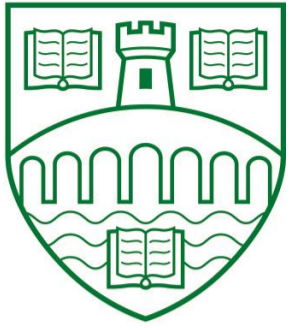# Failure Prediction in Cluster Systems

Muhammed Costan

MSc in Big Data

## UNIVERSITY *of* STIRLING

**Servers** create unstructured syslogs in different formats. Syslog messages usually include information about where, when, and why the event was logged: IP address, timestamp, and the actual log message are the minimal information kept.

**logstash**

Message, timestamp, server ID, severity (the importance level)

**Logstash** is a powerful pipelining tool which collects, parses, enriches and transports any log data from any source. It can load the processed logs into targets: CSV, ES, email and others.

**elasticsearch**

All the structured logs sent by logstash are stored on **elasticsearch** which is a document-oriented database mainly developed for storing all kind of logs.

**kibana**

**Kibana** is used to perform basic analytics and visualisations on syslogs stored on elasticsearch.

The ELK Stack had already been implemented by the company before the project.

## SVM

Data is pre-processed on spark by using **python** and prepared for applying a machine learning algorithm. After pre-processing, the compressed data is exported into a file and the **SVM** algorithm is applied to this file for training a model. SVM performs classification on two class labels (failure or not failure) by maximising the distance between two classes. As a machine learning software, SVM-Light has been chosen.

**python**

**Spark**

**ES hadoop**

Required fields on ES are extracted via **es-hadoop** and fetched into an RDD on **Spark.** Es-hadoop is a two way connector between two technologies. RDD is a resilient and distributed data structure allowing us to perform parallel computations (in memory) on data.

## Model

**Model** file is created by the algorithm.

## Predictions

An extra job is created for retrieving the final time window containing the latest sequence on ES. The same data preparation steps are applied on this sequence. Next, by using the model created, a **prediction** is made on the pre-processed sequence. The job can be scheduled to run in predefined periods and once it predicts a failure, it can raise a notification.
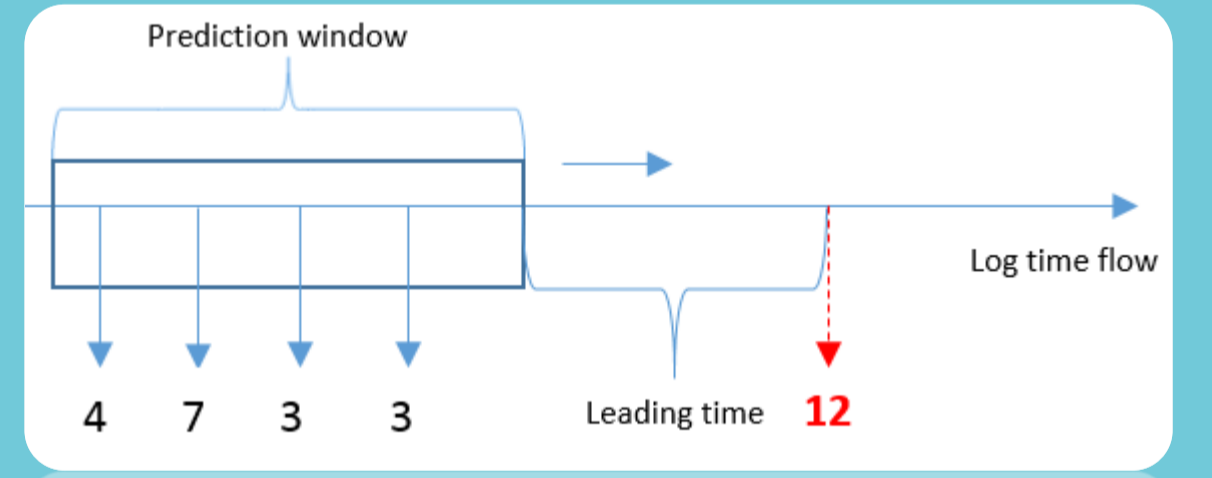
**About The SVM Application:**
- SVM is applied by changing the values of different parameters (like kernel) with the cross validation method and the model with the highest accuracy is chosen.
- Illustration for linear kernel:



- Results of a model trained on a sample dataset with 1506 data points (the accuracy rate is 78.29%):

| Act/Pred | Failure | Non-Failure |
|---|---|---|
| **Failure** | 549 | 204 |
| **Non-Failure** | 123 | 630 |

**Data Pre-processing Steps:**
- Data Cleansing: All letters are transformed into lower case letters, log-specific information (date, IP address…) and all the characters except letters are removed by using regular expressions, duplicate and **redundant*** events are eliminated. For example:
  `SNMP v2 get failed to "123.456.123.456"` turns into `snmp get failed to`
- Clustering Log Messages: Assigning the logs to clusters represented by integers. Clustering the messages is made based on the **edit distance**** between two different messages.
- Cleaning Noises: If a cluster contains insufficient number of log elements, these clusters are removed to fix the distribution in dataset.
- Sequence Extraction: A sequence is a group of logs leading to either a failure or not failure log message. So, after cleaning noises, all the logs are sorted by date and by using time window logic sequences belonging to two different classes are extracted.
- **Time Window:** It represents how many minutes we should go back from a point of failure to extract a sequence. While, **prediction window** is used in machine learning process, the events within **leading time** window are not included to the process to take precautions before a failure:



- Feature Selection: Sequences should be converted into a **frequency matrix*** to be able to run a machine learning algorithm on these sequences.

*Redundant Events: Some machines log the same events more than once within small time intervals. These events are considered as redundant events.

**Edit Distance: It is a way of quantifying how dissimilar two strings are. It is calculated by counting the minimum number of operations required to transform one string into the other.

***Frequency Matrix: It contains the number of occurrences of elements within a sequence. (For example; ClusterID : Frequency) It is usually used in sequential classification problems.

**Pulsant** Business Unlimited