

Exploratory Analysis in Big Data based on PCA and PLS

José Camacho

Departamento de Teoría de la Señal, Telemática y Comunicaciones



Network Engineering & Security Group
<http://nesg.ugr.es>



UGR | Universidad
de Granada

Velocity



Veracity

Big Data

Sensors

Text

Records

Images

Audio

Video

Biological

Variety

VOLUME

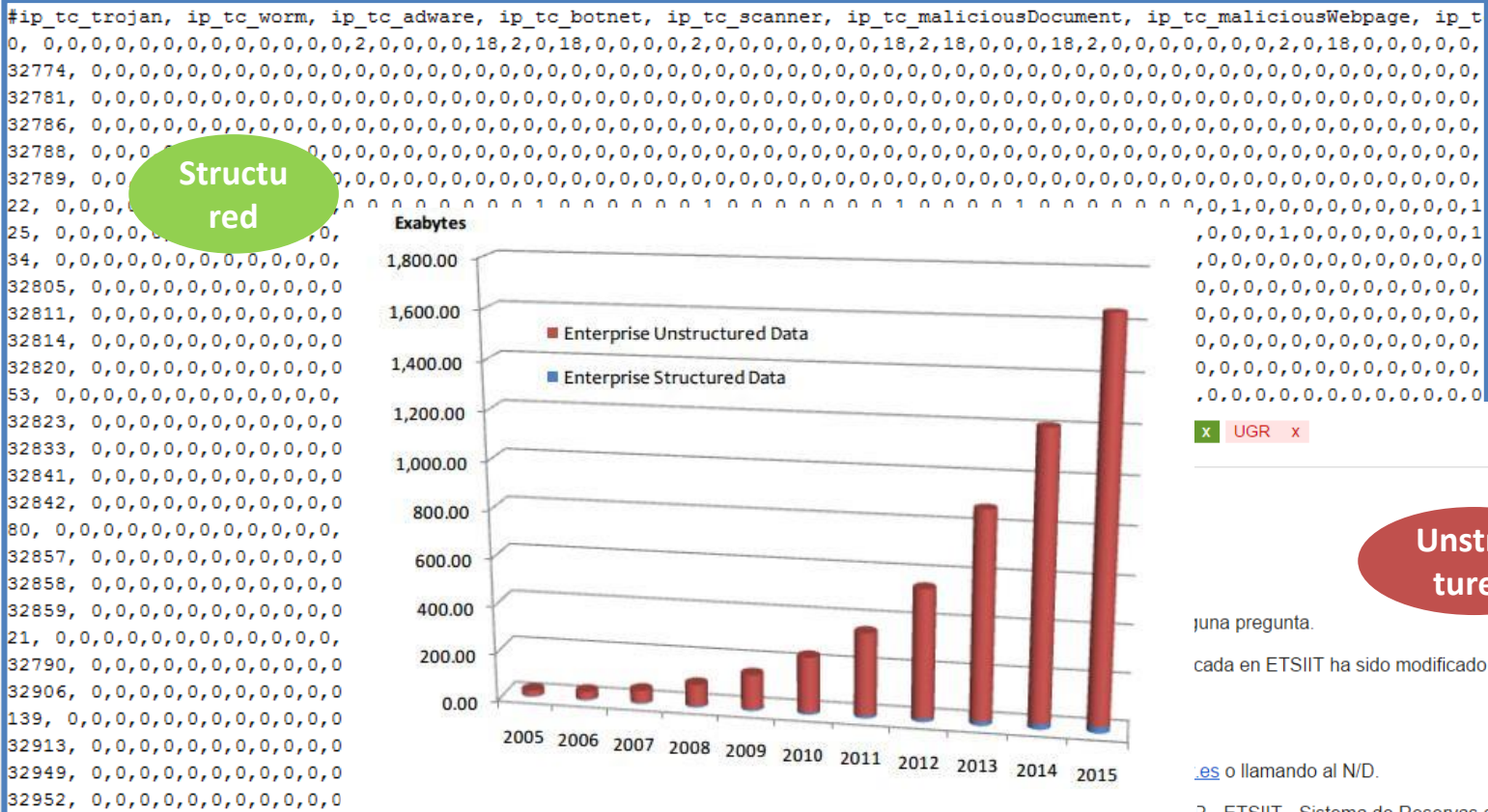
→ Numbers

Unit	Size	Size
Gigabyte	10^9	1.000.000.000 bytes
Terabyte	10^{12}	1.000.000.000.000 bytes
Petabyte	10^{15}	1.000.000.000.000.000 bytes
Exabyte	10^{18}	1.000.000.000.000.000.000 bytes
Zettabyte	10^{21}	1.000.000.000.000.000.000.000 bytes

- ✓ *There was 5 exabytes of information created between the dawn of civilization through 2003, but that much information is now created every 2 days, and the pace is increasing*



RJMetrics co-founder and CEO Robert J. Moore



Unstruc tured

una pregunta. cada en ETSIIT ha sido modificado.

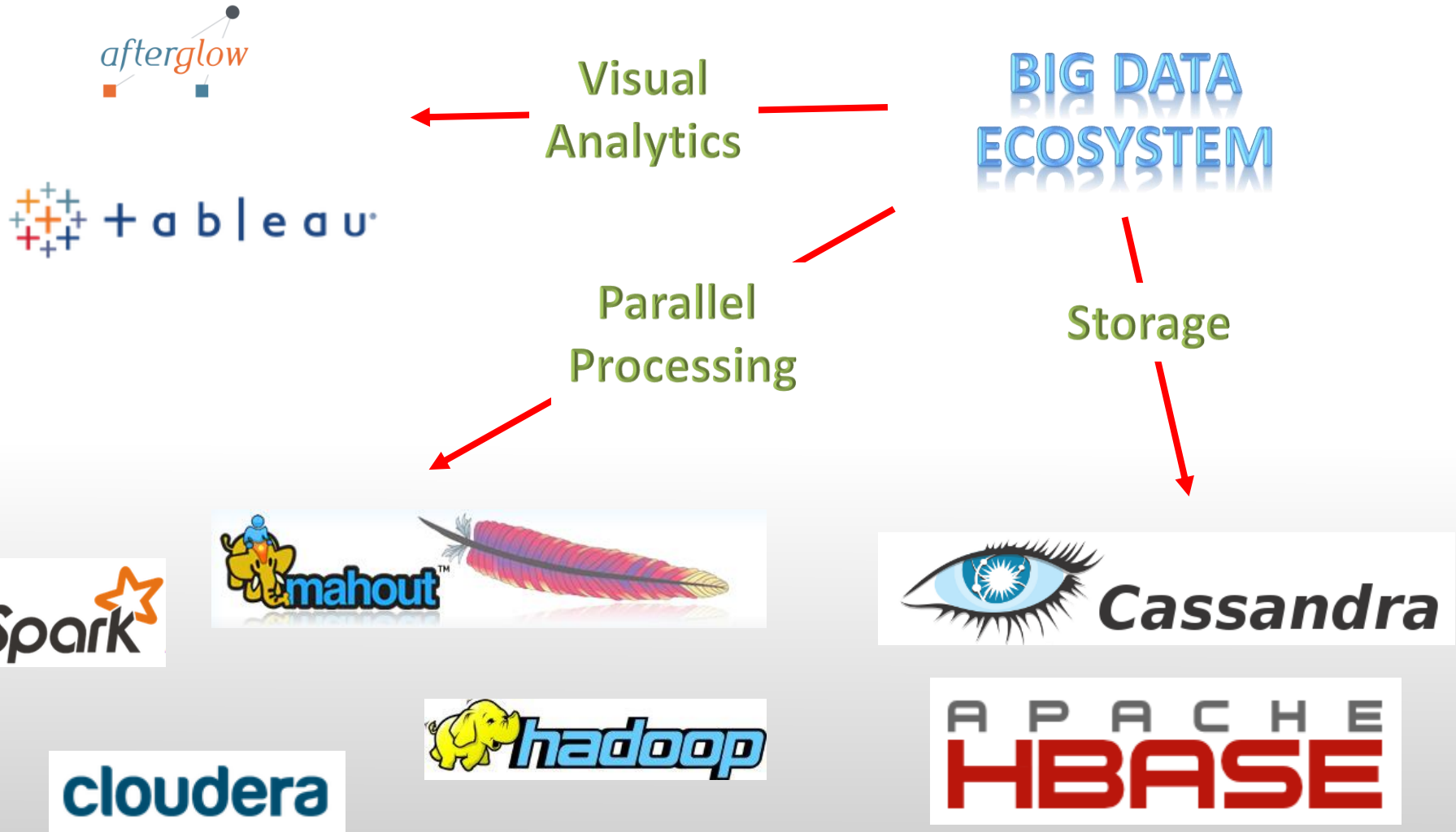
es o llamando al N/D. R - ETSIIT - Sistema de Reservas de Aulas en:

Variety

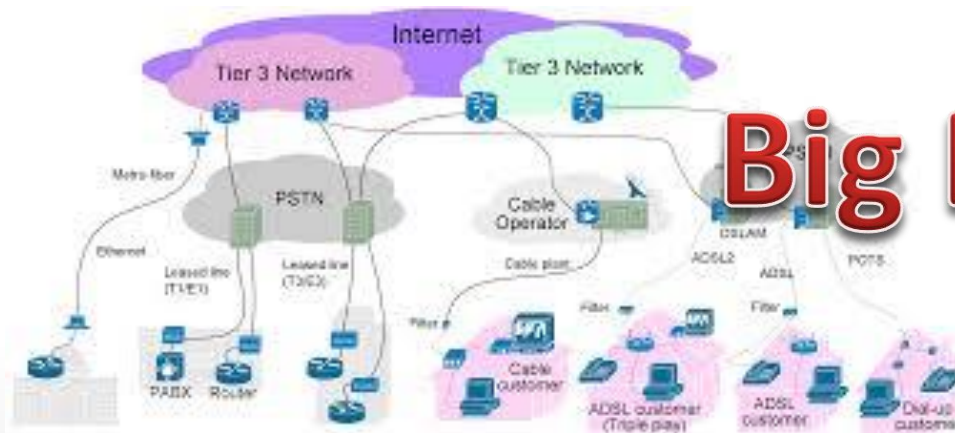


Michael Walker, Data Science Central

Reserva #	Fecha Inicial	Fecha Final	Recurso	Hora de Inicio	Hora de Finalización	Ubicación	Contacto
sc1557544861fc22	10/06/2015	10/06/2015	Sala de juntas	13:00	14:00	ETSIIT	N/D



➔ Problem: Handling communication networks



Big Data



➔ Solution: Multivariate Analysis??

✓ Networkmetrics

→ Networkmetrics

- ✓ Applications for Exploratory Analysis, Optimization, Classification, Anomaly Detection (\cong **Chemometrics**)
- ✓ Big Data by Definition (\neq **Chemometrics**)
 - 4 V's: Tons of data, high speed, from lots of sources, many false alarms....
 - Mostly unstructured → Feature Engineering
- ✓ Complex Data (\cong **Chemometrics**):
 - Fusion
 - High dimensional
 - N-way

Security Monitor

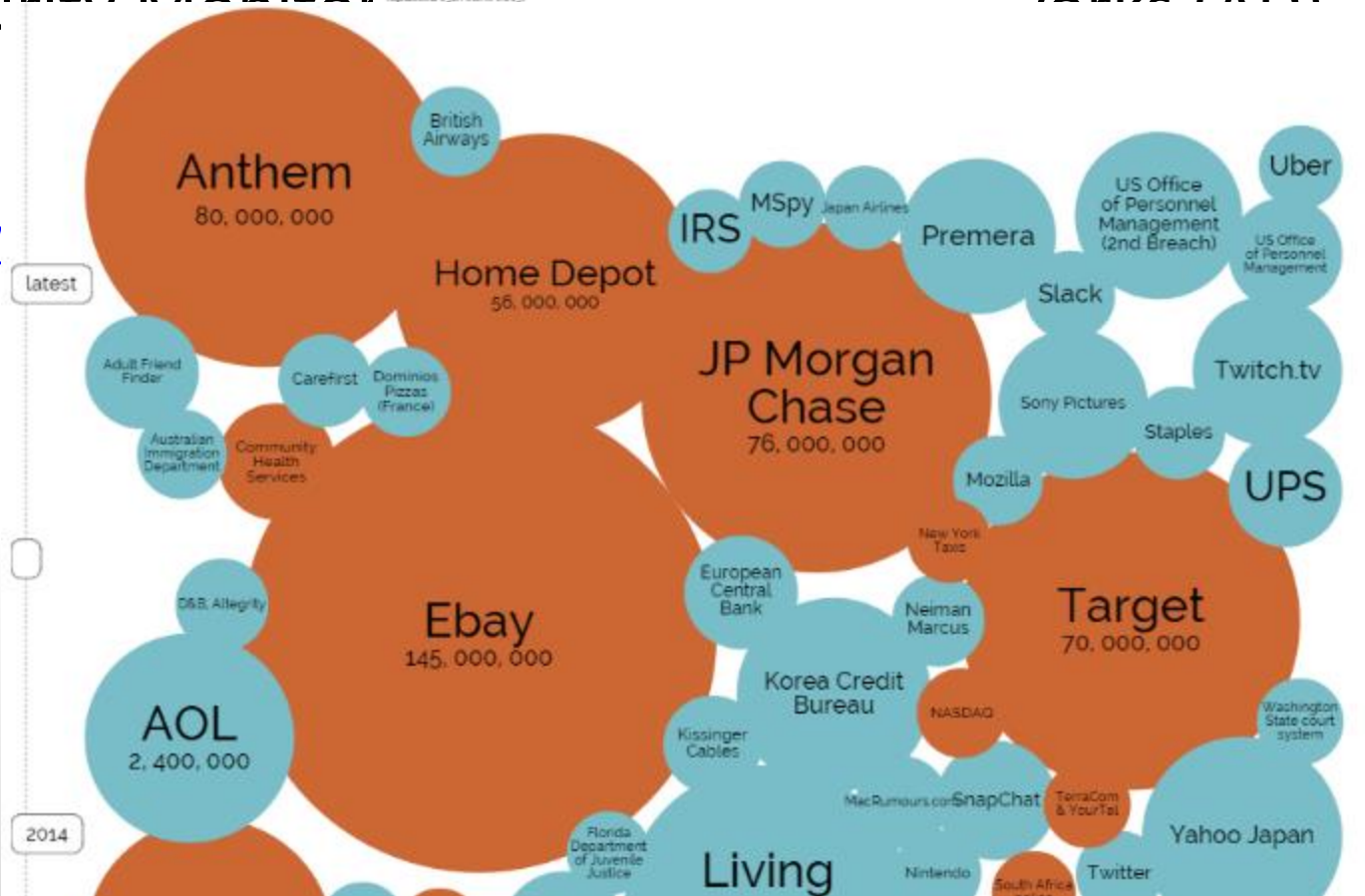
World's Biggest Data Breaches
Selected losses greater than 30,000 records
(updated 13th June 2015)

interesting story

links (AD)

G

<http://>



)

<https://>



→ Secu

jcamacho

pwned?

→ Gr

josecamacho@ugr.es

pwned?

<http://w>

→ Is

<https://>

Oh no — pwned!

Pwned on **1** breached site and found **no** pastes

[Notify me when I get pwned](#) [Donate](#)



Breaches

A "breach" is an incident where a site's data has been illegally accessed by hackers and then released publicly. Review the types of data that were compromised (email addresses, passwords, credit cards etc.) and take appropriate action, such as changing passwords.

Forb

Breaches you were pwned in

A "breach" is an incident where a site's data has been illegally accessed by hackers and then released publicly. Review the types of data that were compromised (email addresses, passwords, credit cards etc.) and take appropriate action, such as changing passwords.



Adobe: The big one. In October 2013, 153 million Adobe accounts were breached with each containing an internal ID, username, email, *encrypted* password and a password hint in plain text. The password cryptography was poorly done and many were quickly resolved back to plain text. The unencrypted hints also disclosed much about the passwords adding further to the risk that hundreds of millions of Adobe customers already faced.

Compromised data: Email addresses, Password hints, Passwords, Usernames

email address or username

pwned?

➔ Security Monitoring in Computer Networks (AD)

➔ Growing tendency of Security Breaches (Leakage of Info)

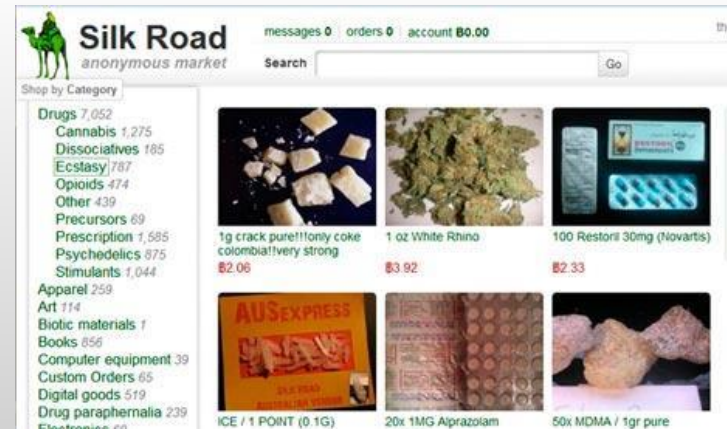
<http://www.informationisbeautiful.net/visualizations/worlds-biggest-data-breaches-hacks/>

➔ Is my info at risk?

<https://haveibeenpwned.com/>



DEEP WEB
(96 %)



➔ Security Monitoring in Computer Networks (AD)

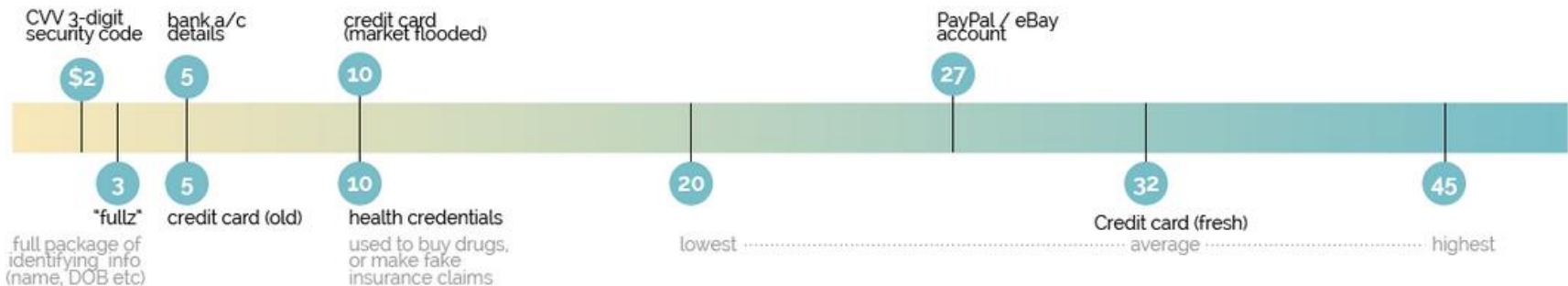
➔ Growing tendency of Security Breaches (Leakage of Info)

<http://www.informationisbeautiful.net/visualizations/worlds-biggest-data-breaches-hacks/>

➔ Is my info at risk?

<https://haveibeenpwned.com/>

How Much is Your Hacked Data Worth? Black market \$ prices

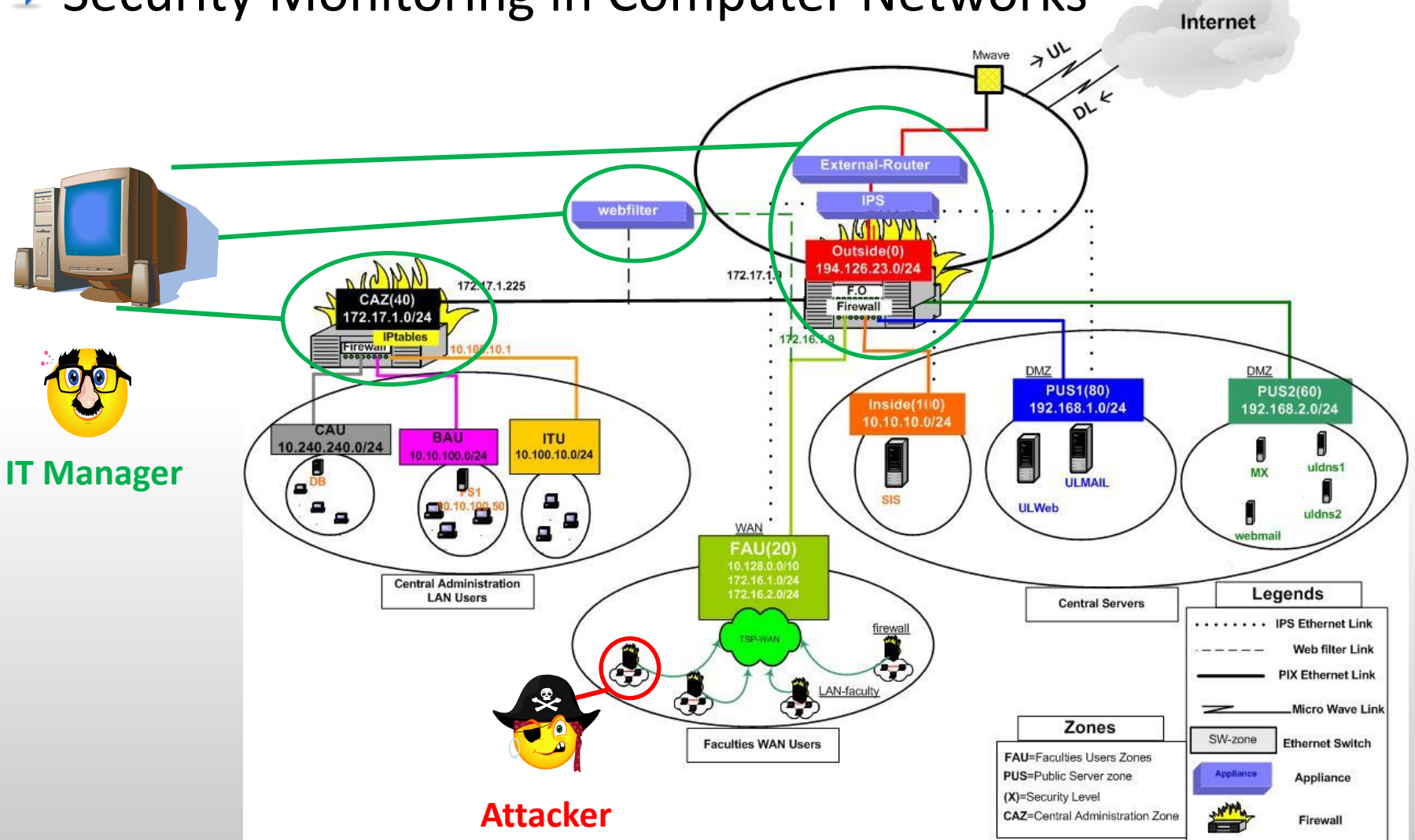


informationisbeautiful.net

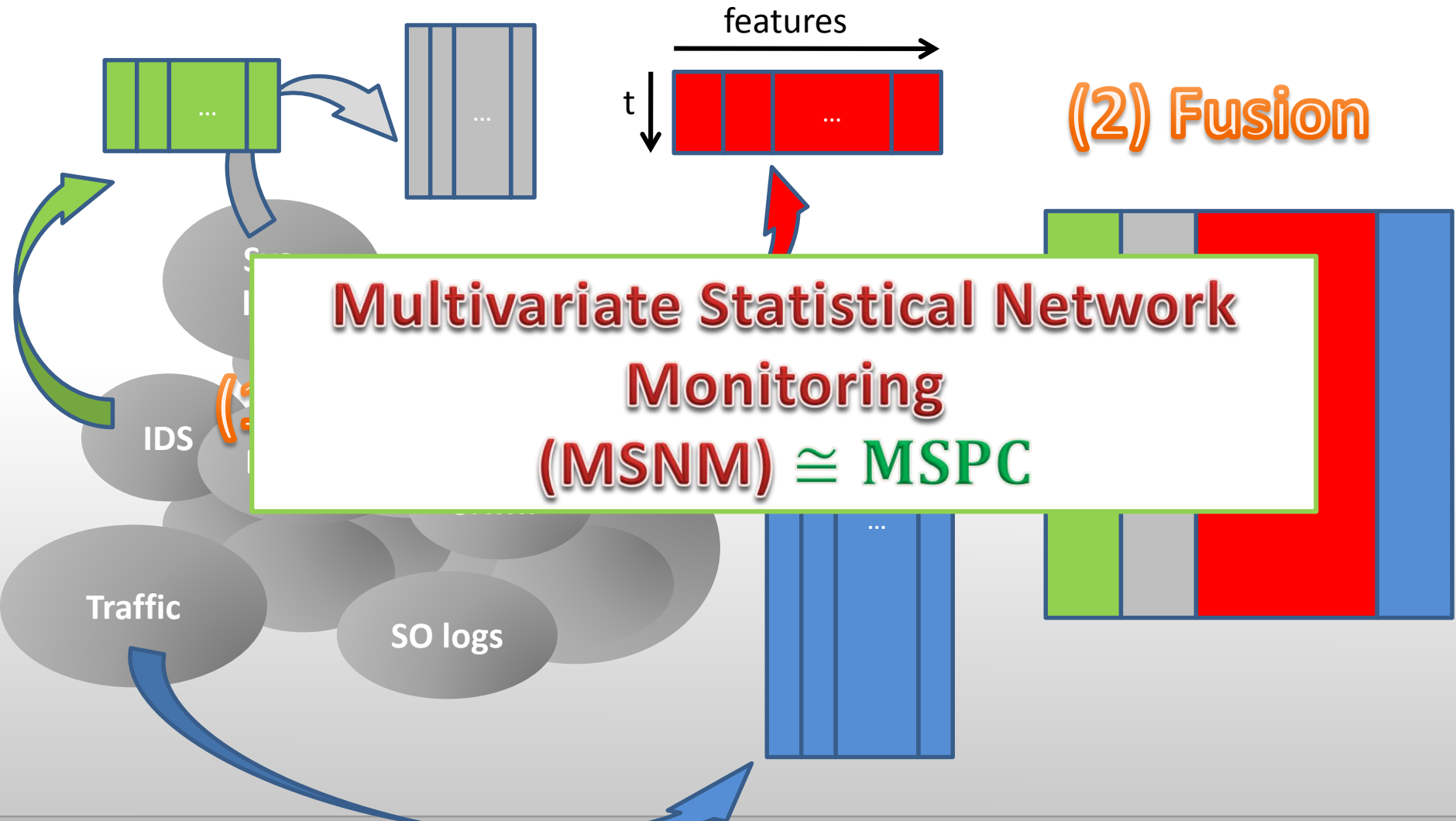
data: bit.ly/bigdatabreaches

sources: Holt & Smirnova (2014), Reuters, Globe & Mail, Rand

Security Monitoring in Computer Networks



How may Multivariate Analysis help???



→ MEDA Toolbox: <https://github.com/josecamachop/MEDA-Toolbox>

✓ Dimensionality:

- Scree plots
- CV (ekf, cekf, ckf)
- SVI Plots

✓ Structure in Variables:

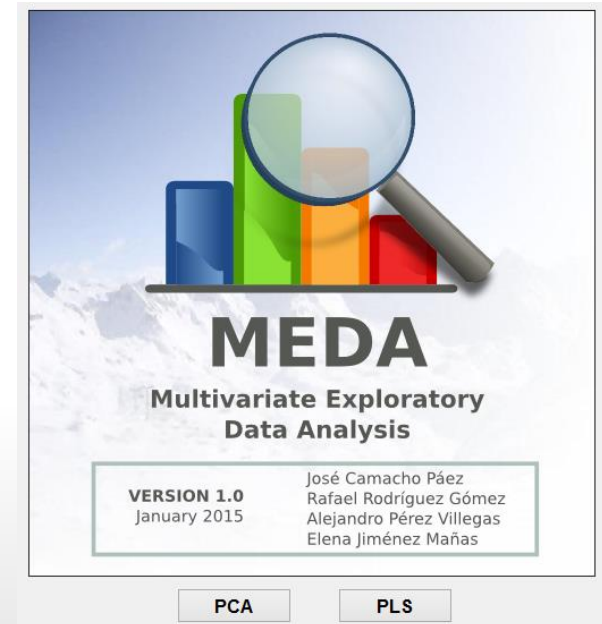
- Loading plots
- MEDA plots

✓ Distribution of Observations

- Score plots
- MSPC: D-st, Q-st
- Covariance MSPC: ADICOV

✓ Observations vs Variables

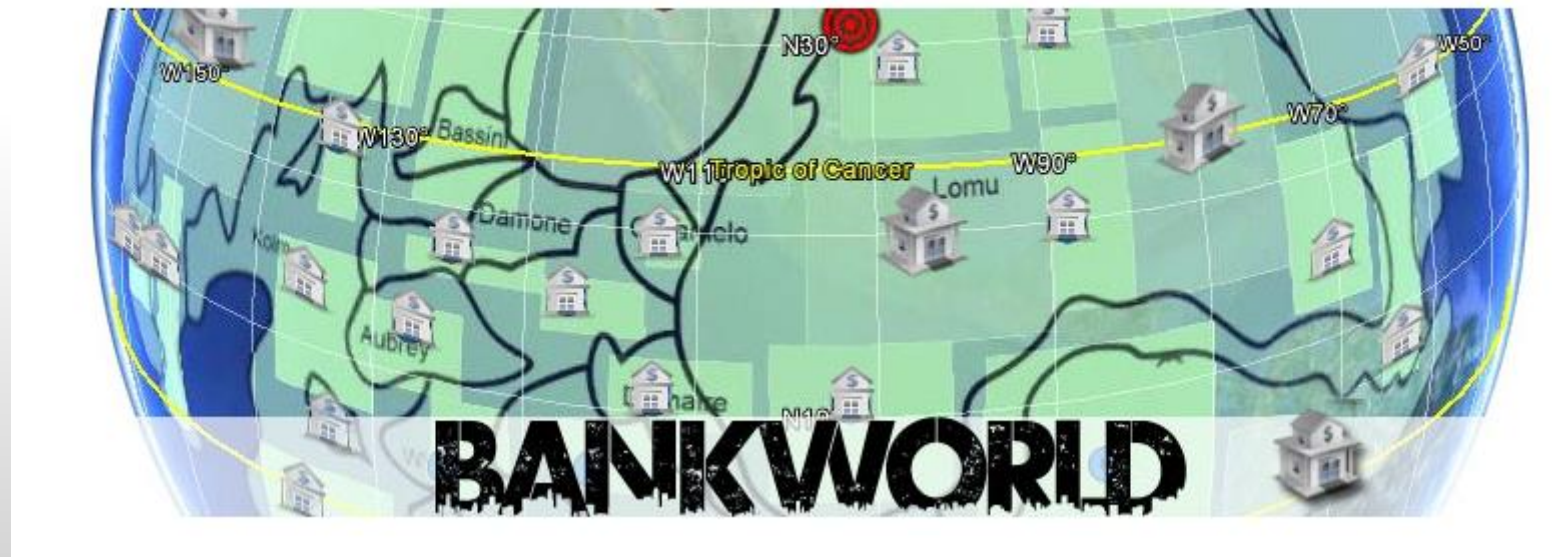
- oMEDA plots



ChemoLab, (2015) 143: 49

→ Networkmetrics: MSNM

VAST Challenge 2012



198.41.0.4
128.9.0.107
192.33.4.12
128.8.10.90
192.203.230.10
192.5.5.241
192.142.22.4

VAST 2012
Challenge 2
Bank of Money
Work Architecture

Structured:
23,711,341
records

Unstructured:
35,984 records

The screenshot displays a network analysis interface. At the top, a list of IP addresses is shown. Below it, a window titled 'syslog-03292012-1hr-parsed' shows a table of log entries with columns for Date/time, Log, Destination, and Action. A green oval highlights a specific record in this table, indicating it is part of the structured data. Below the log window, a network diagram shows various components: 'Websites' (represented by a server icon), 'Workstations' (represented by a computer icon), 'Financial Servers' (represented by a server icon), 'IDS' (represented by a server icon), 'Log Server' (represented by a server icon), and 'Domain Controller / DNS' (represented by a server icon). A red oval highlights a specific record in the log window, indicating it is part of the unstructured data. The log window shows detailed information about a NetBIOS session setup attempt, including classification, priority, and various fields like TTL, TOS, ID, IpLen, DgmLen, Seq, Ack, Win, and TcpLen. It also includes references to security advisories and CVEs.

➔ Parsing & Data Fusion

✓ X = 2,350 obs x 90 features

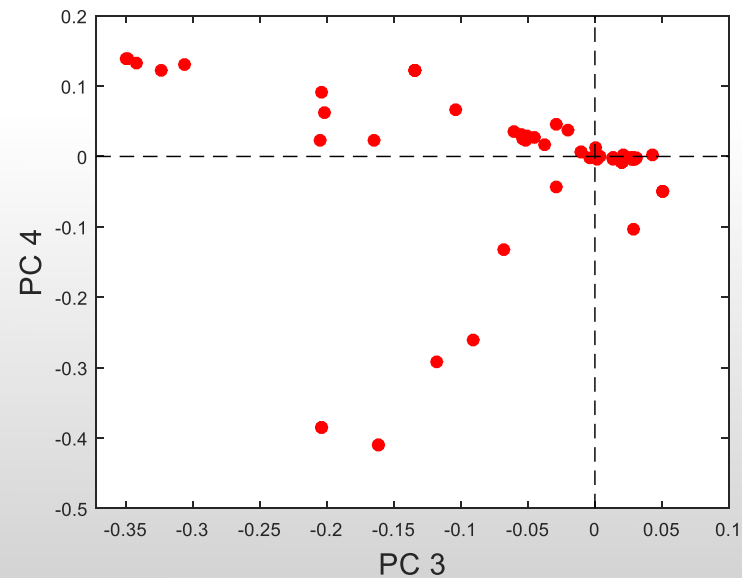
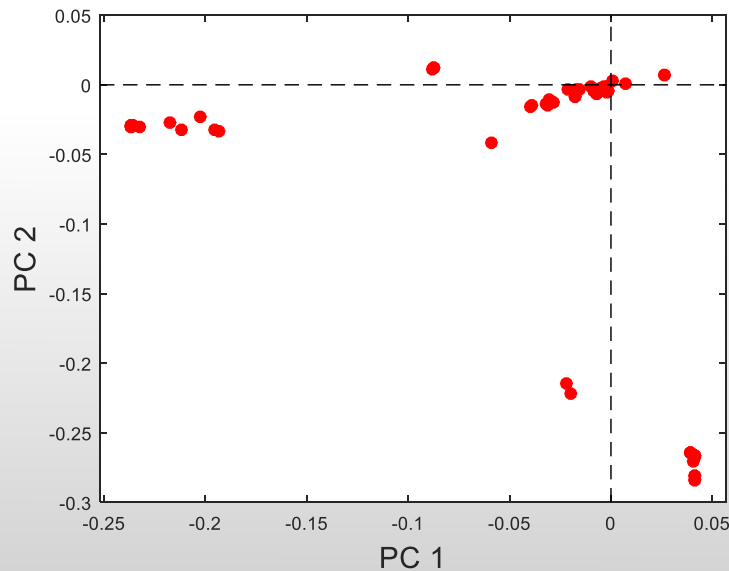
	Feature class	#features
Firewall log	Syslog priority	5
	Operation	6
	Message code	25
	Protocol	3
	IP address	9
	Port number	17
	Direction	2
	Conn. built/teardown	2
	Subtotal	69
IDS log	IP address	9
	Port number	17
	Alert class	5
	Priority	3
	Text label	9
	Subtotal	43

PCA Model

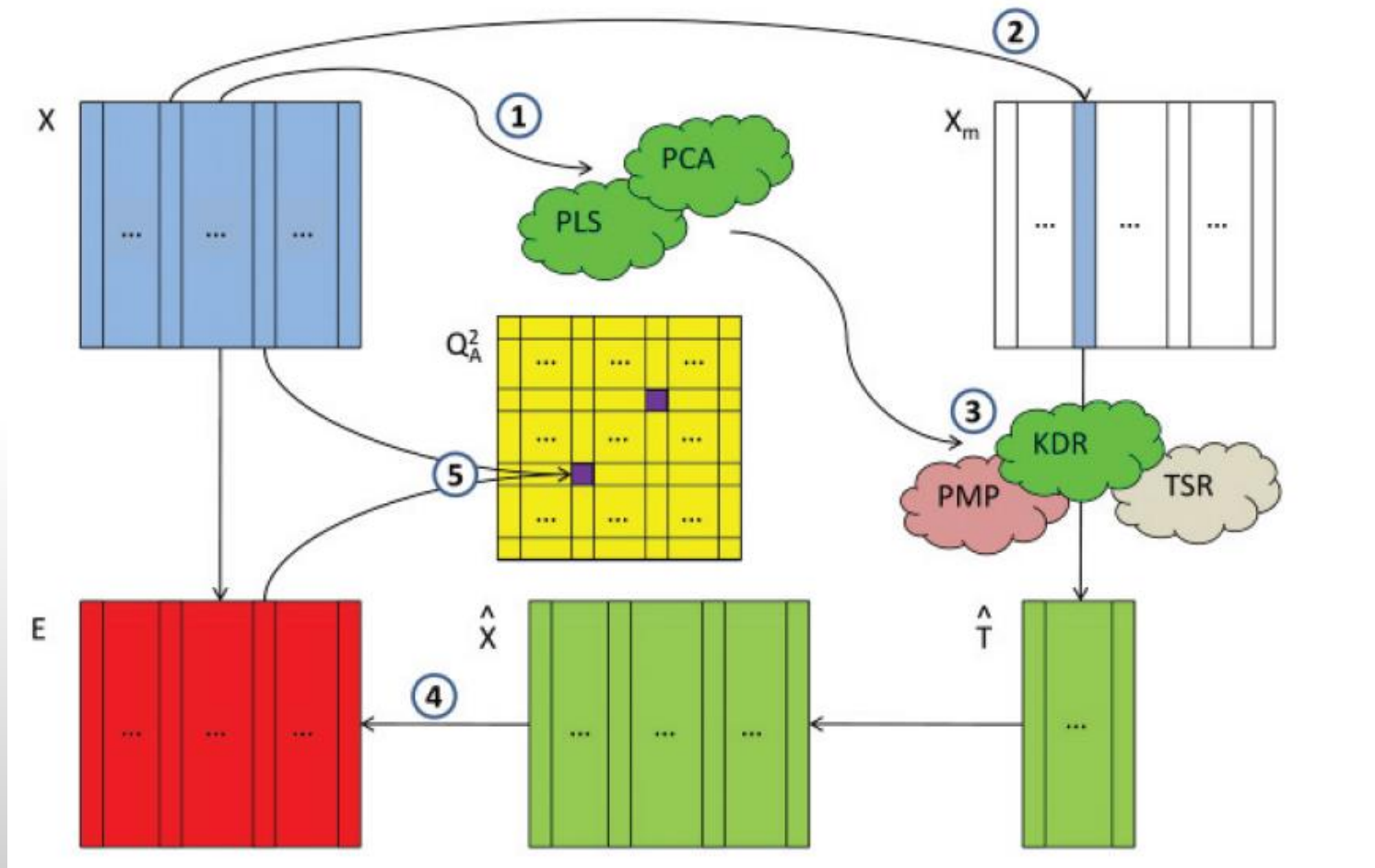
6 PCs and 60% variance app.

➔ Loading plots

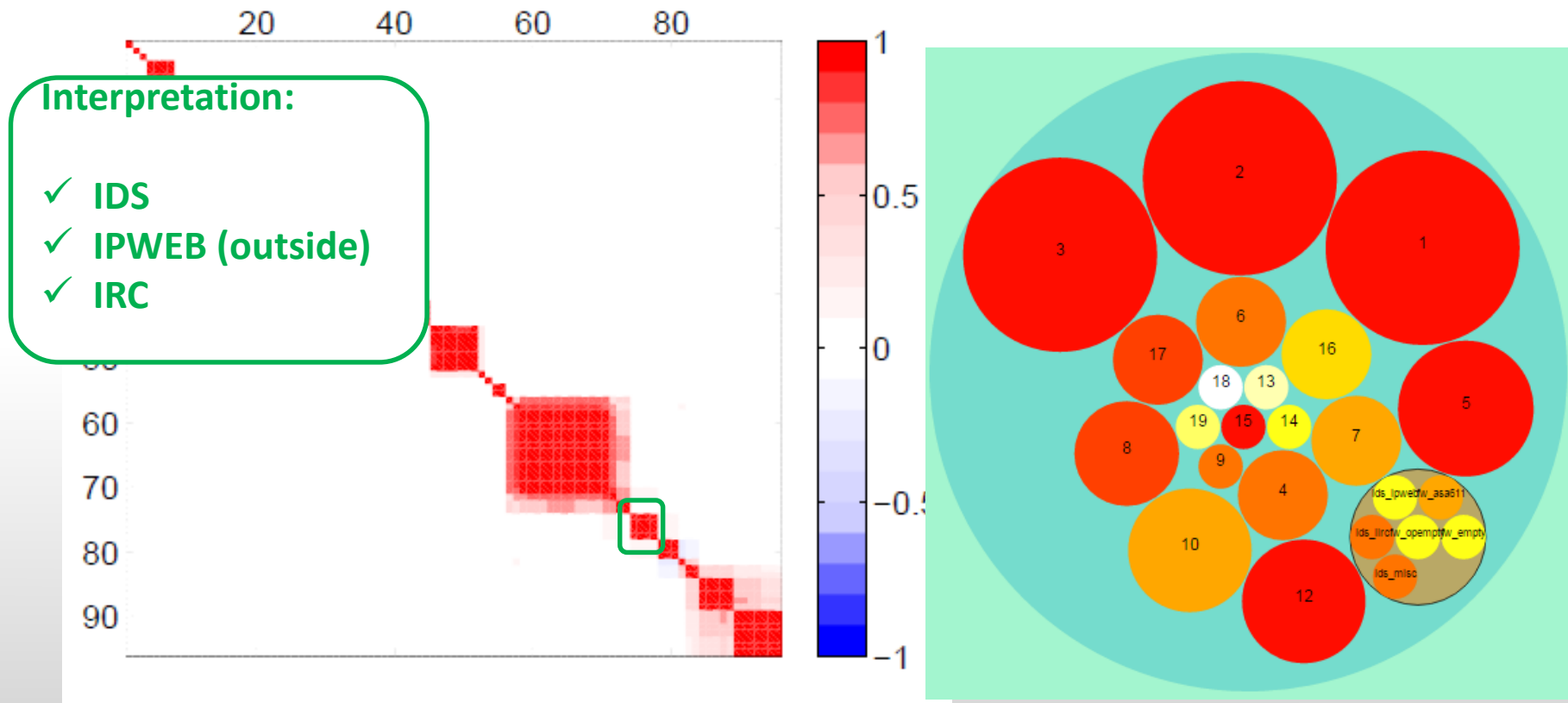
- LP may lead to non-existent correlations, MEDA not
- LP limited to 2/3 PCs, MEDA not
- LP difficult to interpret for many variables



→ MEDA:

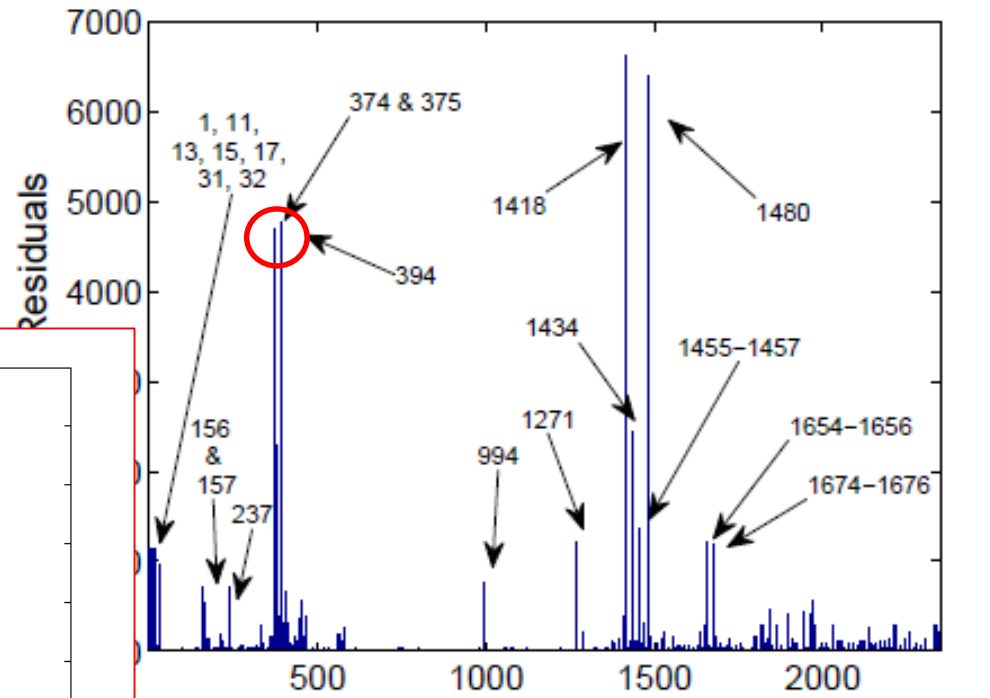
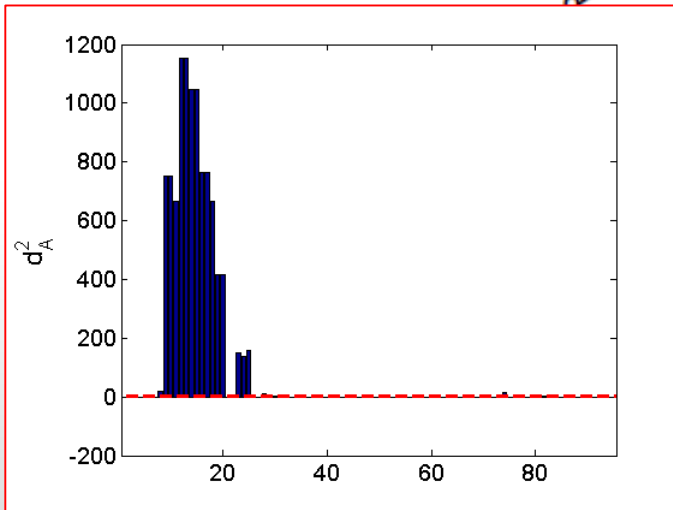


→ MEDA: 6 PCs



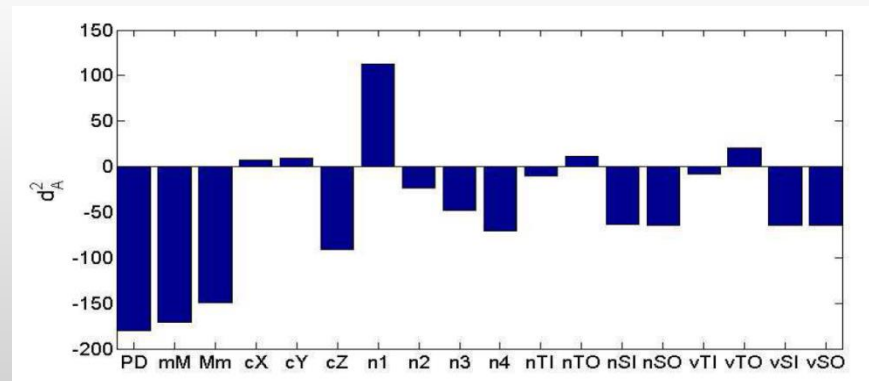
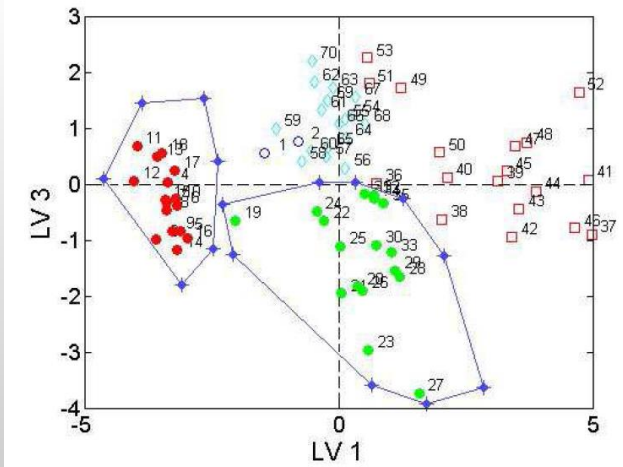
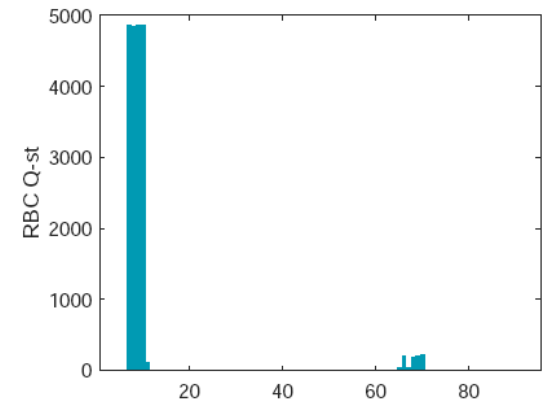
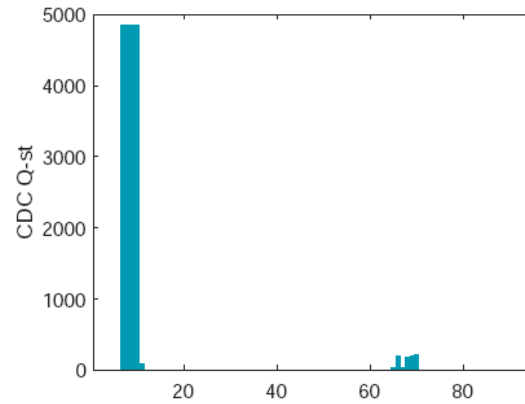
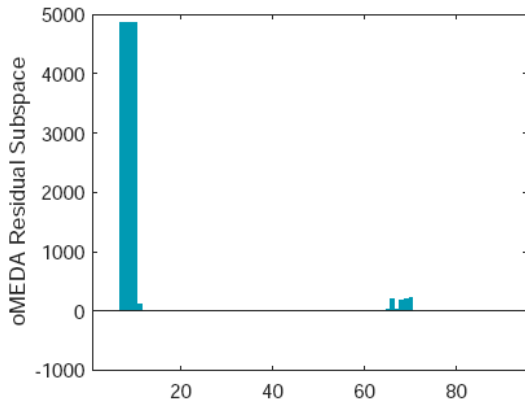
➔ MSPC Q-st

➔ oMEDA

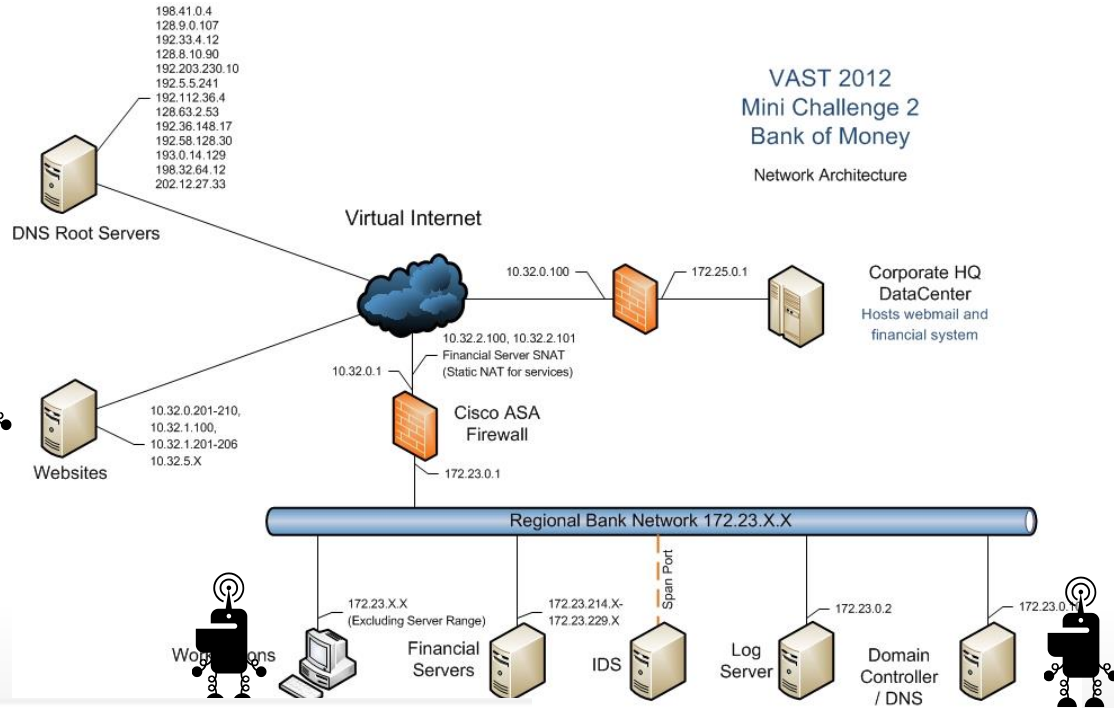
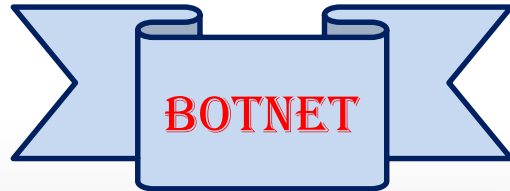


✓ Interpretation: Non-legitime requests of services in access firewall.

oMEDA: Generalization of Contribution plots



➔ Conclusion

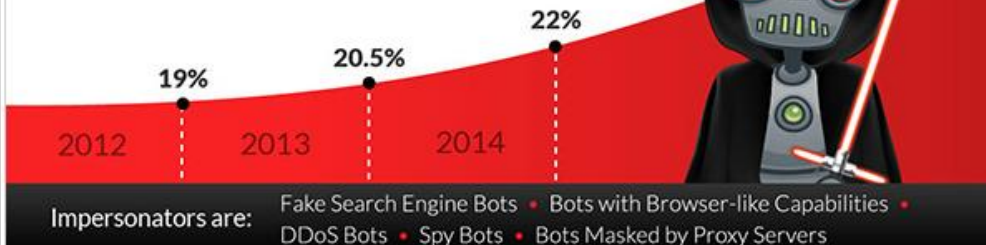


VAST 2012
Mini Challenge 2
Bank of Money
Network Architecture

Coming in 2015

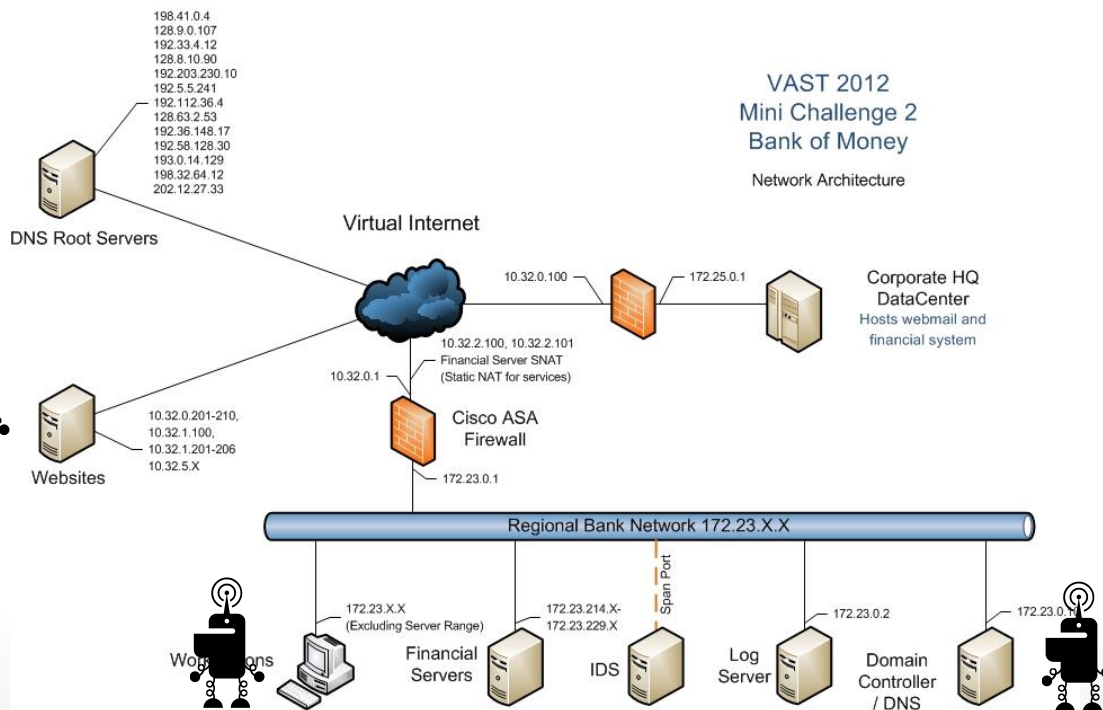
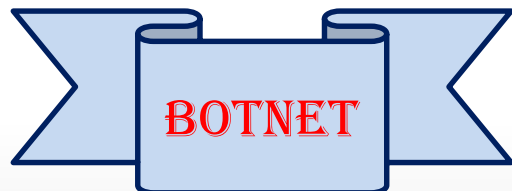
The rise of the impersonator bots

Impersonator bots are the most advanced, malicious non-humans. They are the only group of bots to consistently display growth in the last 3 years.



Incapsula Botnet Report.
2014

➔ Conclusion



Anomaly #	Proposed	Winner [30]	[31]	[32]
Attacks to the DNS / Domain Controller	X			
Access attempts to the firewall	X	P	X	
FTP attempts from inner nodes to outer node	X			X
Background IRC activity	X	X	X	
Parsing errors	X			



"Tackling the Big Data 4 Vs for Anomaly Detection". *INFOCOM'2014 Workshop on Security and Privacy in Big Data*. 2014.

➔ Extensions for Big Data

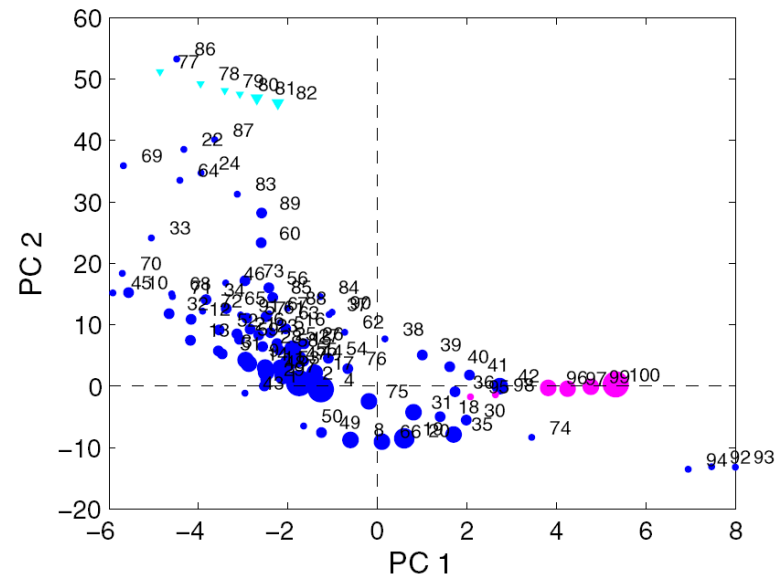
- ✓ For variables ➔ Kernel matrix with EWMA update

$$(X'X)_t = \lambda \cdot (X'X)_{t-1} + \tilde{X}_t' \cdot \tilde{X}_t$$

- Scalable to any size
- PCA/PLS Kernel algorithms
- MEDA from covariance
- oMEDA from accumulates
- ADICOV MSPC

- ✓ For observations ➔ Clustering

- A link with visualization
- Compressed Score Plots
- Compressed MSPC



ChemoLab, (2014) 135: 110

```
clear
load kdd

Lmodel = Lmodel_ini; % Initialization
Lmodel.update = 2; % Change this to 1 for EWMA and 2 for Iterative
Lmodel.type = 2; % Change this to 1 for PCA and 2 for PLS
Lmodel.lv = 3; % Initial number of LVs
Lmodel.prep = 2; % X-block prepr. 0: None, 1: Mean-center, 2: Auto-scaling
Lmodel.prepy = 2; % Y-block prepr. 0: None, 1: Mean-center, 2: Auto-scaling
Lmodel.nc = 100; % Number of clusters

lambda = 1-1e-4; % Forgetting factor in EWMA
step = 0.01;

%% Model building (EWMA or Iterative)

if Lmodel.update == 1
    Lmodel = update_ewma(short_list, '', Lmodel, lambda, step, 1); % EWMA
else
    Lmodel = update_iterative(short_list, '', Lmodel, 20, step, 0, '', 1); % Iterative
end

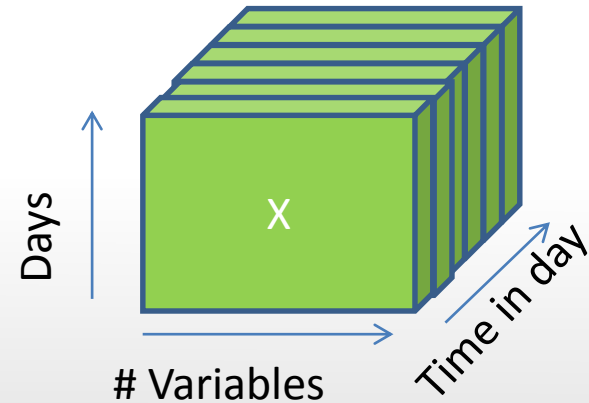
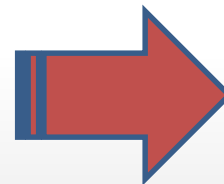
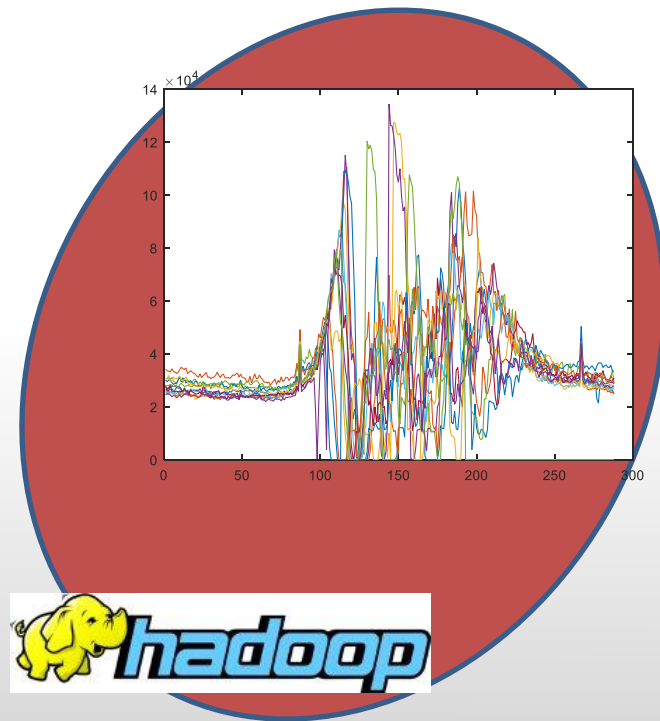
%% Data Analysis

if Lmodel.type==2, % for PLS

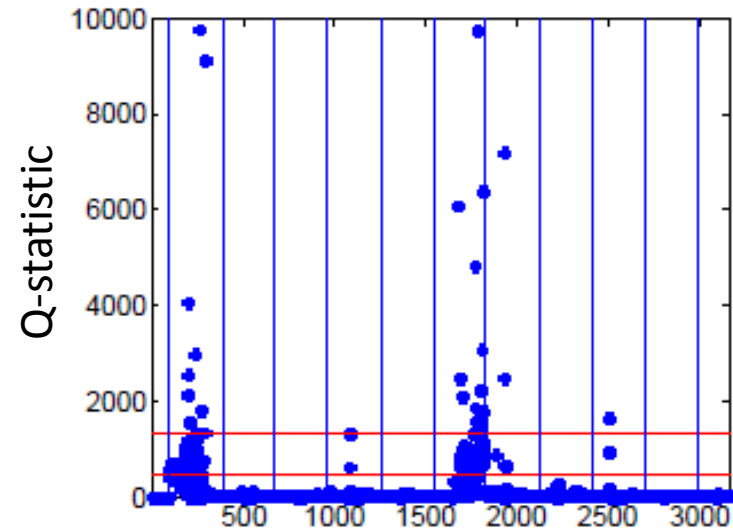
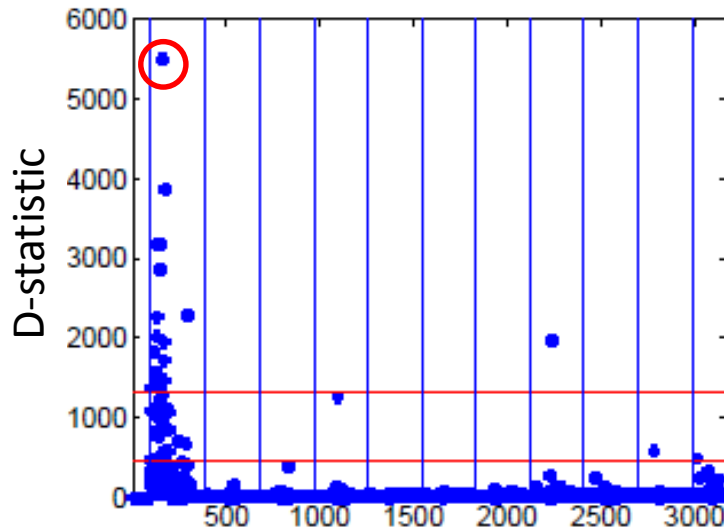
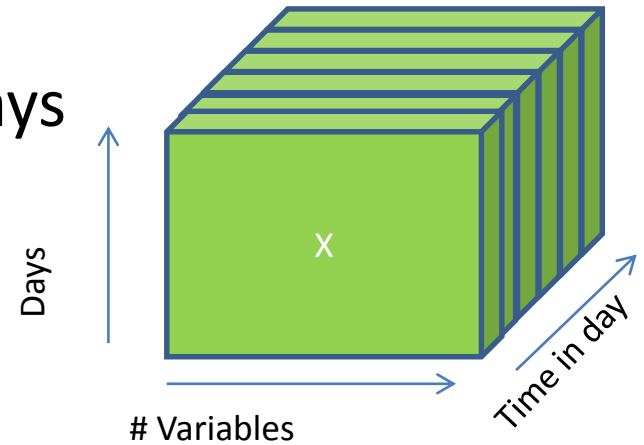
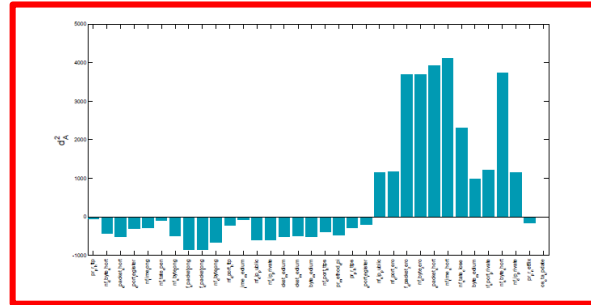
    % Score plot
    scores_Lpls(Lmodel, 1:2);

    % MEDA
    map = meda_Lpls(Lmodel, 1:2, 0, 3);
```

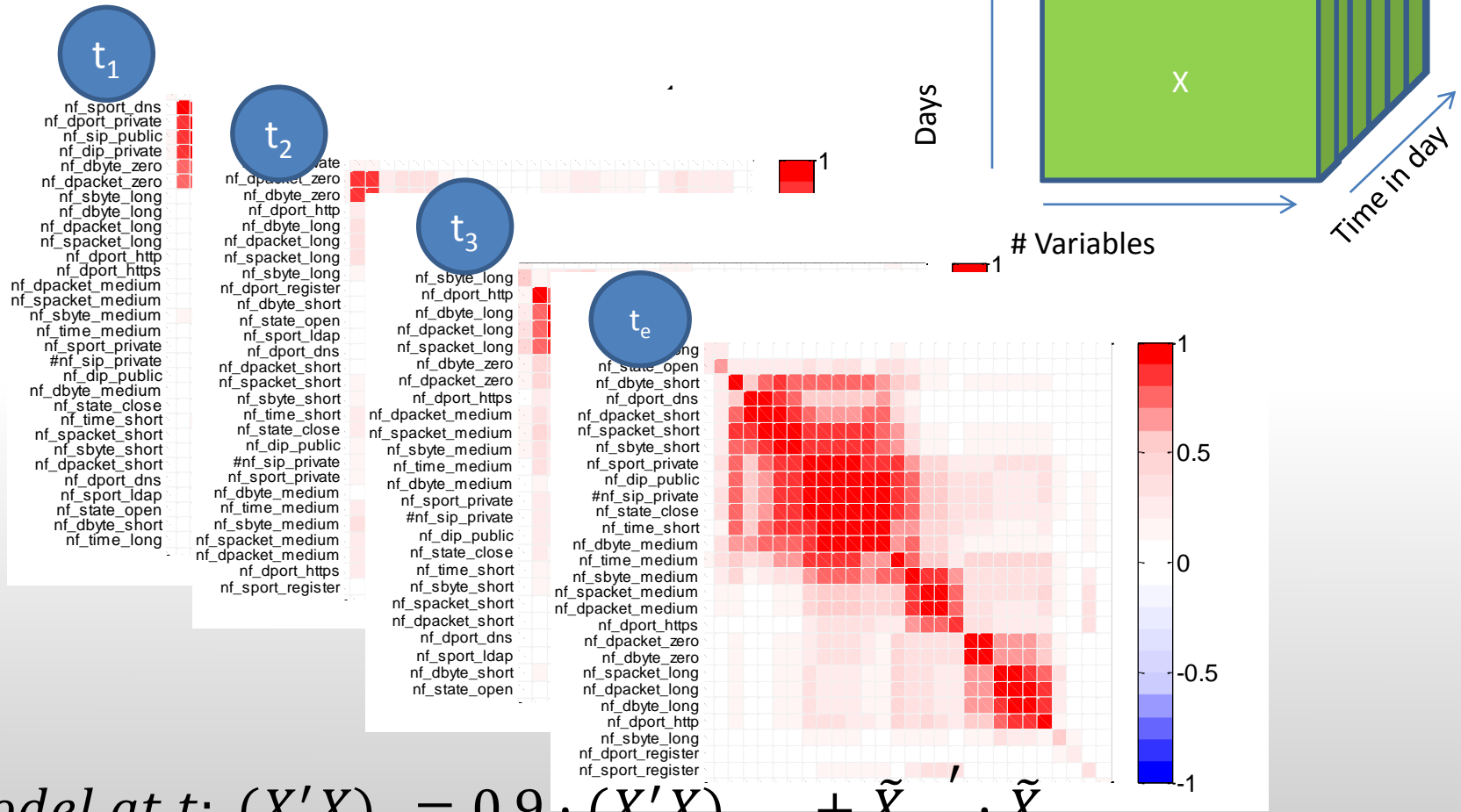

→ Confidential Network: 12 Working days



➔ Confidential Network: 12 Working days

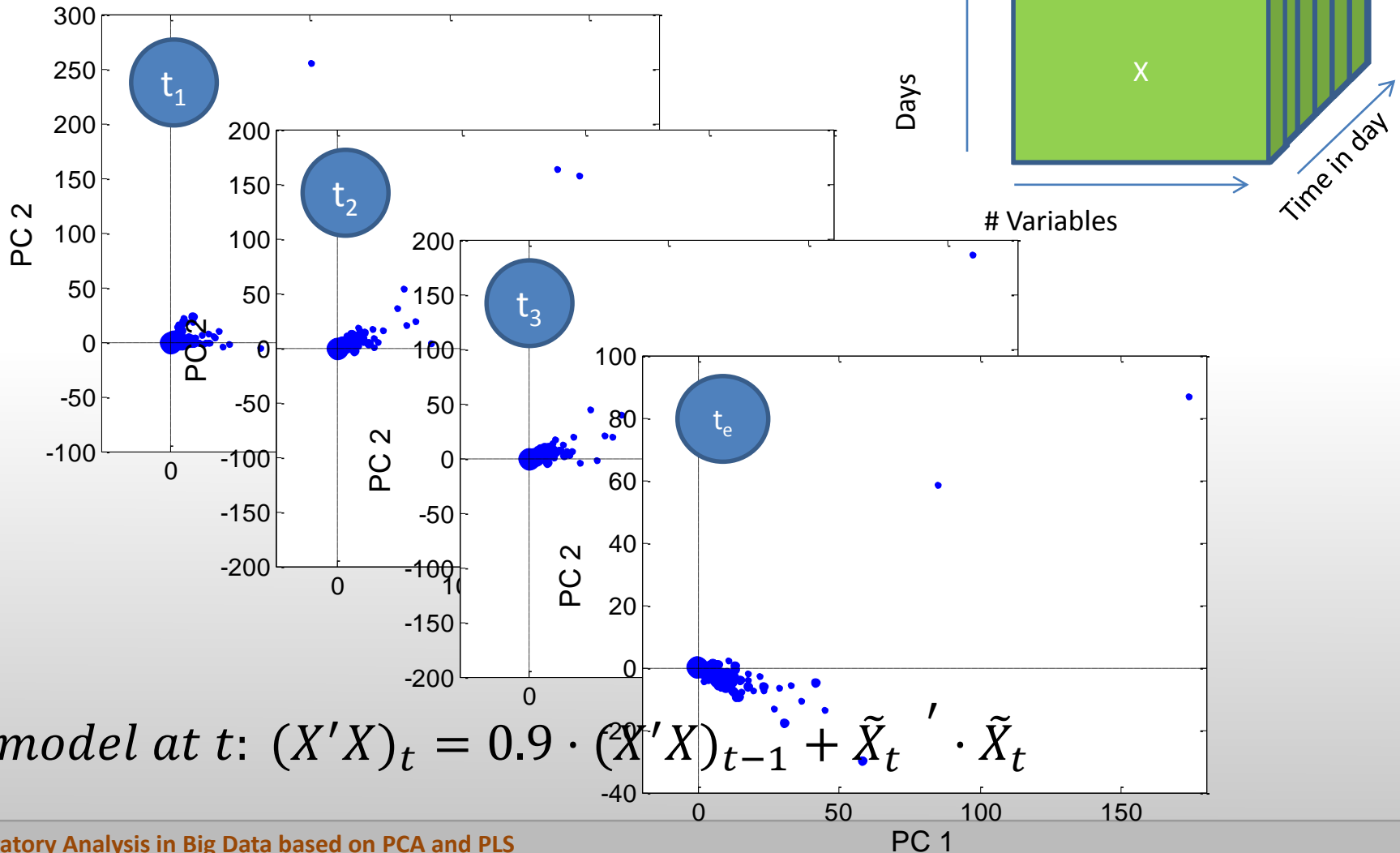


Confidential Network: 12 Working days

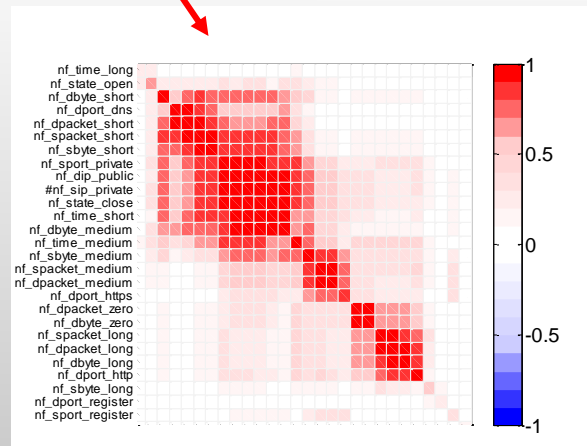
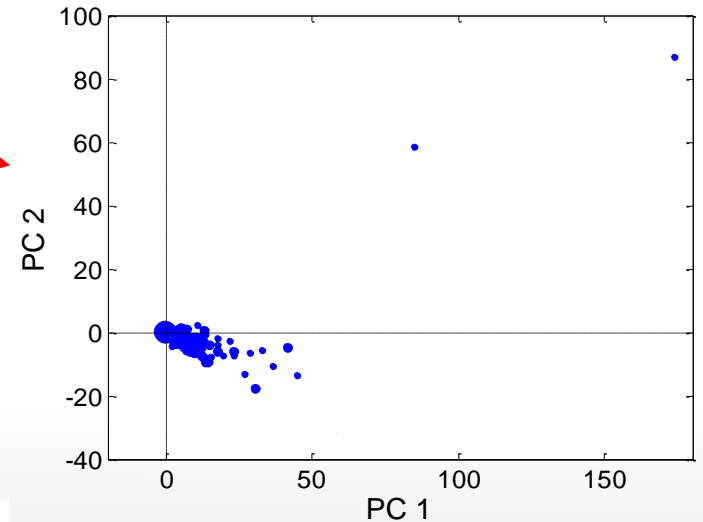
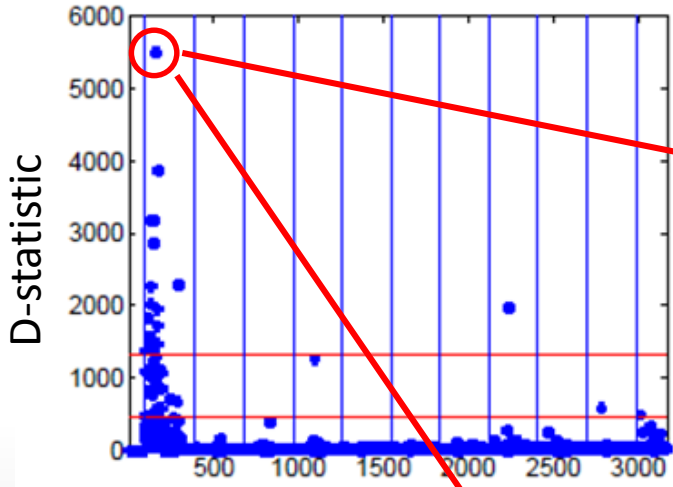


$$Lmodel \text{ at } t: (X'X)_t = 0.9 \cdot (X'X)_{t-1} + \tilde{X}_t' \cdot \tilde{X}_t$$

Confidential Network: 12 Working days



➔ Confidential Network: 12 Working days



- ➔ Multivariate Analysis has a lot to do in Big Data and Visual Analytics

- ➔ Multivariate Analysis tools can be extended to Networking for anomaly detection, optimization, ...
 - ✓ with new and interesting particularities and challenges
 - ✓ with challenges already solved in chemometrics

- ➔ Who should lead this?

Exploratory Analysis in Big Data based on PCA and PLS

José Camacho

Departamento de Teoría de la Señal, Telemática y Comunicaciones



Network Engineering & Security Group
<http://nesg.ugr.es>



UGR | Universidad
de Granada