# Data Classification using Python and R – Industry Use

**Meenal Pant**

Morpho Detection (MDI), 7151 Gateway Blvd, Newark, CA 94560

## Introduction

1. Why is "data classification" needed ?
Today the world is data driven world … data classification and data mining is critical to understanding patterns and behavior. This enables data prediction and projection.

2. What is data classification ? … and Data mining ?
**Data mining** is typically done on abstract data think social networking data. **Data Classification** is typically done on known data.

At MDI, we manufacture complex, state of the art baggage scanning equipment. The many components of these EDS's (Explosive Detection System) produce a large amount of known data. We *classify, process & predict* this data to provide proactive maintenance
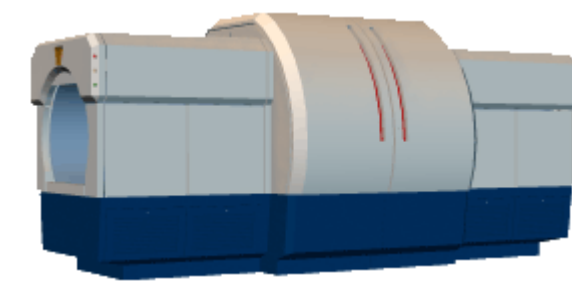
**Figure 1.** Explosive Detection System (EDS)

## Process

- Collect data from the airport(s)
- Analyze large datasets daily and publish on a web dashboard
- Define critical parameters
- Use rules engine and knowledge base to create alerts
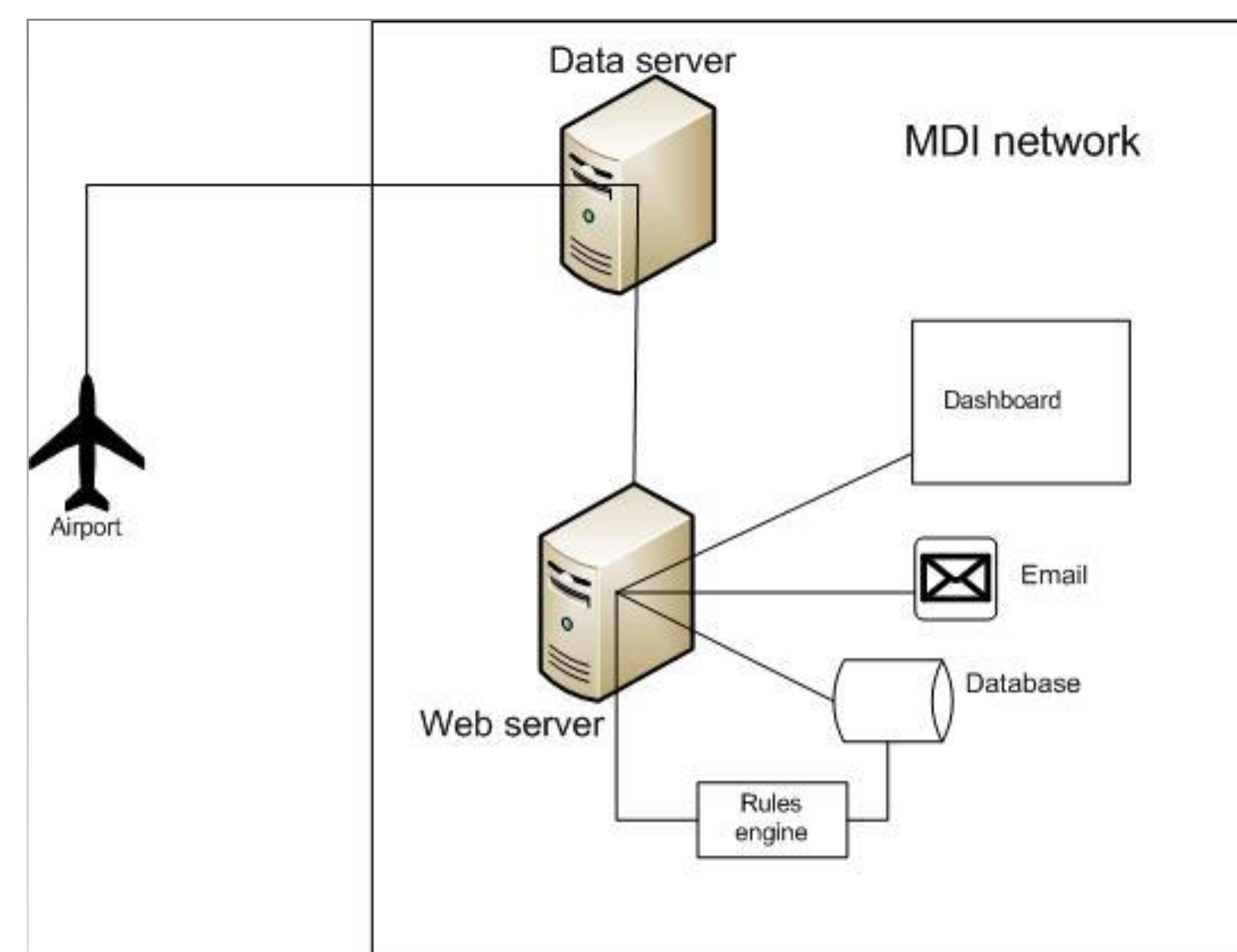- Email alerts to relevant staff

**Figure 2.** Top level design for the real time performance dashboard.

## Technology

**Decision Tree Learning**
From Wikipedia, Decision Tree Learning used in statistics, data mining and machine learning, uses a decision tree as a predictive model which maps observations about an item to conclusions about the item's target value.
In data mining, a decision tree describes data but not decisions; rather the resulting classification tree can be an input for decision making.
http://en.wikipedia.org/wiki/Decision_tree_learning

Classification And Regression Tree (CART) model is used for "rules engine" here.

```
gData <-read.csv("enclerr.csv", sep=",", header=TRUE)
fit <- rpart(Result ~ err + errf, method="class", data=gData)
```
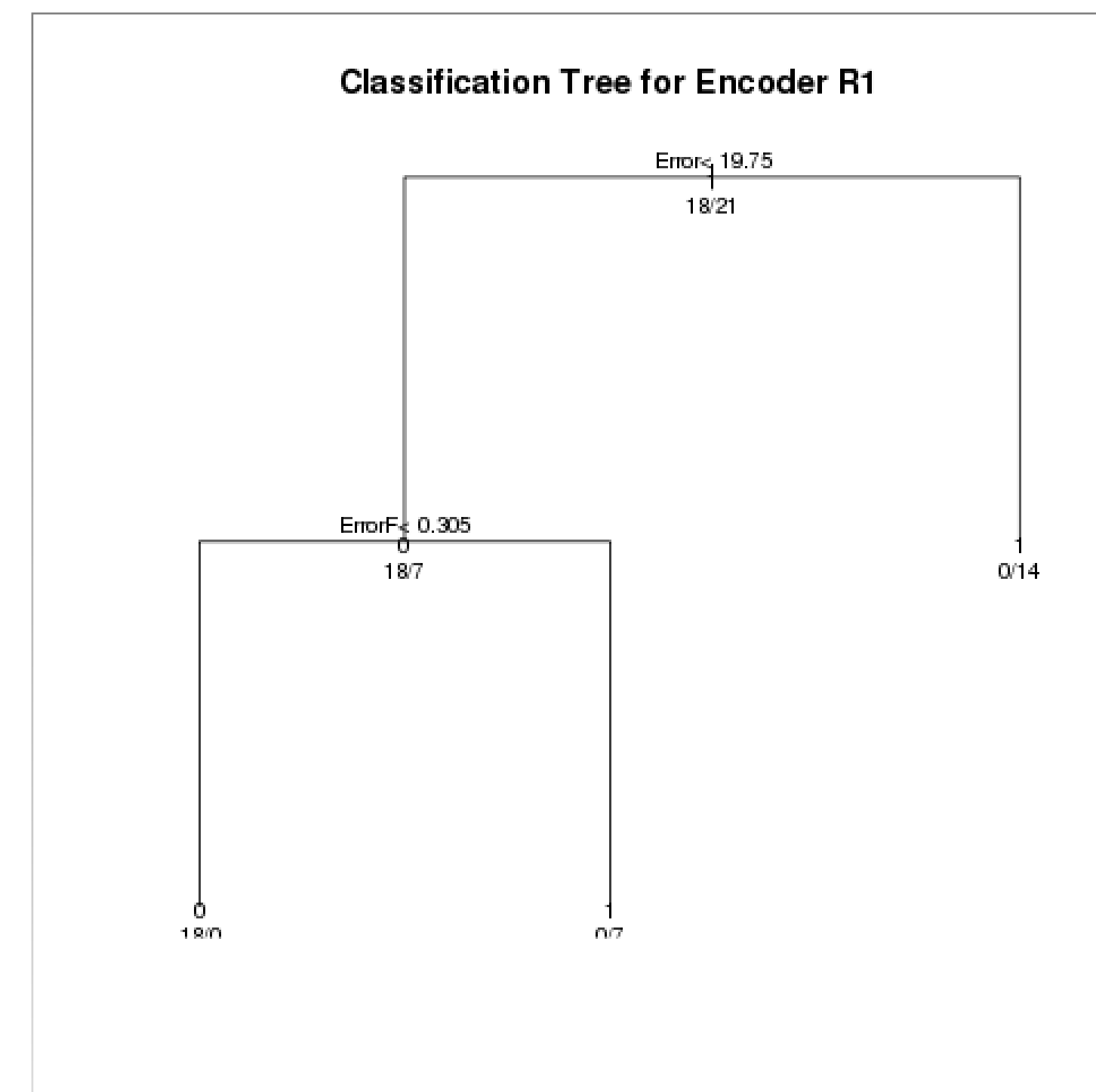
**Figure 3.** Decision Tree Classification for one of the key components in the MDI EDS.

R is a free software package for statistical computing and graphics. rpart is the package for producing CART models
**To create an R decision tree:**
- Get training data samples
- 'Fit' data using rpart
- Score using validation data
- Predict test data
- Additional steps, pruning a tree etc.

## Implementation

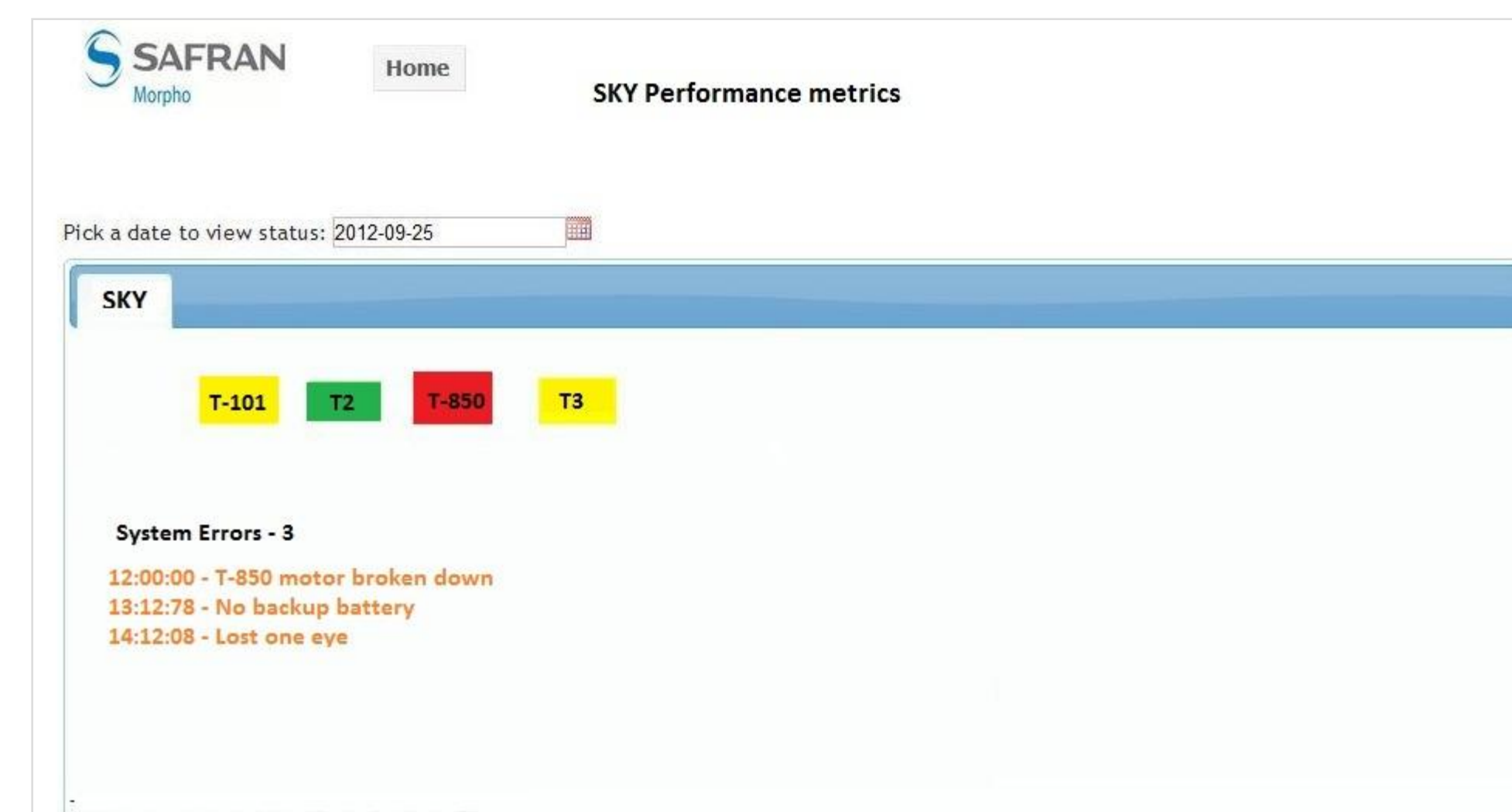Color codes: Yellow = warning, Red = Error , Green = Normal

**Figure 4.** Example of what a status page looks like on the web based dashboard.

## Sample Code

```
#rpy script to predict and score
import os
from config import ROOT_PATH
from rpy2 import robjects
from rpy2.robjects import r
rpart = r.library("rpart")
class Rules(object):
    #Create rules based on CART model
    def __init__(self):
        self.csvpath = os.path.join(ROOT_PATH, "/train/")
    def rltree(self):
        #Fitting the tree based on training sample
        self.elpath = os.path.join(self.csvpath,
"elerr.csv")
        self.rlcmd = 'elData <-read.csv("%s", sep=",",
header=TRUE)'%(self.elpath)
        robjects.r(self.rlcmd)
        robjects.r('elfit <- rpart(res ~ err + ef,
method="class", data=elData)')
    def rlval(self):
        #Validate and predict new data
        robjects.r('validationData <-
read.csv("/tmp/elval.csv", sep=",",   header=TRUE)')
        robjects.r('test <- predict(elfit,interval =
"prediction",newdata=validationData, type="vector")')
        self.scored = robjects.r('test')
        return self.scored
```

✓ Returns a score of 1 or 0 and classifies the data
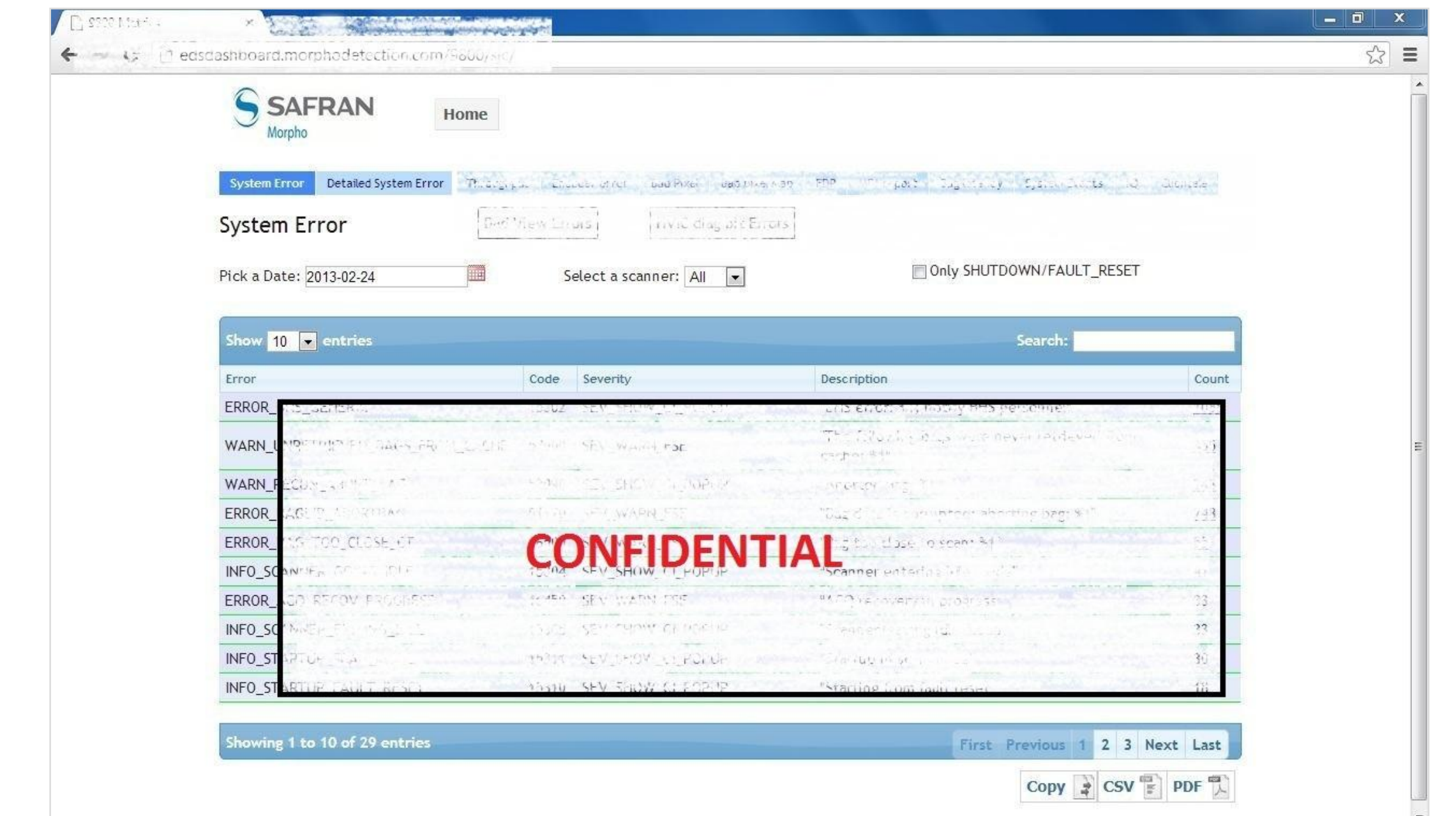✓ Data analysis can now be performed by applying rules on this classified data.
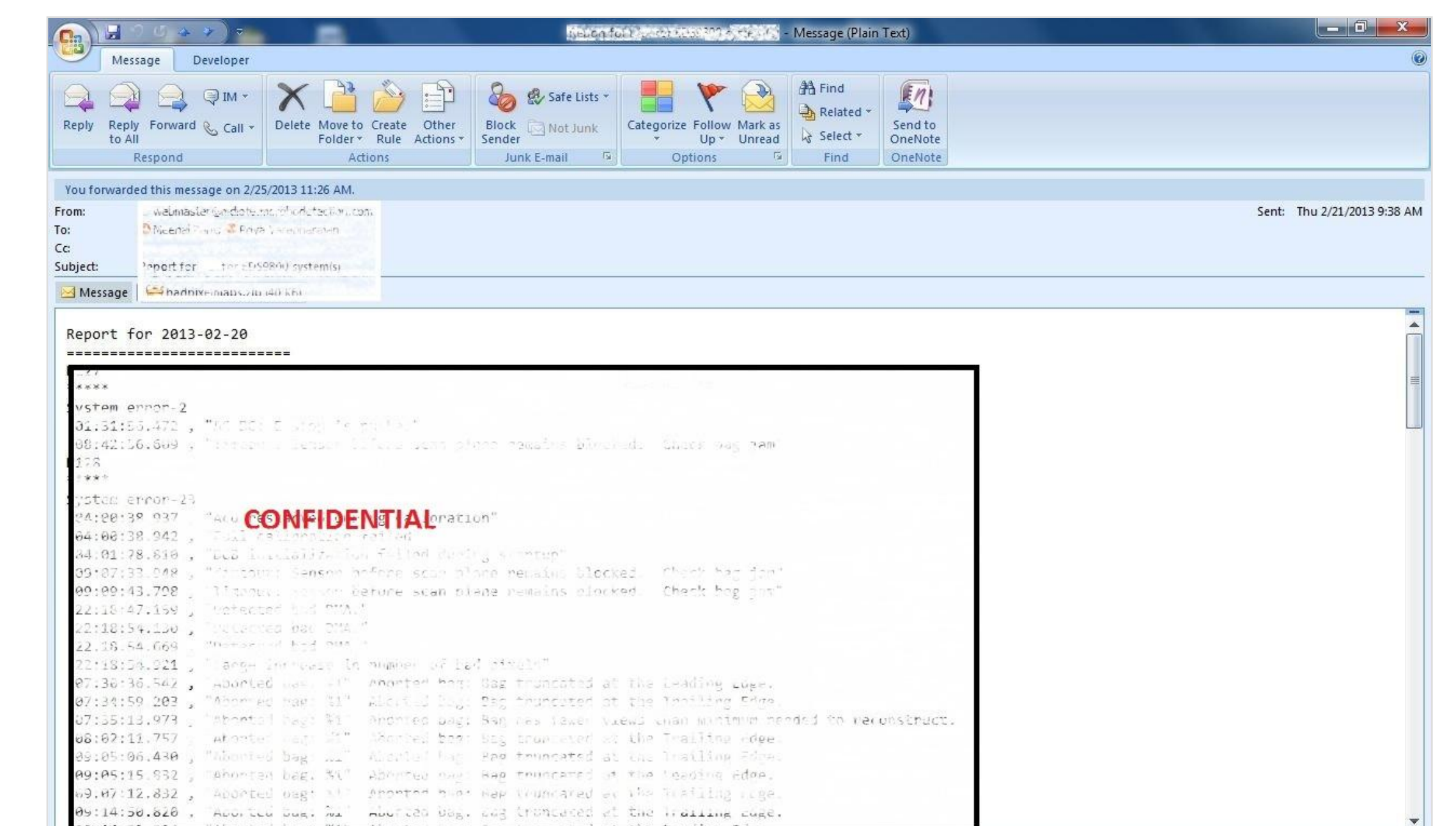
**Figure 5.** Detailed reports example

**Figure 6.** Daily email to the maintenance support staff

## Benefits

- Direct cost benefits due to reduction in downtime
- Can identify and troubleshoot issues faster
- Improves personnel efficiency
- Always available status and results
- Archived results to learn trends and patterns

## Acknowledgments

## Questions and feedback

Please contact *mpant@morphodetection.com*