# Exploratory Analysis on Big Data using the MEDA Toolbox
## Present and Future

**José Camacho,**

**Roberto Therón,**

**Roberto Magán**

*Departamento de Teoría de la Señal, Telemática y Comunicaciones*
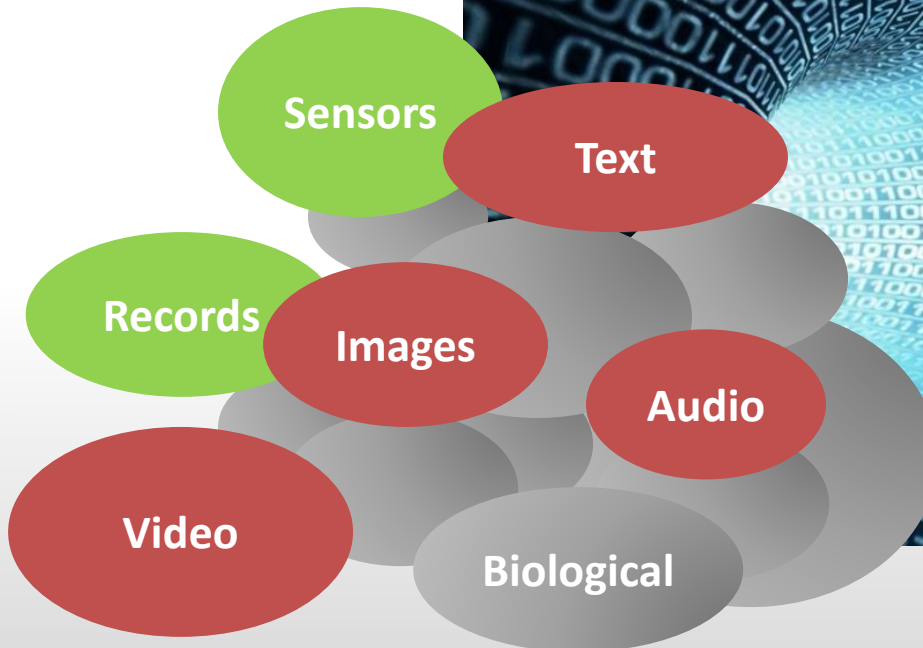
*Universidad de Granada*

MINIARCTIC 2016

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

*ugr* | Universidad de Granada
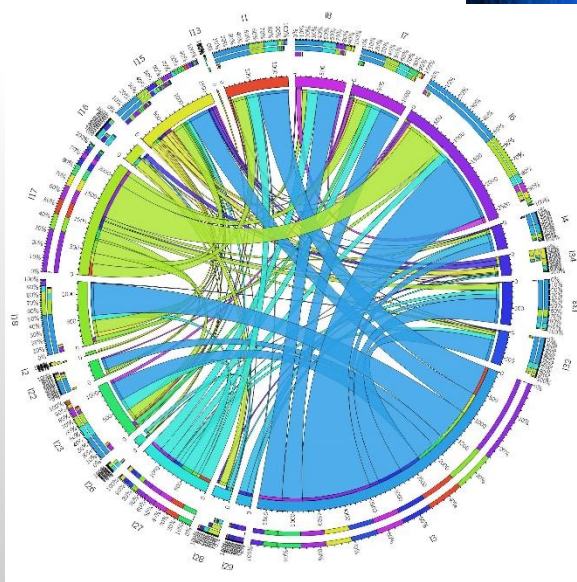
Data-Driven Documents

Visual Analytics

BIG DATA ECOSYSTEM

Parallel Processing

Storage

Cluster of Computers

Spark

mahout

Cassandra

cloudera

hadoop

APACHE HBASE

## MEDA Toolbox

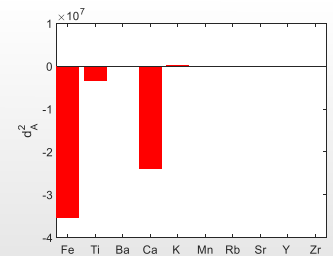https://github.com/josecamachop/MEDA-Toolbox

- ✔ Models: PCA, PLS-DA, SPLS, GPCA, GPLS

- ✔ Dimensionality:
  - Scree plots
  - CV & D-CV
  - 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

- ✔ Data simulation
  - simuleMV



**MATLAB**

MathWorks

ChemoLab, (2015) 143: 49

https://github.com/josecamachop/MEDA-Toolbox

➡ Extensions for Big Data

✔ For variables ➔ (linear) kernel calibration 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, MEDA, oMEDA
- GPCA, GPLS
- ADICOV MSPC

✔ For observations ➔ Clustering

- Scalable to any size
- Compressed Score Plots
- Compressed MSPC



ChemoLab, (2014) 135: 110

```matlab
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);
```

# BIG DATA SETS

```
Lmodel = update_iterative(short_list,'',Lmodel,20,step,0,'',1); % Iterative
```

## MEAN

$$\mathbf{M}_t^x = \mathbf{M}_{t-1}^x + \mathbf{X}_t$$
$$N_t = N_{t-1} + B_t$$

$$\mathbf{m}_t^x = (1/N_t) \cdot \mathbf{M}_t^x$$

## SCALE

$$(\sigma_t^x)^2 = (\sigma_{t-1}^x)^2 +$$
$$\sum_{i=1}^{B_t} \left(\mathbf{x}_t^i - \mathbf{m}_t^x\right)^2$$

$$\sigma_t^x = \sqrt{(1/(N_t-1)) \cdot (\sigma_t^x)^2}$$

## CROSS-PROD

$$\tilde{\mathbf{x}}_t^i = \left(\mathbf{x}_t^i - \mathbf{m}_t^x\right) \oslash \sigma_t^x$$

$$\mathbf{XX}_t = \mathbf{XX}_{t-1} + \tilde{\mathbf{X}}_t^T \cdot \tilde{\mathbf{X}}_t$$

ChemoLab, (2014) 135:110

PCA(ED)

PLS(XX,XY)

MEDA, Loading plots, CV, ...

# BIG DATA SETS

```
Lmodel = update_iterative(short_list,'',Lmodel,20,step,0,'',1); % Iterative
```

## Compressed Scores

$$\mathbf{C} \leftarrow Cluster(\mathbf{X}, \mathbf{K}^{-1}):$$

$$\mathbf{C} = []$$

$$\boldsymbol{\mu} = []$$

$$[\mathbf{X}_1, ..., \mathbf{X}_T] \leftarrow partition(\mathbf{X})$$

for each packet $\mathbf{X}_t$,

$$\mathbf{C} \leftarrow [\mathbf{C}, \mathbf{X}_t]$$

$$\boldsymbol{\mu} \leftarrow [\boldsymbol{\mu}, 1_{B_t}]$$

$$[\mathbf{C}, \boldsymbol{\mu}] \leftarrow merge(\mathbf{C}, \boldsymbol{\mu}, \mathbf{K}^{-1})$$

end

ChemoLab, (2014) 135:110

# BIG DATA SETS

```
Lmodel = update_iterative(short_list,'',Lmodel,20,step,0,'',1); % Iterative
```
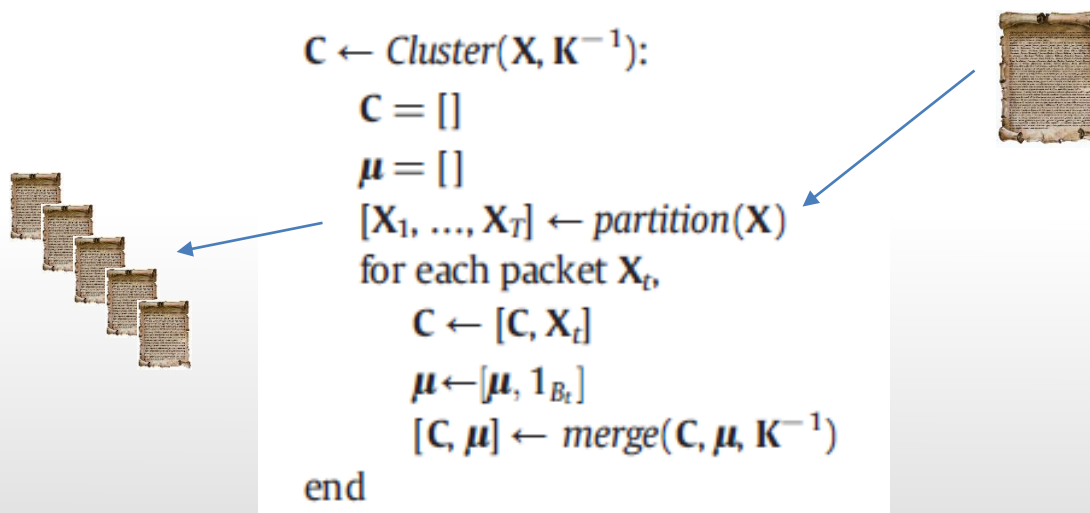
## Compressed Scores

$C \leftarrow Cluster(\mathbf{X}, \mathbf{K}^{-1})$:

   $\mathbf{C} = []$

   $\boldsymbol{\mu} = []$

   $[\mathbf{X}_1, ..., \mathbf{X}_T] \leftarrow partition(\mathbf{X})$

   for each packet $\mathbf{X}_t$,

      $\mathbf{C} \leftarrow [\mathbf{C}, \mathbf{X}_t]$

      $\boldsymbol{\mu} \leftarrow [\boldsymbol{\mu}, 1_{R_t}]$

      $[\mathbf{C}, \boldsymbol{\mu}] \leftarrow \boxed{merge(\mathbf{C}, \boldsymbol{\mu}, \mathbf{K}^{-1})}$

end

$[\mathbf{C}, \boldsymbol{\mu}] \leftarrow merge(\mathbf{C}, \boldsymbol{\mu}, \mathbf{K}^{-1})$:

   $L := \# (\mathbf{C})$

   $\mathbf{C} := [\mathbf{c}_1, ..., \mathbf{c}_L]$

   $\boldsymbol{\mu} := [\mu_1, ..., \mu_L]$

   while $(\# (\mathbf{C}) > L_{end})$,

      $[\mathbf{c}_i, \mathbf{c}_j] \leftarrow min\_dist(\mathbf{C}, \mathbf{K}^{-1})$

      $\mathbf{c}_i \leftarrow centroid(\mu_i \cdot \mathbf{c}_i, \mu_j \cdot \mathbf{c}_j)$

      $\mathbf{C} \leftarrow [\mathbf{c}_1, ..., \mathbf{c}_{j-1}, \mathbf{c}_{j+1}, ..., \mathbf{c}_L]$

      $\mu_i \leftarrow \mu_i + \mu_j$

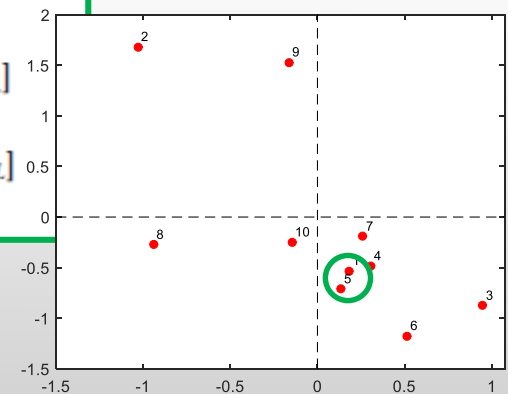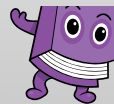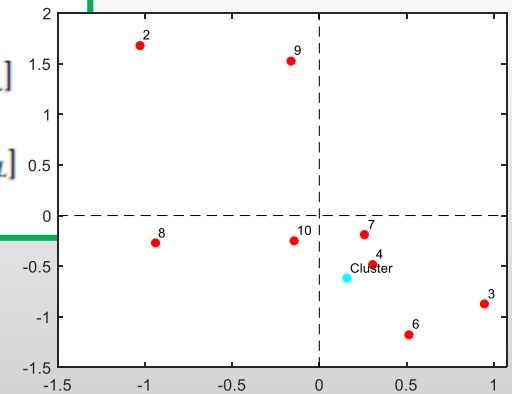      $\boldsymbol{\mu} \leftarrow [\mu_1, ..., \mu_{j-1}, \mu_{j+1}, ..., \mu_L]$

end



ChemoLab, (2014) 135:110

# BIG DATA SETS

```
Lmodel = update_iterative(short_list,'',Lmodel,20,step,0,'',1); % Iterative
```

## Compressed Scores

$C \leftarrow Cluster(\mathbf{X}, \mathbf{K}^{-1})$:

$\mathbf{C} = []$

$\boldsymbol{\mu} = []$

$[\mathbf{X}_1, ..., \mathbf{X}_T] \leftarrow partition(\mathbf{X})$

for each packet $\mathbf{X}_t$,

$\mathbf{C} \leftarrow [\mathbf{C}, \mathbf{X}_t]$

$\boldsymbol{\mu} \leftarrow [\boldsymbol{\mu}, \mathbf{1}_{R_t}]$

$[\mathbf{C}, \boldsymbol{\mu}] \leftarrow \boxed{merge(\mathbf{C}, \boldsymbol{\mu}, \mathbf{K}^{-1})}$

end

$[\mathbf{C}, \boldsymbol{\mu}] \leftarrow merge(\mathbf{C}, \boldsymbol{\mu}, \mathbf{K}^{-1})$:

$L := \#(\mathbf{C})$

$\mathbf{C} := [\mathbf{c}_1, ..., \mathbf{c}_L]$

$\boldsymbol{\mu} := [\mu_1, ..., \mu_L]$

while $(\#(\mathbf{C}) > L_{end})$,

$[\mathbf{c}_i, \mathbf{c}_j] \leftarrow min\_dist(\mathbf{C}, \mathbf{K}^{-1})$

$\mathbf{c}_i \leftarrow centroid(\mu_i \cdot \mathbf{c}_i, \mu_j \cdot \mathbf{c}_j)$

$\mathbf{C} \leftarrow [\mathbf{c}_1, ..., \mathbf{c}_{j-1}, \mathbf{c}_{j+1}, ..., \mathbf{c}_L]$

$\mu_i \leftarrow \mu_i + \mu_j$

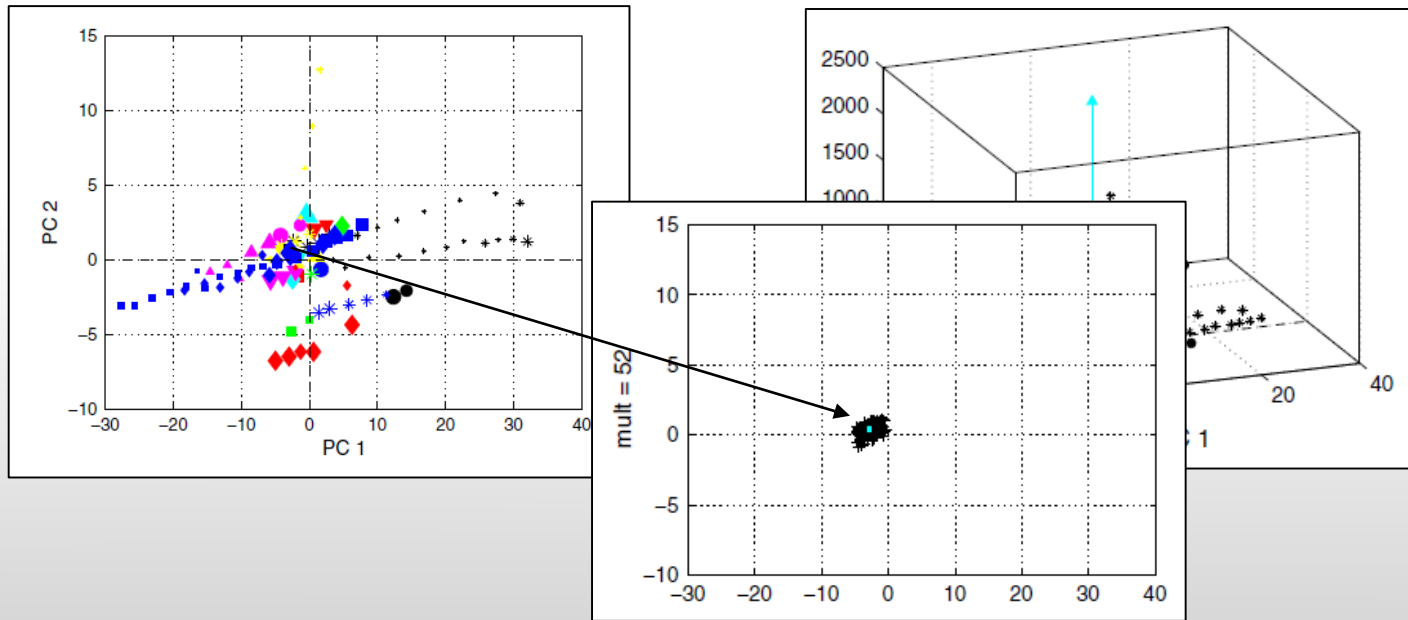$\boldsymbol{\mu} \leftarrow [\mu_1, ..., \mu_{j-1}, \mu_{j+1}, ..., \mu_L]$

end



ChemoLab, (2014) 135:110

# BIG DATA SETS

```
Lmodel = update_iterative(short_list,'',Lmodel,20,step,0,'',1); % Iterative
```

## Compressed Scores

$\mathbf{C} \leftarrow Cluster(\mathbf{X}, \mathbf{K}^{-1})$:

$\quad \mathbf{C} = []$

$\quad \boldsymbol{\mu} = []$

$\quad [\mathbf{X}_1, ..., \mathbf{X}_T] \leftarrow partition(\mathbf{X})$

$\quad$ for each packet $\mathbf{X}_t$,

$\quad\quad \mathbf{C} \leftarrow [\mathbf{C}, \mathbf{X}_t]$

$\quad\quad \boldsymbol{\mu} \leftarrow [\boldsymbol{\mu}, 1_{R_t}]$

$\quad\quad [\mathbf{C}, \boldsymbol{\mu}] \leftarrow merge(\mathbf{C}, \boldsymbol{\mu}, \mathbf{K}^{-1})$

end

$[\mathbf{C}, \boldsymbol{\mu}] \leftarrow merge(\mathbf{C}, \boldsymbol{\mu}, \mathbf{K}^{-1})$:

$\quad L := \#(\mathbf{C})$

$\quad \mathbf{C} := [\mathbf{c}_1, ..., \mathbf{c}_L]$

$\quad \boldsymbol{\mu} := [\mu_1, ..., \mu_L]$

$\quad$ while $(\#(\mathbf{C}) > L_{end})$,

$\quad\quad [\mathbf{c}_i, \mathbf{c}_j] \leftarrow min\_dist(\mathbf{C}, \mathbf{K}^{-1})$

$\quad\quad \mathbf{c}_i \leftarrow centroid(\mu_i \cdot \mathbf{c}_i, \mu_j \cdot \mathbf{c}_j)$

$\quad\quad \mathbf{C} \leftarrow [\mathbf{c}_1, ..., \mathbf{c}_{j-1}, \mathbf{c}_{j+1}, ..., \mathbf{c}_L]$

$\quad\quad \mu_i \leftarrow \mu_i + \mu_j$

$\quad\quad \boldsymbol{\mu} \leftarrow [\mu_1, ..., \mu_{j-1}, \mu_{j+1}, ..., \mu_L]$

end

$d_{\mathbf{K}}\left(x_i, x_j\right) = \left\|x_i - x_j\right\|_{\mathbf{K}} =$
$= \left(\left(\mathbf{x}_i - \mathbf{x}_j\right)^T \mathbf{K}^{-1}\left(\mathbf{x}_i - \mathbf{x}_j\right)\right)^{1/2}$

$\mathbf{K}^{-1} = \mathbf{P} \cdot \mathbf{K}_{PCA}^{-1} \cdot \mathbf{P}^T,$

$\mathbf{K}^{-1} = \mathbf{R} \cdot \mathbf{K}_{PLS}^{-1} \cdot \mathbf{R}^T$

ChemoLab, (2014) 135:110

# BIG DATA SETS

```
Lmodel = update_iterative(short_list,'',Lmodel,20,step,0,'',1); % Iterative
```
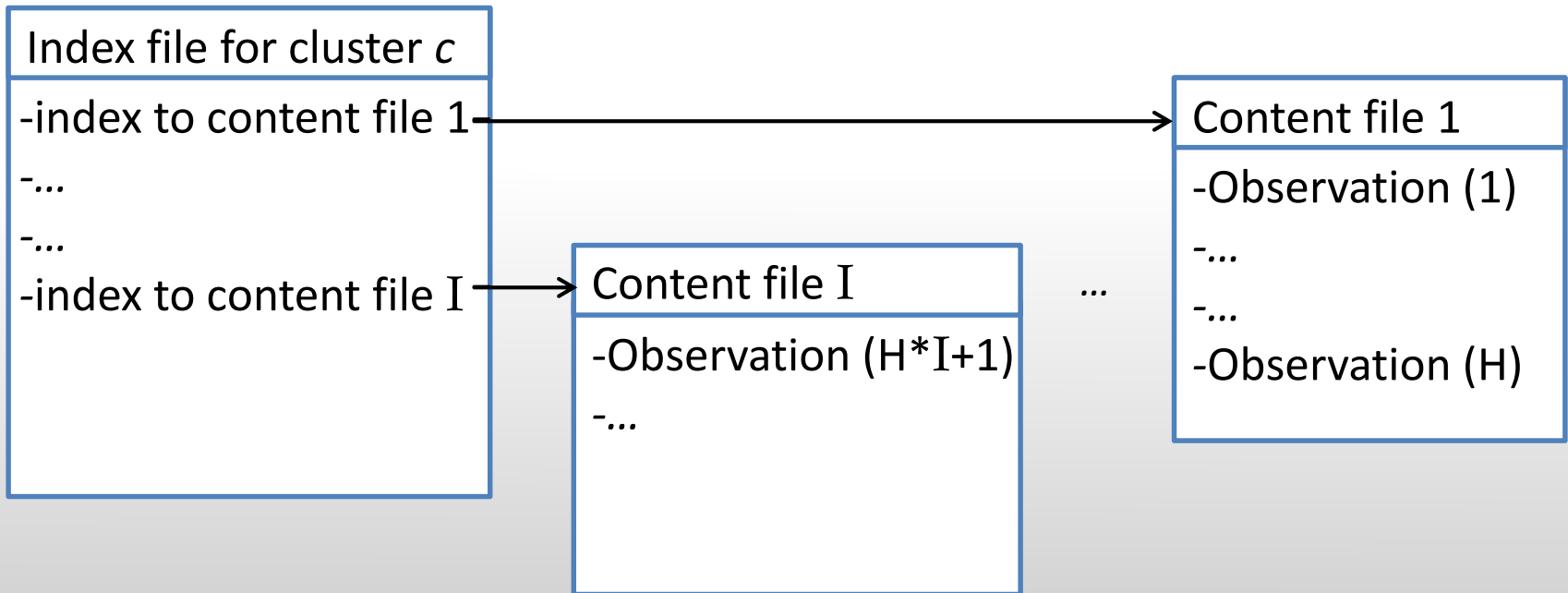
## Compressed Score Plot (CSP)
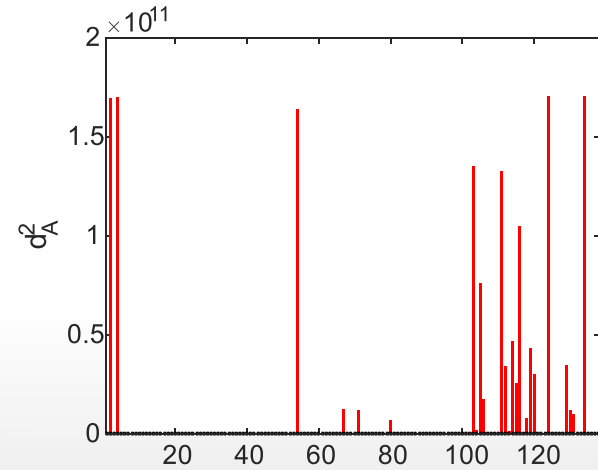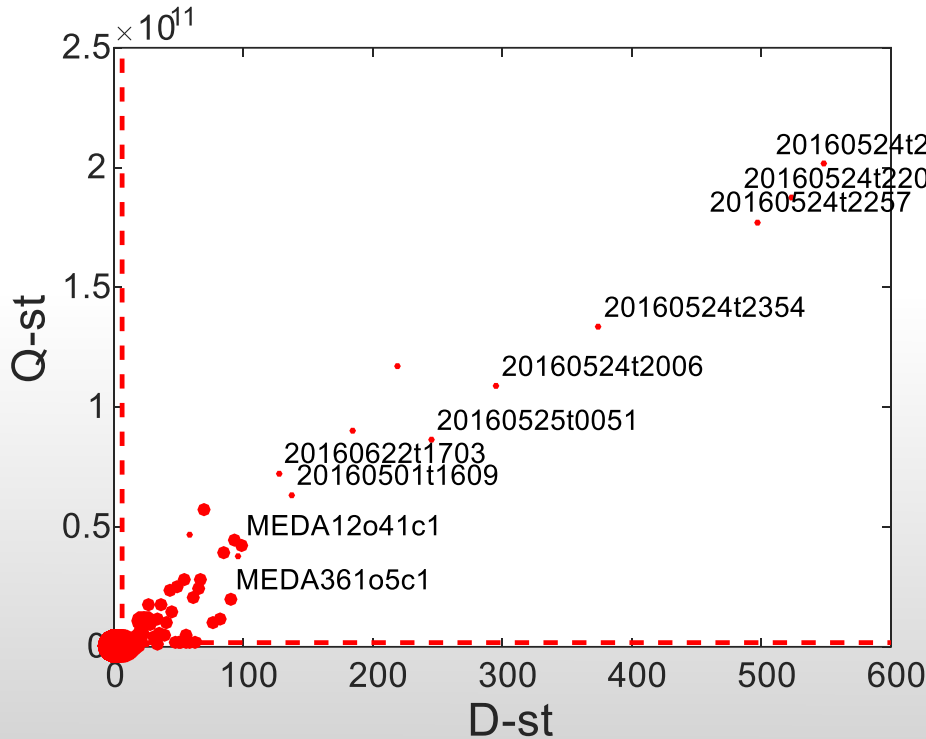


ChemoLab, (2014) 135:110

# BIG DATA SETS

```
Lmodel = update_iterative(short_list,'',Lmodel,20,step,1,'./output/',1); % Iterative
```

## Compressed Score Plot (CSP)

| Index file for cluster $c$ |
|---|
| -index to content file 1 |
| -... |
| -... |
| -index to content file I |

Content file 1
- -Observation (1)
- -...
- -...
- -Observation (H)

Content file I
- -Observation (H*I+1)
- -...

...

# BIG DATA SETS

## Example: vw PCA-MSPC in Big Data (Networkmetrics)



SPAM ATTACK

# Data Streams

```
Lmodel = update_ewma(short_list,'',Lmodel,lambda,step,1); % EWMA
```

$$\mathbf{M}_t^x = \lambda \cdot \mathbf{M}_{t-1}^x + \mathbf{X}_t$$

$$\mathbf{m}_t^x = (1/N_t) \cdot \mathbf{M}_t^x$$

$$N_t = \lambda \cdot N_{t-1} + B_t$$

$$(\sigma_t^x)^2 = \lambda \cdot (\sigma_{t-1}^x)^2 + \sum_{i=1}^{B_t} \left(\mathbf{x}_t^i - \mathbf{m}_t^x\right)^2$$

$$\sigma_t^x = \sqrt{(1/(N_t-1)) \cdot (\sigma_t^x)^2}$$

$$\tilde{\mathbf{x}}_t^i = \left(\mathbf{x}_t^i - \mathbf{m}_t^x\right) \oslash \sigma_t^x$$

$$\mathbf{XX}_t = \lambda \cdot \mathbf{XX}_{t-1} + \tilde{\mathbf{X}}_t^T \cdot \tilde{\mathbf{X}}_t$$

**PCA(ED)**

**PLS (XX,XY)**

**MEDA,
Loading plots,
CV, ...**

J.P.C., (1997) 7:169
ChemoLab, (2014) 135:110

# Data Streams



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

➡ Data mining / Machine learning

➡ Chemometrics / Exploratory Data Analysis

INTERACTION

✓ Cross-validation, …
✓ Variable Selection
✓ Outlier Identification
✓ …

➡ EDA + Visual Analytics

FULL INTERACTION

OCTAVE

iMEDA

INTERNET

SMALL DATA

CITIC-UGR

Web Server

Web Frontend

Data-Driven Documents

iMEDA Dashboard 1.1

Interactive visualization for EDA using the MEDA-Toolbox.

Multivariate Software

python

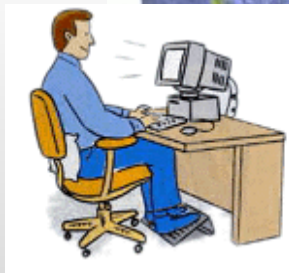Middelware

Spark

Cluster of Computers

Data Pathfinder

INTERNET

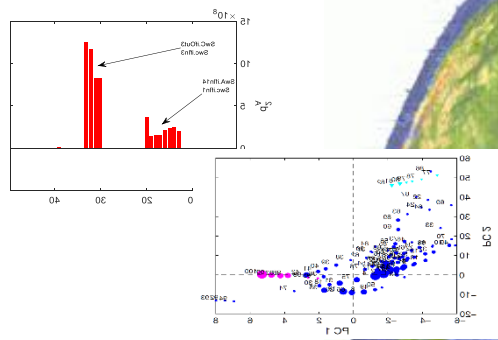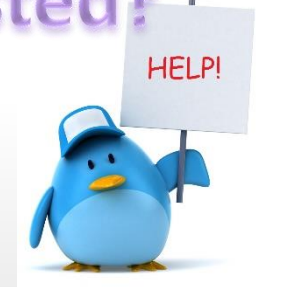BIG DATA

CITIC-UGR

Is this the way
for MA + BD?

Interested?

HELP!

IS IT TO ACHIEVE
FULL DATA INTERACTIVITY?

# Exploratory Analysis on Big Data using the MEDA Toolbox
## Present and Future

**José Camacho,**

**Roberto Therón,**

**Roberto Magán**

*Departamento de Teoría de la Señal, Telemática y Comunicaciones*

*Universidad de Granada*

MINIARCTIC 2016

*ugr* | Universidad de Granada