



Deliverable 1. Problem Scope.

Description of the Project Context and Definition of Learning Goals.

Grupo de Sistemas Inteligentes

Departamento de Ingeniería de Sistemas Telemáticos

Universidad Politécnica de Madrid.

Project Report

Madrid, September 2012

Authors:

Adrián Pérez Orozco

Álvaro Carrera Barroso

Carlos A. Iglesias Fernández

Executive Summary

This document describes the scope and potential outcomes of the project *Trainmining*. This project is developed by Thales Spain in collaboration with the research group *Grupo de Sistemas Inteligentes* (Intelligent Systems Group) of the Universidad Politécnica de Madrid. The goal of the project is the application of machine learning techniques for predicting alarms in the maintenance systems that have been developed by Thales Spain. Machine learning techniques will be based on historical data coming from Thales Spain railway maintenance system which operates in several railway lines throughout Spain.

The document provides a general overview of the railway maintenance system developed by Thales Spain, as well as a first insight in the available historical data and how it will be treated. In addition, the document describes potential learning goals for the application of machine learning techniques. Finally, the document describes which are the learning goals Thales Spain is more interested in, and details potential machine learning techniques to achieve them. The main learning goal is learning alarm prediction techniques, and their usage is illustrated with a use case.

Contents

Executive Summary	i
Contents	ii
List of Figures	iii
List of tables	iv
1 Introduction	1
2 Database description	1
3 Reduced representation of alarms	7
4 Statistic analysis	8
4.1 Alarm classification	8
4.1.1 Hourly timeline	9
4.1.2 Daily correlation	12

List of Figures

1	Simplified diagram of the maintenance systems architecture . .	6
2	Alarm information for Antequera	8
3	Alarm information for Segovia	9
4	Alarm information for Sevilla	9
5	Hourly distribution for Antequera (stacked)	10
6	Hourly distribution for Segovia (stacked)	10
7	Hourly distribution for Sevilla (stacked)	11
8	Daily correlation for Antequera	12
9	Daily correlation for Segovia	13
10	Daily correlation for Sevilla	13

List of Tables

1	Detail of fields on table ER_ERRORS	2
2	Detail of fields on table IG_INSTALLATIONGENERAL . . .	3
3	Detail of fields on table IG_NODO_INSTALLATION	3
4	Detail of fields on table ERS_ERRORS_SAM_ENCE	4
5	Description of values for the field EVENT_TYPE	5

1 Introduction

2 Database description

In order to properly process the data provided in form of database backups, it is of essential importance that we completely understand how data is represented in databases. We will analyse the structure and how data is represented in the provided databases: Antequera, Segovia and Sevilla. Each of these database corresponds to a single *maintenance station*, which comprises a whole railway line with several elements along it. The elements with diagnosis systems which can raise alarms are called *installations*, and have different sets of sensors and other systems to control *field elements*. An schematic representation of this architecture is represented in figure 1. The detailed description of available systems and subsystems is of few interest to us. Initially we will only need to differentiate between maintenance stations and installations.

Each *maintenance station* has its own unique database, which is of great convenience in order to treat different stations independently. We will start analysing the structure of the main tables of said databases. Due to the high complexity of the maintenance stations, there are a vast amount of tables with configuration parameters and other operational values which are not of interest for our purposes. With assistance from Thales engineers, we have reduced the tables only to those which characterise registered alarms. A total of 4 different tables is used in order to register this information, which are the following:

Table ER_ERRORS This table contains an entry for every alarm received by the maintenance station. Its fields are detailed on table 1.

Table IG_INSTALLATIONGENERAL This table contains information on all the installations managed by the maintenance station. Its fields are detailed on table 2

Table IG_NODO_INSTALLATION This table gathers additional information on installations which are nodes. Nodes are installations which

Field name	Description
DVNI_ERRORNUMBER	Alarm identifier
DVNS_ERRORTIME	Time-stamp for the alarm
DVNI_INSTALLATIONCODE	Code of the installation in which the alarm was raised
DVNI_SENDERINSTALLATIONCODE	Code of the installation from which the alarm was sent (might be different from the one which raised it)

Table 1: Detail of fields on table ER_ERRORS

can raise alarms but need a parent installation to send them to the maintenance station. Its fields are detailed on table 3

ERS_ERRORS_SAM_ENCE This table contains detailed information about the alarms. Its fields are detailed on table 4

ERH_ERRORS_HSL1 This table is equivalent to ERS_ERRORS_SAM_ENCE. Maintenance stations use one or the other depending on how they receive the alarms. Its only difference with ERS_ERRORS_SAM_ENCE is that registers the method used to receive the alarm. For our purposes it will be treated exactly as its equivalent, and therefore its structure can also be reviewed in table 4

Concluding, for each alarm we will have a timestamp and an alarm identifier in table ER_ERRORS. Alarm identifier is a foreign key which points to table ERS_ERRORS_SAM_ENCE (or equivalent) in which further details of the alarm are saved. Among these details, we can find an installation identifier which specifies which installation has produced the alarm. That identifier is also a foreign key pointing to table DVNI_INSTALLATIONCODE, in which further details about the installation are stored. Further details on all the database fields are given in tables 1, 2, 3 and 4.

Field name	Description
DVNI_INSTALLATIONCODE	Installation identifier
DVNI_SYSTEMCODE	Type of system, as defined in the “SG_SYSTEMSGENERAL” table
DVNI_VERSION	System version
DVAC_SHORTNAME	Short name of the installation
DVAC_INSTALLATIONNAME	Name of the installation
DVAC_LOCATION	Location for the installation
CHK_IS_NODE	Whether it is a node (doesn’t directly send alarms, only raise them) or not

Table 2: Detail of fields on table IG_INSTALLATIONGENERAL

Field name	Description
IG_NODO_INSTALLATION	Identifier of the installation which is a node
DVNI_FATHER_INSTALLATION	Identifier of the parent installation

Table 3: Detail of fields on table IG_NODO_INSTALLATION

Field name	Description
DVNI_ERRORNUMBER	Alarm identifier
MESSAGE_ID	Unique alarm identifier
MESSAGE_TYPE	Type of alarm, always set as “notification” (not relevant)
INVOKE_TYPE	Tells whether the alarm has generated itself due to a connection or disconnection (if type is “node”) or is generated by a diagnosis system (“saml”) or energy system (“energy”)
INVOKE_NAME	Irrelevant, always set to “diagnosis”
EVENT_TYPE	Defines the type of alarm which has been generated. Its possible values are listed in table 5.
ADDITIONAL_TEXT	Alarm code
ADDITIONAL_INFOS	Additional parameters to be shown in error message
DVNI_ERRORCATEGORY	Alarm severity. Values from 1 to 5 indicating importance of the alarm, or -1 if the alarm indicates recovery from a previous failure.

Table 4: Detail of fields on table ERS_ERRORS_SAM_ENCE

Field name	Description
fieldElementAlarm	Alarm related to a field element
fieldElementFailure	Failure in a field element
operatorInformation	Information to the operator
imCpuAndCommunications	Related to IM CPU or IM communications
internalDiagnosis	Internal diagnosis of a system
operationsDiagnosisCommunications	Communication error in Operation and Diagnosis systems
ImFecVersions	IM or FEC version
internalTraces	Internal traces of a system
operatorCommandAnswer	Answer to an operator command
CommProblem	Undefined communication problem
Information	Information message: versions, etc.
CommunicationsAlarm	Procedures and processes to carry information from one point to other
QualityOfServiceAlarm	Loss of quality of service
ProcessingErrorAlarm	SW or processing error
EquipmentAlarm	Equipment failure
EnvironmentAlarm	Related to the environment where the system is located
other	Other

Table 5: Description of values for the field EVENT_TYPE

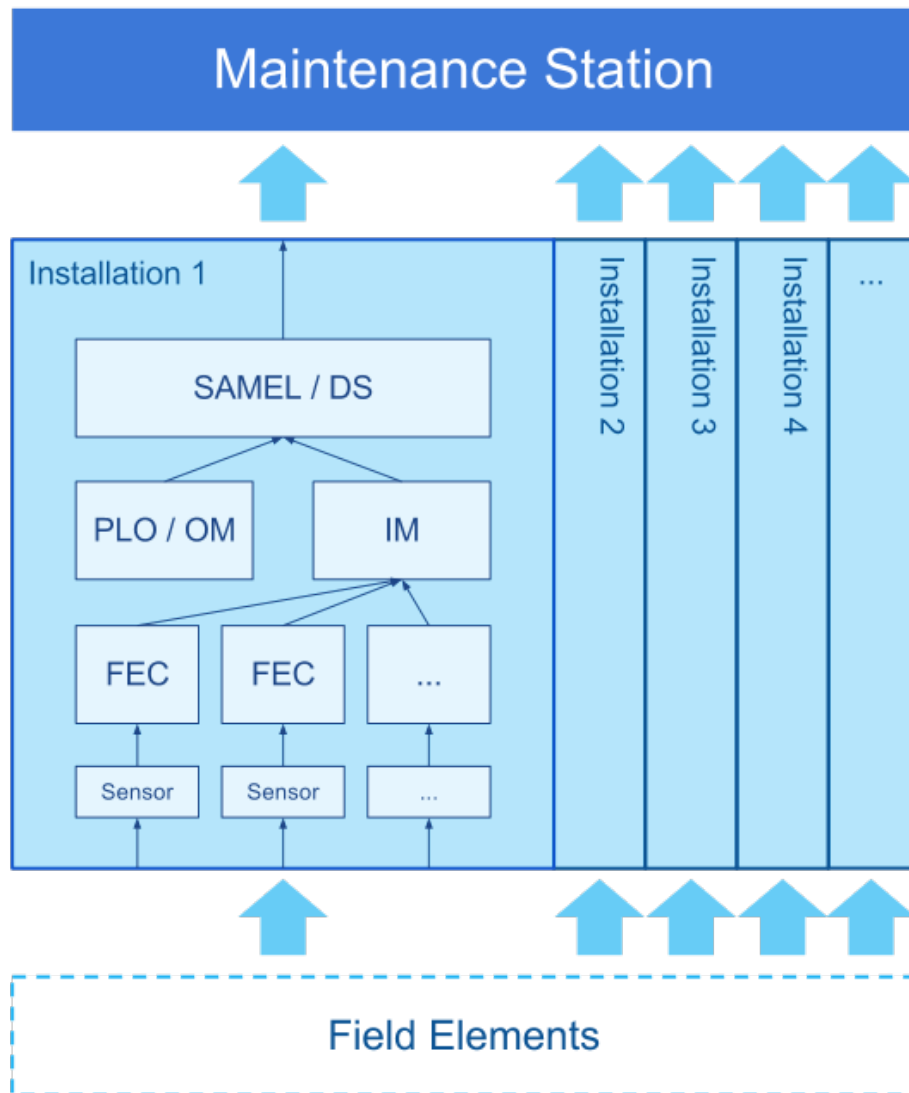


Figure 1: Simplified diagram of the maintenance systems architecture

3 Reduced representation of alarms

In section 2 we have seen a deep definition of all the tables characterising registered alarms. Each of these tables contain several fields, which in total makes an inconvenient large number of variables. While all of them are necessary for correct system function and maintenance purposes, not all of them will be necessary for us to work with alarms.

In order to characterise an event, the main things we need to know can be reduced to three variables:

- What has happened
- When has it happened
- Where has it happened

In section 2 we have seen other variables which can provide additional information which - although not essential - can be useful. Specifically, we think the following data can be of possible interest:

- How severe the event is
- Which type of event has happened

These variables can help us to classify alarms or give more importance to those which are more severe. As this information is already provided on given databases, we will keep it and use it for better alarm classification and filtering. However, none of them are essential in order to characterise alarms, as both of them give information which is already implicit in our previous “*what has happened*” variable. Specifically, this information will be of great help in order to make a preliminary statistical insight on the events of the databases, for which a generalisation in terms of severity and category can help us have a better overview of the situation.

We have to identify which fields on our database corresponds to each of the variables we want to obtain. A direct relation is not possible, as details on *what* has happened is registered in several fields of the database.

This is necessary for maintenance purposes and better alarm handling in the maintenance station, but for our purposes we should identify *what* has happened with a single variable.

In our database, we have unique alarm identifiers for each of the alarms. For better handling and understanding of what is happening, we will use the textual identifier of the events to identify them. This identifier is gathered on the *ADDITIONAL_TEXT* field, and can be translated to a full comprehensive human-readable message by the maintenance station. Furthermore, there is additional data to fill in details about the message. For example, we can have an alarm such as “Communication channel with *X* down”, being *X* an additional parameter saved in the *ADDITIONAL_INFOS* field. Here we can follow two different approaches: disregard the information about *X*, and just treat it as a “Channel down” error; or easily build a compact representation including both variables, such as

4 Statistic analysis

4.1 Alarm classification

In order to have a better insight of the provided databases and the mentioned descriptions, a preliminary insight was made, quantitatively analysing some of the parameters which seemed more relevant for alarm definition. Specifically, the chosen parameters are the following:

- EVENT_TYPE
- INVOKE_TYPE
- DVNI_ERRORCATEGORY (Error Category)

The proportion of each kind of alarms in each of the provided databases (Antequera, Camas, Segovia and Sevilla) is as follows:

In Sevilla, we observe an additional error category marked as “other”. If we make a deeper insight on those errors, we find that is a group formed by 77 alarms of the same type.

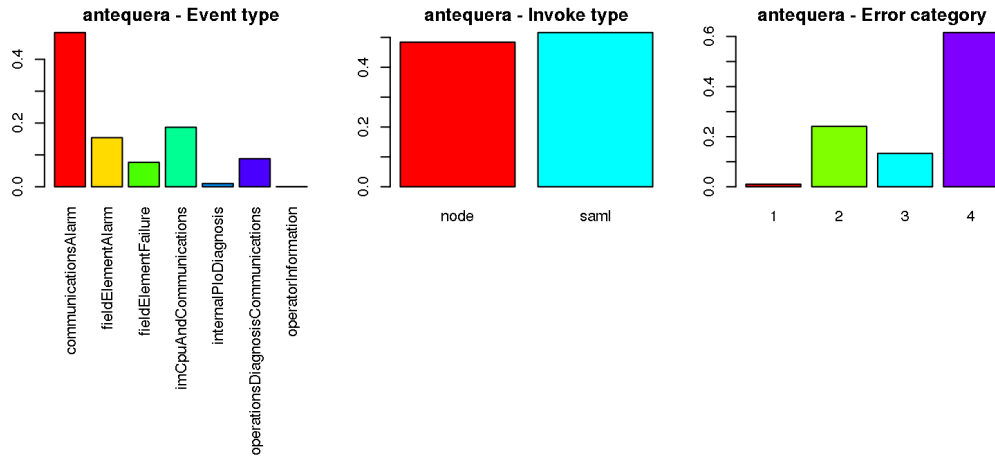


Figure 2: Alarm information for Antequera

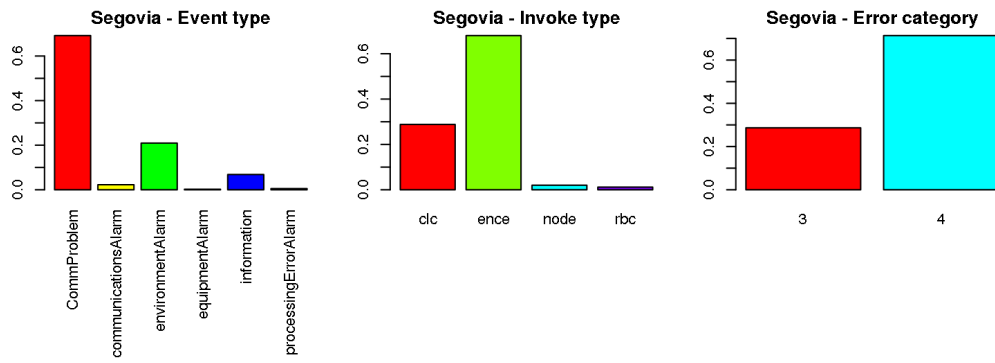


Figure 3: Alarm information for Segovia

4.1.1 Hourly timeline

In order to make a first approach to data analysis, we decided to analyse the alarms on a hourly distribution, checking which types of alarms are more likely to happen in different hours during the day. The result is the following:

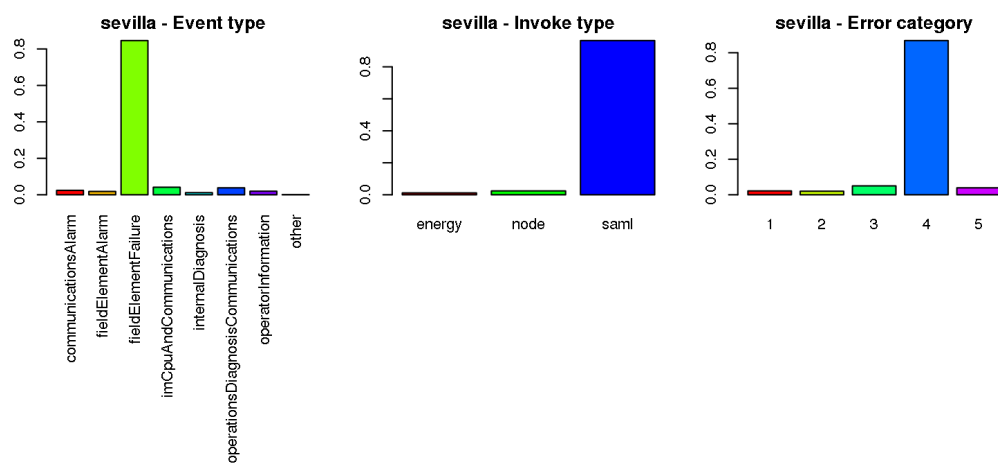


Figure 4: Alarm information for Sevilla

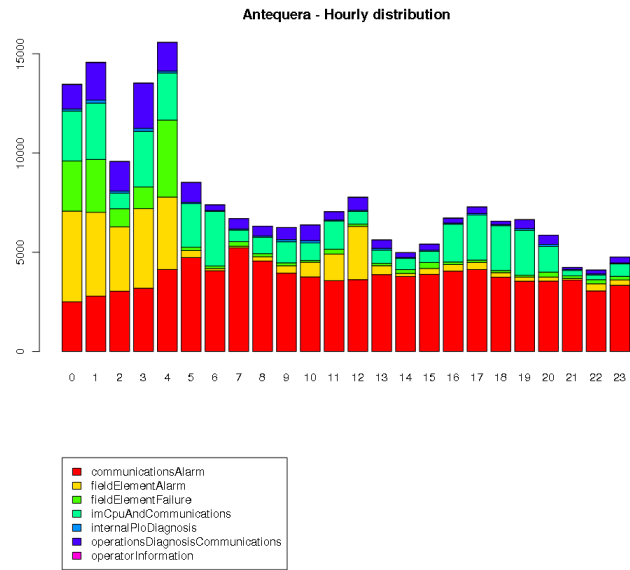


Figure 5: Hourly distribution for Antequera (stacked)

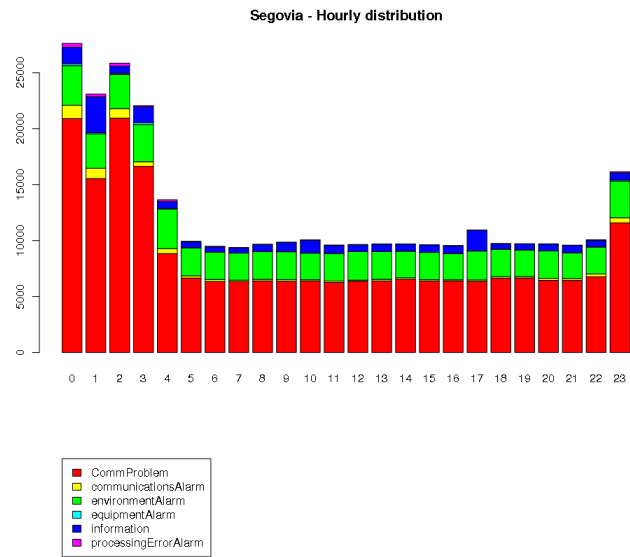


Figure 6: Hourly distribution for Segovia (stacked)

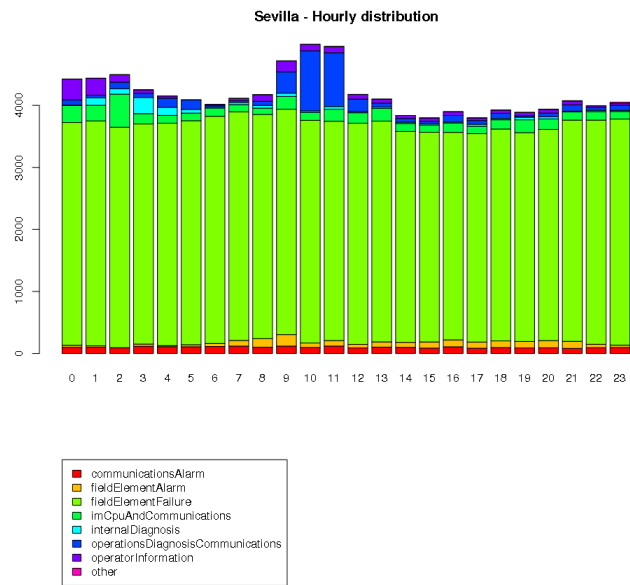


Figure 7: Hourly distribution for Sevilla (stacked)

4.1.2 Daily correlation

We have also generated graphics for correlation between number of alarms of each type during the day, and occurrences of other types of alarms. The result is as follows:

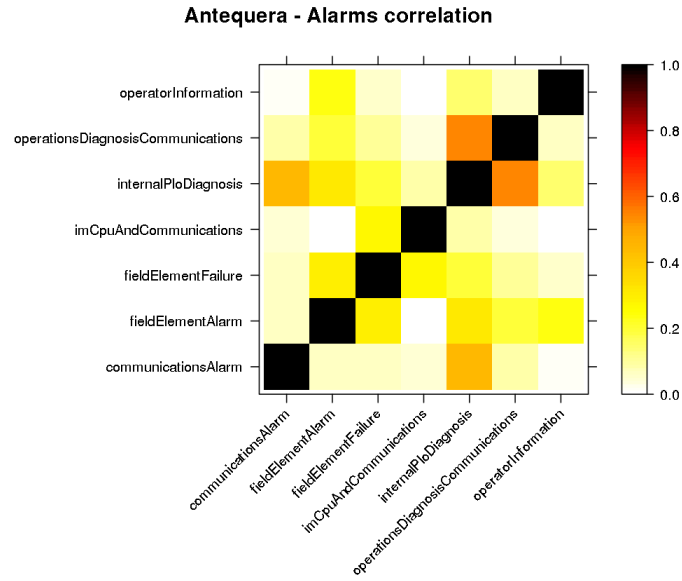


Figure 8: Daily correlation for Antequera

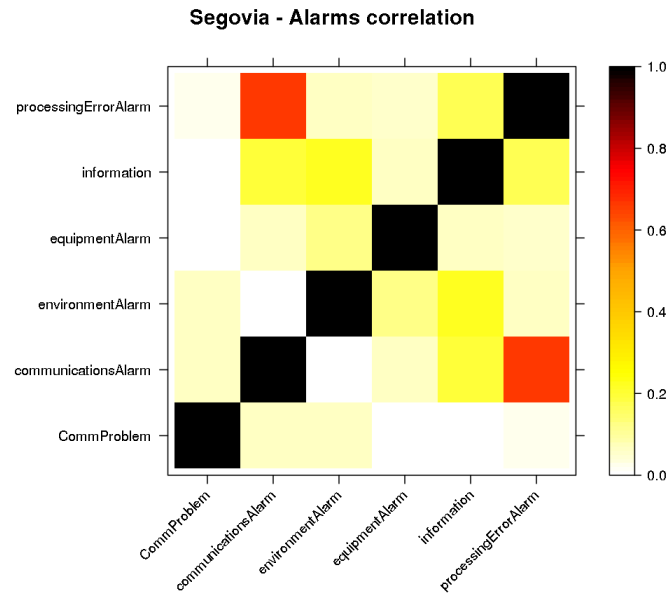


Figure 9: Daily correlation for Segovia

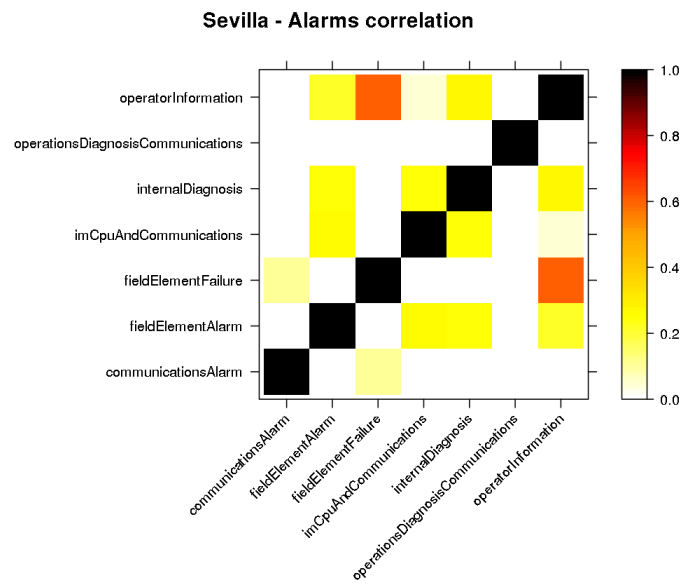


Figure 10: Daily correlation for Sevilla

References

- [1] UM Feyyad. Data mining and knowledge discovery: Making sense out of data. *IEEE expert*, 1996.