Month	Latitude	Course steered.	WEATHER.			Rate,	Guns,	
			Direction.	Fog.	Longitude			
			Height in feet above the sea.	At Sea Bath. Waves at sea.	For Cables by telegraph.	Per cent of clouds, by specimen.	Per cent of the time.	Wind at end of route.
A. M.								
1	4° 45' S.	ESE 2	280.50	90.98	85.00	80.000	Rank 0	S. Stem down
2	5° -	NW 2	-	85.95	85.95	-	0	-
3	5° -	West	1	-	-	-	0	-
4	5° -	"	1	-	-	80	0	-
5	5° -	S. W.	1	-	-	80.96.91	-	0
6	2° -	S. W.	1	280.50	90.98.91	-	0	-
7	5° -	S. W.	1	280.50	90.98.91	-	0	-
8	5° -	"	1	-	85.95	90.95	0	225.000.000
9	4° E. 2	ESE 2	280.50	90.98	85.00	80.000	Rank 0	-
10	4° 0' -	S. E.	1	280.50	90.98.91	80	0	-
11	4° "East"	Calm	1	280.50	90.98.91	80.000	0	-
Now.	4° 45' E. 2	"	0	-	85.95	90.95	0	-

# Data Analysis

## A multidisciplinary field



# Data Analysis

## Traditional Stack

- Spreadsheets (e.g.: Excel, OpenOffice)
- RDBMS (e.g.: Oracle, PostgreSQL, MySQL)
- Statistical Packages (e.g.: R, Matlab)
- GIS Packages (e.g.: QGIS)
- Scripting and programming languages (e.g.: Python, Java)
- Libraries

# Big Data

## What changed in recent years?

- Differences in the way global business and transportation are done have exploded the volume of traditional data sources.
- Widespread increase in data from sensors.



# Big Data

## Smart Citizen Project



# Big Data

## Smart Citizen Project

- Platform to generate participatory processes of people in the cities, by connecting data, people and knowledge.



<https://smarcticitizen.me/>

# Big Data

## Smart Citizen Project

- Platform to generate participatory processes of people in the cities, by connecting data, people and knowledge.
- Based on geolocation, Internet and free hardware and software for data collection and sharing.



# Big Data

## Smart Citizen Project

- Platform to generate participatory processes of people in the cities, by connecting data, people and knowledge.
- Based on geolocation, Internet and free hardware and software for data collection and sharing.
- Relies on the production of objects to connect people with their environment and their city.



# Big Data



Based on reports, Fillipo Arrieta published maps describing a multi-stage containment plan designed to limit the plague in Bari (1864).



The iLab used the ushahidi platform to collect and display crowd-sourcing information about Ebola in Liberia (2014).

# Big Data

- A great deal of this data is actually geo-located (e.g.: satellite navigate coordinates, ip addresses).



Based on reports, Fillipo Arrieta published maps describing a multi-stage containment plan designed to limit the plague in Bari (1864).



The iLab used the ushahidi platform to collect and display crowd-sourcing information about Ebola in Liberia (2014).

# Big Data

- A great deal of this data is actually geo-located (e.g.: satellite navigate coordinates, ip addresses).
- Geography has finally the opportunity to switch from being based on guesses and samples, to become a truly data-driven science.



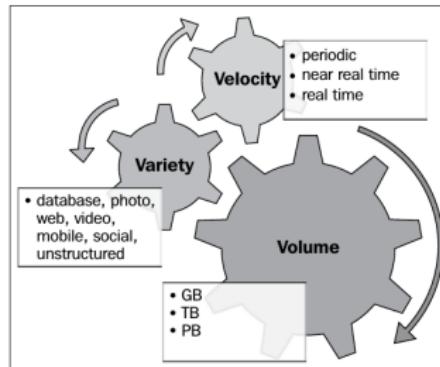
Based on reports, Fillipo Arrieta published maps describing a multi-stage containment plan designed to limit the plague in Bari (1864).



The iLab used the ushahidi platform to collect and display crowd-sourcing information about Ebola in Liberia (2014).

# Big Data

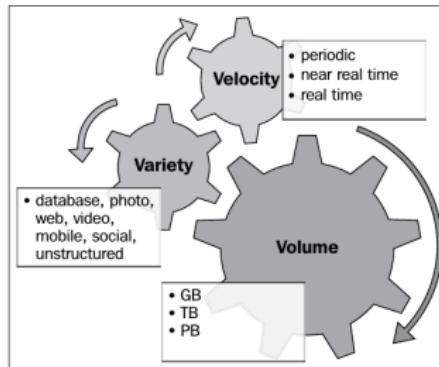
## What characterizes Big Data?



# Big Data

## What characterizes Big Data?

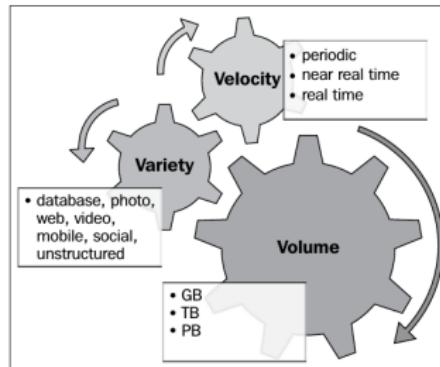
- Volume: Large amounts of data.



# Big Data

## What characterizes Big Data?

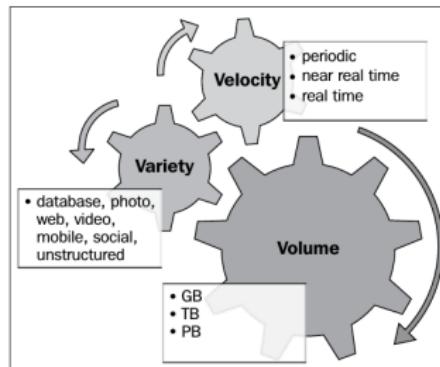
- Volume: Large amounts of data.
- Variety: Different types of structured, unstructured, and multi-structured data.



# Big Data

## What characterizes Big Data?

- Volume: Large amounts of data.
- Variety: Different types of structured, unstructured, and multi-structured data.
- Velocity: Needs to be analyzed quickly.



# Technological Challenges

These characteristics map into challenges:

- Scalability
- Heterogeneity
- Low latency

When the traditional stack is no longer enough, a paradigm shift is required.



# RDBMS vs NoSQL



# RDBMS vs NoSQL



- NoSQL databases trade away some capabilities of relational databases (SQL), in order to improve scalability.

# RDBMS vs NoSQL



- NoSQL databases trade away some capabilities of relational databases (SQL), in order to improve scalability.
- Advantages: NoSQL databases are simpler, can handle semi-structured and denormalized data and have an higher scalability.

# RDBMS vs NoSQL



- NoSQL databases trade away some capabilities of relational databases (SQL), in order to improve scalability.
- Advantages: NoSQL databases are simpler, can handle semi-structured and denormalized data and have an higher scalability.
- Disadvantages: loss of abstraction provided by the query optimizer, that increases the complexity of the applications.

# RDBMS vs NoSQL



- NoSQL databases trade away some capabilities of relational databases (SQL), in order to improve scalability.
- Advantages: NoSQL databases are simpler, can handle semi-structured and denormalized data and have an higher scalability.
- Disadvantages: loss of abstraction provided by the query optimizer, that increases the complexity of the applications.
- Recently, tools were developed that bring back the full power of SQL language to the NoSQL ecosystem (e.g.: Apache Drill, Hive).

# RDBMS vs NoSQL

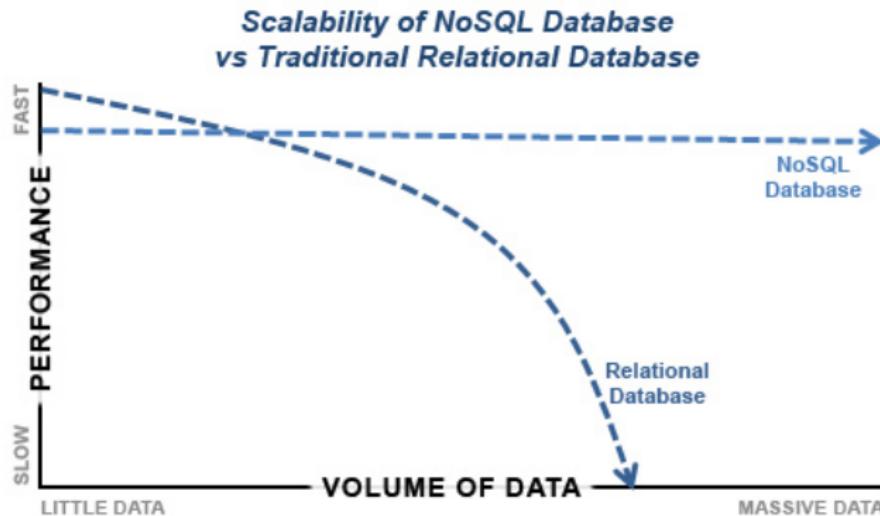


Image Credit: DataJobs.com

# Examples



cassandra

*hadoop*

RethinkDB.

HYPERTABLE<sup>TM</sup>

# MapReduce

Programming model for processing and generating large data sets with a parallel, distributed algorithm on a cluster.

- Map(): performs filtering and sorting.
- Reduce(): performs a summary operation.

# MapReduce

Programming model for processing and generating large data sets with a parallel, distributed algorithm on a cluster.

- Map(): performs filtering and sorting.
- Reduce(): performs a summary operation.
- The framework coordinates the processing, by marshalling the distributed servers, running the tasks and parallel and managing all communications and data between the various parts of the system.

# MapReduce

Programming model for processing and generating large data sets with a parallel, distributed algorithm on a cluster.

- Map(): performs filtering and sorting.
- Reduce(): performs a summary operation.
- The framework coordinates the processing, by marshalling the distributed servers, running the tasks and parallel and managing all communications and data between the various parts of the system.
- There are many libraries that implement MapReduce.

# A Big Data Approach

**Does this mean we have to throw away our traditional tools and methods?**

# A Big Data Approach

**Does this mean we have to throw away our traditional tools and methods?**

- First ensure that you **really** have, or will have Big Data at some point in the future.

# A Big Data Approach

**Does this mean we have to throw away our traditional tools and methods?**

- First ensure that you **really** have, or will have Big Data at some point in the future.
- Then identify the stages of the workflow that are bottlenecks, in terms of the current technologies.

# A Big Data Approach

**Does this mean we have to throw away our traditional tools and methods?**

- First ensure that you **really** have, or will have Big Data at some point in the future.
- Then identify the stages of the workflow that are bottlenecks, in terms of the current technologies.
- It is possible to mix & and match.

## Analysing geo-located Tweets



## Analysing geo-located Tweets

- The purpose of this use case was to analyse the stream of geo-located Tweets, as a sensor for citizen presence.



## Analysing geo-located Tweets

- The purpose of this use case was to analyse the stream of geo-located Tweets, as a sensor for citizen presence.
- Number of tweets in Catalunya in around 3 months: +- 6 million.



## Analysing geo-located Tweets

- The purpose of this use case was to analyse the stream of geo-located Tweets, as a sensor for citizen presence.
- Number of tweets in Catalunya in around 3 months: +- 6 million.
- Continuous stream of data.



## Analysing geo-located Tweets

- The purpose of this use case was to analyse the stream of geo-located Tweets, as a sensor for citizen presence.
- Number of tweets in Catalunya in around 3 months: +- 6 million.
- Continuous stream of data.
- This amount of data is not easily assimilated by the "human-eye", so we decided to create clusters of Tweets.



# An Use Case

## Clustering

# An Use Case

## Clustering

- It is a descriptive data mining technique, often used for dimensionality reduction.

# An Use Case

## Clustering

- It is a descriptive data mining technique, often used for dimensionality reduction.
- It groups a set of objects in such a way that objects in the same group are more similar to each other, than to object in other groups.

# An Use Case

## Clustering

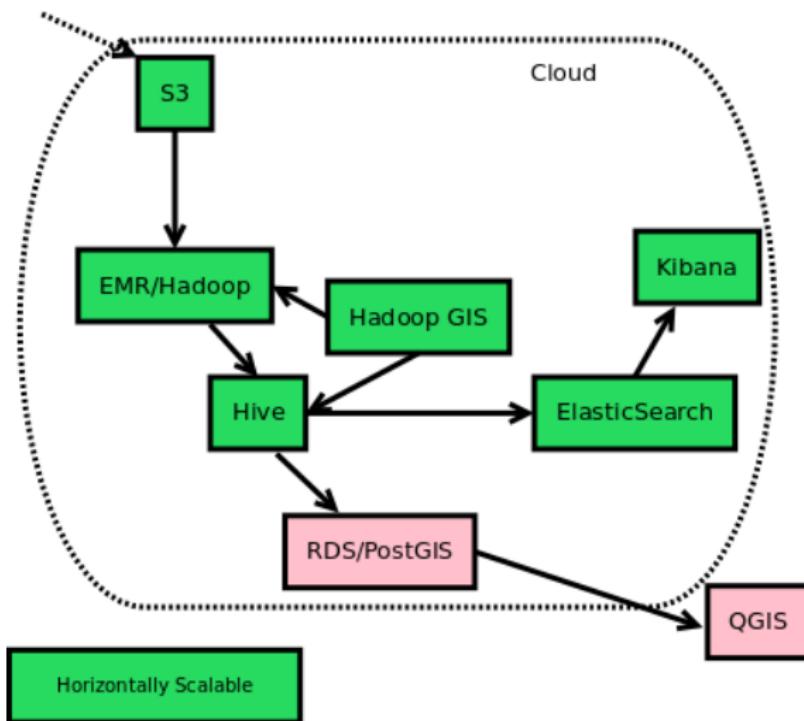
- It is a descriptive data mining technique, often used for dimensionality reduction.
- It groups a set of objects in such a way that objects in the same group are more similar to each other, than to object in other groups.
- Strictly, it corresponds to a family of unsupervised machine learning algorithms.

# An Use Case

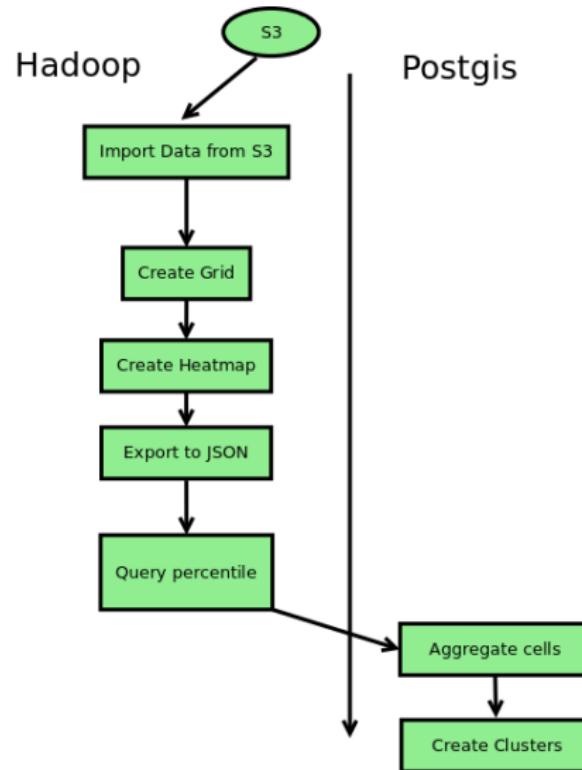
## Clustering

- It is a descriptive data mining technique, often used for dimensionality reduction.
- It groups a set of objects in such a way that objects in the same group are more similar to each other, than to object in other groups.
- Strictly, it corresponds to a family of unsupervised machine learning algorithms.
- We wanted to implement this concept using only Hadoop, and apply it to spatial attributes.

# Technological Stack



# Workflow



# Results

Using this workflow we were able to turn the original raw tweets, first into a density grid, and then into clusters.



As we identified and solved the bottlenecks with the relevant tools, in theory this algorithm is scalable to any Petabytes of data.

# Final Remarks

## Final Remarks

- It is ok to use non scalable tools, at certain points of the workflow.

## Final Remarks

- It is ok to use non scalable tools, at certain points of the workflow.
- We are working on the edge of existing technologies; some functions were not implemented yet, and bugs are common.

## Final Remarks

- It is ok to use non scalable tools, at certain points of the workflow.
- We are working on the edge of existing technologies; some functions were not implemented yet, and bugs are common.
- Since it is a niche, this is even more true for spatial technologies.

## Final Remarks

- It is ok to use non scalable tools, at certain points of the workflow.
- We are working on the edge of existing technologies; some functions were not implemented yet, and bugs are common.
- Since it is a niche, this is even more true for spatial technologies.
- There are not ready-made solutions: the particular stack and workflow should be compiled for the specific case study.

## Final Remarks

- It is ok to use non scalable tools, at certain points of the workflow.
- We are working on the edge of existing technologies; some functions were not implemented yet, and bugs are common.
- Since it is a niche, this is even more true for spatial technologies.
- There are not ready-made solutions: the particular stack and workflow should be compiled for the specific case study.
- this is one stack to solve this problem; it is not the only one, and it may not even be the "best" one;

## Acknowledgements

I would like to thank Ellen Friedman (MapR, Apache Mahout, Apache Drill), for her inspiring work and for taking the time for kindly reviewing this presentation.



## References

- Cuesta, H. "Practical Data Analysis". Packt Publishing (2013)
- Dunning, T. and Friedman, E. "Time Series Databases: New Ways to Store and Access Data". O'Reilly Media; 1 edition (October, 2014)
- Myatt, G and Johnson, W. "Making Sense Of Data I: A Practical Guide to Exploratory Data Analysis and Data Mining". O'Reilly : 2 edition (2014).