# Using ML clustering tools to improve data transfer management operations

L Rinaldi[1] , L Clissa[2], L Morganti[3]

*1)Bologna University and INFN, 2)INFN-Bologna, 3)INFN-CNAF*

The **GRID** computing paradigms adopted by the main HEP experiments are based on the **distribution** of experimental data on computer resources located all over the world

All the **data transfer** processes are tracked in **log files** produced by the various services involved and such log files represent a source of **information** which is largely **underutilized**

An approach based on **unsupervised ML** techniques is used to automatically process information stored in **log files** with the aim of grouping the error messages and speed up the procedures for **detecting errors** and **solving problems**

## Clustering pipeline

**Step 1: text pre-processing**
- Lowercase transformation and punctuation/stopword stripping
- Tokenization
- URL split

**Step 2: vectorization**
- Transformation of the pre-processed text into numeric information to map each message to a point in a vectorial subspace (embedding)
- *word2vec* language model adopted

**Step 3: clustering**
- k-means++ algorithm: intuitive approach and good performance in a wide range of applications
- The number of clusters $k$ is set $k \in [12, 15, 20, 30]$ at each clustering stage based on a grid search and geometrical criteria:
  - Within cluster Sum of Squared Errors (elbow method)

$$\text{WSSE}(\text{dist}, k) = \sum_{j=1}^{k} \sum_{x_i \in C_j} \text{dist}(x_i - \bar{x}_j)$$ ,

$x_i$ generic data point
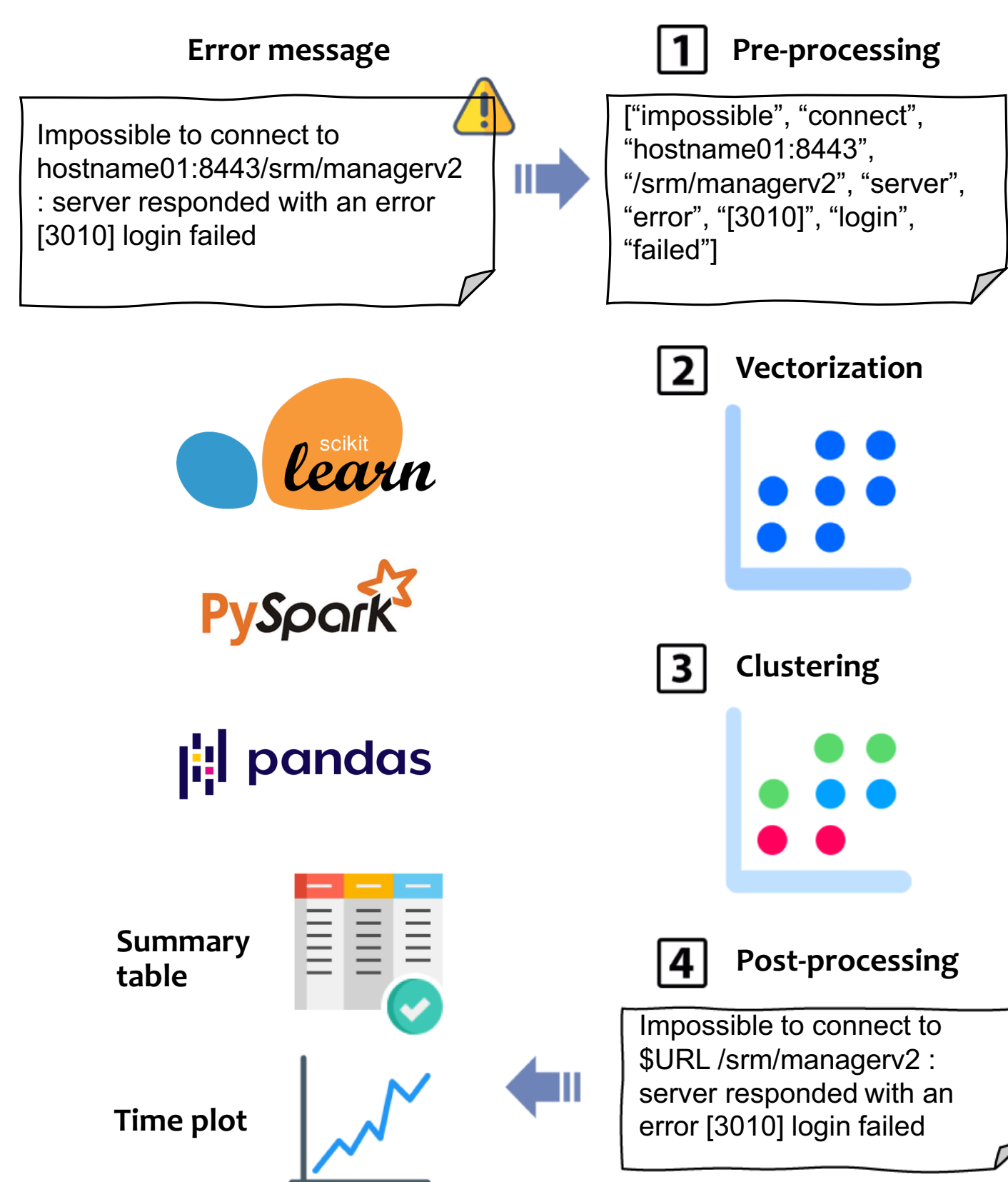$\bar{x}_j$ centroid of a generic cluster $C_j$

  - Average Silhouette Width (max)

$$\text{ASW}(\text{dist}, k) = \frac{1}{n}\sum_{i=1}^{n} \frac{b_i - \bar{a}_i}{\max(\bar{a}_i, b_i)}$$

$\bar{a}_i$ distance of $x_i$ from same cluster points
$b_i$ distance of $x_i$ from other clusters points

**Step 4: post-processing**
- Organize results according to the use case



## Analysis of File Transfer Service (FTS) errors



**Cluster Description** (tabular format)

Analysis of the error messages in FTS log files (1 day)

The three most frequent triplets of <pattern>-<source>-<destination> reported in descending order for each cluster
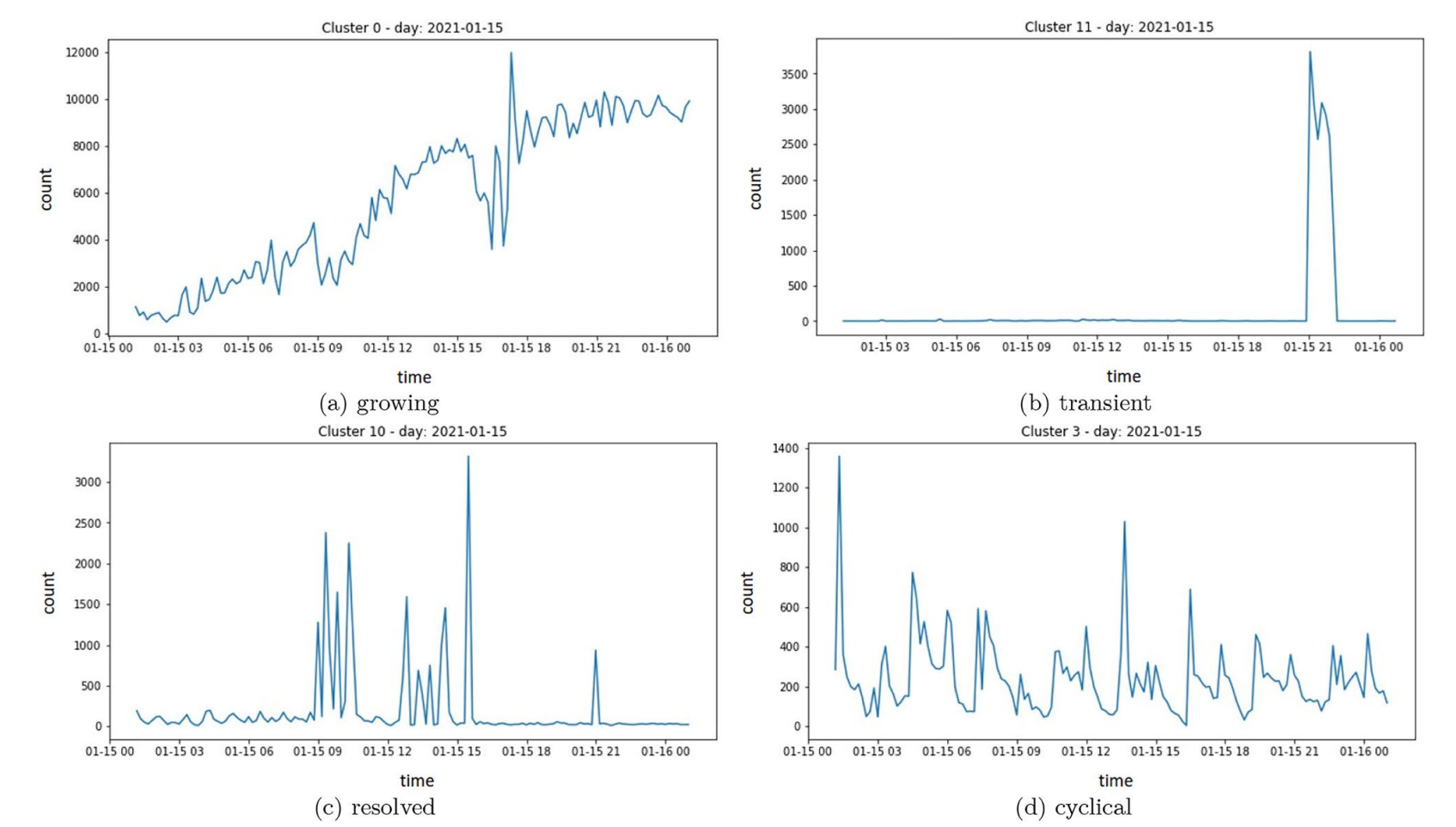
Precious insights for spotting:
- what type of errors and where they occur
- amount of errors, both absolute and relative to the error group

**Time Series**

Temporal trend of the number of errors generated by each cluster
- escalating or cyclical failures →require immediate actions
- transient or in resolution



| N. clusters | ASW | WSSE | Perfect match | Fuzzy match | Partial match | False positives | False negatives |
|---|---|---|---|---|---|---|---|
| 15 | 0.89 | 17107 | 7 | 3 | 2 | 3 | 1 |

Summary of the cross-check between clusters and incidents reported in GGUS. Most of the groups discovered are linked to reported issues, with only 3 false positives and 1 false negative

Extensive testing using incidents reported in **GGUS** as a benchmark (17 days window)
- overlapping between discovered clusters and the reported issues

*Computing and Software for Big Science (2022) 6:16 https://doi.org/10.1007/s41781-022-00089-z*

## Analysis of StoRM log files

Analysis performed on a 7 days interval Dataset (daily log files) from INFN-CNAF StoRM system

Messages of a single log file are grouped by StoRM process (separately for explicit *error* messages)  and then passed to the clustering pipeline

Higher values of the optimal number of clusters may give a hint of an anomaly (day 5 in this example)



**Optimal number of clusters (daily)**



**Clusters on day 1**



**Clusters on day 5**



The variation of the daily number of clusters and clustered error messages could give more information on the cause of the problem