

Measuring the Accuracy Levels Regarding the Dual Business Function Criticality Classifier

ATHANASIOS PODARAS 

Department of Informatics, Faculty of Economics, Technical University of Liberec, 46117 Liberec, Czech Republic

e-mail: athanasios.podaras@tul.cz

This work was supported by the Czech Ministry of Education, Youth and Sports under the Technical University of Liberec Institutional Support Research Project entitled “Smart Data for the Prevention and Identification of Malfunction in the System” –Project Code: EF-SMARTDATA.

ABSTRACT This paper illustrates the most recent results regarding the criticality ranking decision support classifier for an individual business function. The validated classifier is part of the business continuity points approach which estimates the recovery complexity of an individual business function and is based on the use case points method for estimating the software complexity. The business continuity points method is utilized in order to estimate specific recovery complexity parameters of a given business function. A part of the approach concerns the business function criticality ranking which is based on the recovery complexity parameters. In this paper, we measure the accuracy in the criticality ranking classifier by comparing results between the speedy and the detailed criticality ranking of a business function. The measurement is performed via the R-Studio software and the confusion matrix technique. The results are based on a learning data set prepared in MS Excel which includes the empirical calculations for constructing the specific classifier.

INDEX TERMS Business continuity points, business function, criticality ranking, classifier, R-Studio.

I. INTRODUCTION

One successful definition regarding business continuity management has been provided by domain experts, stating that “business Continuity is the management of a sustainable process that identifies the critical functions of an organization and develops strategies to continue these functions without interruption or to minimize the effects of an outage or a loss of service provided by these functions” [1]. However, apart from its identification, a criticality ranking of any single business function for determining its recovery priority should also be implemented. This task is currently executed based on the experience of managers. Neither a standard mathematical method for classifying an individual business function as critical/non-critical, nor a software tool which supports such a solution have been proposed so far. In order to fill this gap, a standard mathematical method for classifying individual business functions, entitled business continuity testing points (or simply business continuity points) [2] has been recently developed.

The method involves the execution of the two following tasks: firstly, the estimation of the recovery complexity as

well as the recovery time effort for an individual interrupted business function, and secondly, the proposal of the appropriate recovery exercise category. The first aforementioned task includes the calculation of various complexity parameters, similar to the use case points [3] method, such as actors, processes, technical and environmental factors and estimates the unadjusted points, the adjusted points and the recovery time effort for an individual given function. Moreover, for the estimation of the recovery time, different types of recovery scenarios, namely simple, average and complex are considered.

The Business Continuity Points method, hereinafter BCPTs, is based on the concept that the recovery complexity of any function is inversely proportional to its recovery time. Due to the fact that the restoration time is strongly related to the criticality ranking of a business function or an interrupted information system Fasolis *et al.* [4], Gibson [5] constructed a data set to be utilized as a classifier for the criticality ranking of an individual business function.

Every newly proposed classifier, according to the data mining theory, has to be validated with the help of the best

possible data mining techniques, theoretical methods and widely tested software tools. Decision trees and rule based data mining techniques have been successfully applied for decision making against unexpected disruptions that can significantly affect critical business activities as the supply chain management [19] and the industrial safety management [27]. Disaster information management experts claim that “data mining and information retrieval techniques help impacted communities better understand the current disaster situation and how the community is recovering.” [28]. “A decision tree is a classifier expressed as a recursive partition of the instance space. The decision tree consists of nodes that form a rooted tree, meaning it is a directed tree with a node called root that has no incoming edges. All other nodes have exactly one incoming edge. A node with outgoing edges is called an internal or test node. All other nodes are called leaves (also known as terminal nodes or decision nodes)” [20].

A dual business function criticality classification approach had been proposed after implementing the necessary empirical calculations regarding the Unadjusted Points and the Recovery Time Effort. The first classification path is based on the Unadjusted Points and is entitled *speedy classification* while the second is entitled *detailed classification* and stems from the Recovery Time as well as the Adjusted Points values.

Another characteristic of the BCPTs approach is the necessity of testing various different recovery scenarios. The initial calculations led us to the proposal of the *representative recovery scenario* (see Table 3) in order to derive approximate classification results. Moreover, two different subcategories of recovery scenarios are included namely the *default recovery scenario* and the *alternative recovery scenario*. Detailed explanation of these terms is provided in section C.

The initial measurement regarding the first part of the classifier’s validation had been recently conducted. The empirical calculations prompted the Podaras [17] to the creation of specific business rules which permit standard BF classification. The specific rules have been illustrated in the form of inducted decision trees.

However, the rules which had been initially conducted for the BCPTs, had been based on primary calculations which did not permit the thorough rule validation based on all the dimensions of the approach. Specifically, the initial accuracy measurement permitted the rule validation by comparing exclusively the “speedy criticality ranking with the detailed criticality ranking for a business function by applying the *default recovery scenario*”.

The primary goal of the present article is to extend the initial rule validation proposal and measure the accuracy of the early (or speedy criticality classifier) via the comparison of the speedy classification results with the detailed criticality ranking by applying the alternative recovery case. Additionally, the proposal of a new classifier which determines the

selection of the recovery scenario for a given function by considering the Unadjusted Points value is also included in the current paper. For the achievement of the aforementioned goal, the paper has been organized as follows:

1) INITIAL MEASUREMENT OF THE CRITICALITY CLASSIFIER

In this introductory subsection the initial measurement of the criticality ranking classifier is summarized and the accuracy rate between the speedy and the detailed default recovery case are explained.

2) TOOLS AND METHODS SECTION

The specific section reports background information. At first, the utilized data mining techniques as well as the BCPTs approach, the core terminologies and the important equations are delineated for the comprehension of the achieved results. Additionally, the primary advantages of the method are explained. Secondly, a brief explanation of the initially derived measurement based on the default recovery case is included and the steps followed to validate the initial classifier are analyzed. Finally, the specific part inevitably includes a reference to the most widely utilized and tested data mining classification approaches along with the explanation towards the selection of the most appropriate one. Furthermore an explanation for selecting specific software tools which support the classification techniques is included at the same section.

3) RESULTS, DISCUSSION AND CONCLUSIONS SECTIONS

The specific section is devoted to the classification of a business function when applying the detailed alternative recovery case as well as the measurement of the accuracy levels between the detailed default with the detailed alternative recovery approaches for the same individual business function. Moreover the derived measurement is thoroughly analyzed in the discussion section and a new classifier for selecting the appropriate recovery scenario based on the Unadjusted Points value is proposed. The final part of the current work summarizes the current achievements and includes implications for conducting future research for the further investigation of the BCPTs as well as the conducted criticality ranking results.

A. INITIAL MEASUREMENT OF THE CRITICALITY CLASSIFIER

Via the specific measurement the classifier has been proved to be highly accurate (almost 90% accuracy between the two different recovery classification approaches). The specific part of the research has been conducted by importing the learning data set to the R-Studio [6], [15] software package. The learning data set has been divided into 2 subsets, namely the training set and the testing set (70% results rate and 30% results rate respectively) in order to be validated. The overall

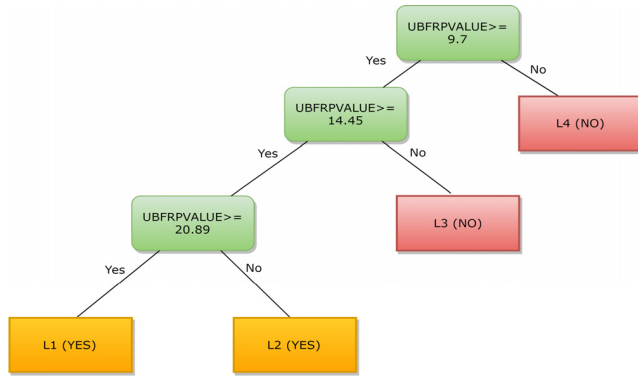


FIGURE 1. The decision tree for the speedy criticality ranking (Source: Author).

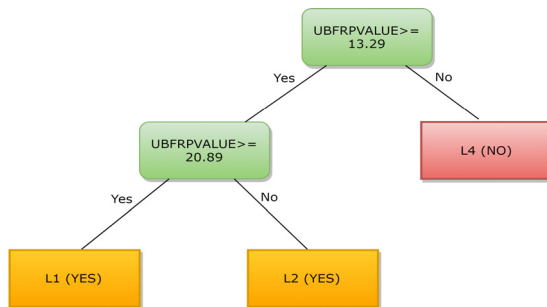


FIGURE 2. The decision tree for the detailed criticality ranking (Default_RS) (Source: [17]).

number of records was 47 different business functions. For each function only the default recovery case had been applied for the speedy (see Fig. 1) and the detailed (see Fig. 2) recovery scenarios.

The induced decision trees have been based on empirical calculations of the initial data set and are considered efficient criticality ranking classifiers for the default recovery case. The recently published confusion matrix validation check indicated 89.36% accuracy rate of the speedy criticality ranking classifier [17]. The confusion matrix is a widely utilized “tool to measure the performance of a classification system” [7]. It has been characterized by data mining researchers as, firstly, the “most common descriptor for assessing the classification accuracy” [23] and, secondly, a technique that can even improve the performance of ensemble (combination of more than one) classifiers [22]. Additionally, the confusion matrix is highly recommended by scientists who utilize the R – Studio package for classification along with the Random Forest approach [24].

II. TOOLS AND METHODS

A. DECISION TREES AND R-PACKAGE

Multiple data mining experts have underlined the importance and the accuracy of decision trees as a supervised learning classification method. Some of them state that “decision trees are an efficient nonparametric method that can be applied

towards classification or to regression tasks” [11], while others mention that “decision trees are one of the most popular classification algorithms used in data mining and machine learning to create knowledge structures that guide the decision making process. The creation of a good knowledge structure is the main step in the development of a decision making system” [21].

The core idea behind the induction of a decision tree is the discovery of the ideal split. “Each split partitions the sample into two or more parts and each subset of the partition has one or more classes in it. If there is only one class in a subset, then it is pure, else it is impure. The purer the partition is, the better it is” [35]. The most important splitting measures regarding the decision tree induction are the following [35]:

Information gain (1):

$$Gain(S, A) = Entropy(S) - \sum (|S_v| - |S|) Entropy(S_v) \quad (1)$$

Where $v \in Values(A)$, $v = Instances\ of\ Features\ A$.

Gain Ratio (2),(3):

$$Gain_Ratio(S, A) = Gain(S, A) / Split_Information(S, A) \quad (2)$$

$$SplitInformation(S, A) = \sum_{i=1}^c (|S_i|/|S|) \log_2 (|S_i|/|S|) \quad (3)$$

Nevertheless, despite their broad acceptance from the data mining scientists, some researchers have highlighted a number of drawbacks regarding the decision trees. Some of these researchers refer to the problem of over-sensitivity of the training set [20] or their ineffectiveness when changes occur in the data set [21], while other refer to the problem of the overfitting, stating that “the decision tree algorithm can produce more branches and leaves of tree and most of branches are over-fitting for training data samples so that there are some bias to predict new data in the application” [26]. Other researchers [29] refer to the dimension problem of the traditional decision trees which is caused by redundant features. According to these researchers the most commonly utilized technique for overcoming the overfitting problems is known as pruning, while the solution to the problems of high variance and bias of the training data is defined as ensemble classifiers for which the term *boosting* is used [33], [34]. Despite their major importance in deriving accurate results the pruning and boosting implementation is not a subject of the current contribution since deep explanation and thorough analysis is required for the selection of the ideal approach for the presented classifier. The application of these techniques to the investigated classifier will be analyzed in a future publication. A reference to these techniques and their possible application to the current decision trees is included in the discussion section..

For the induction of the decision trees the R Package has been selected. The R package [6] is considered to be an ideal software tool for the decision tree induction procedure,

due to the fact that it is a free [24] and it supports various algorithms [16], i.e. ID3, C4.5 and CART, compared to other packages.

B. CLASSIFICATION AND REGRESSION TREES

Classification and Regression Tree (CART) [12–14] “has been found ideal for the current model due to the fact that it supports classifications for both binary and continuous variables” [17]. Moreover, experts from the information security domain, which is strictly bound to the business continuity management, have proposed CART as the ideal decision algorithm to be combined with Fuzzy Logic towards the accuracy of Intrusion Detection Systems [25].

Multiple data mining researchers have thoroughly cited and analyzed other important decision tree algorithms like ID3 [20] as proposed by Quinlan [30], and its improved version namely C4.5 [26]. ID3 algorithm has not been selected for the induction of the BCPTs decision trees due to the fact that it cannot handle neither numeric attributes (i.e. Unadjusted Points in the BCPTs model) nor missing values. Moreover, the C4.5 algorithm which is a more complete version of ID3, has the disadvantage of producing large decision trees which may include errors of misclassification [31]. Moreover, in many cases where the C4.5 has been compared to the CART algorithm, the latter demonstrated better performance regarding the classification accuracy [33].

The CART algorithm supports the construction of binary trees which are ideal for the BCPTs classifiers, since the main target classification of an individual business function is binary (critical/non-critical) no matter the 4 possible *Impact Value Levels (IVLs)*. Moreover, the CART algorithm implements exhaustive search regarding the discovery of the best splitting attributes by utilizing the impurity measurement for each attribute. The impurity is calculated with the help of the *Gini (diversity) index* (4) [14, 26].

$$Gini = 1 - \sum_{j=1}^n p(j)^2 \tag{4}$$

where, p(j) is the relative frequency of class j in T, and T is the dataset which contains examples of n classes. “Gini index is an impurity-based criteria that measures the divergence between the probability distributions of the target attribute’s values.”[20]

C. THE BUSINESS CONTINUITY POINTS METHOD

The method is based on the principles the Karner’s [3] use case points theory which has been proposed as an approach that determines the complexity of the software development process. A brief demonstration of the BCPTs method compared with use case points is required for the interpretation of the way that can be applied to estimate the recovery complexity of a given business function (see Table 1).

TABLE 1. comparison between the use case points and the business continuity points [8].

| | Use Case Points | Business Continuity Points |
|---------------------------------|---|--|
| Estimated Complexity Type | Software Complexity Estimation | Business Function Recovery Complexity Estimation |
| Use Cases vs Business Functions | Use Cases are classified as Simple, Average and Complex (according to the number of involved transactions), utilized to calculate Unadjusted Use Case Weights | Business Functions are classified as Simple, Average and Complex (according to the number of involved processes), utilized to estimate Unadjusted Business Function Weights (UBFW) |
| Actors Classification | Actors’ classification (Unadjusted Actor Weights – UAW) | Separate Classification of Human Level Actors and Application Level Actors involved in the Process (Total Unadjusted Actor Weights – TUAW) |
| Unadjusted Points Estimation | Unadjusted Use Case Points : UCP = UAW + UUCW | Unadjusted Business Function Recovery Points: UBFRP = TUAW + UBFW |
| Technical Factors | 13 Technical Factors (Limited Number) | Unlimited Number of Technical Recovery Factors (TRF) |
| Environmental Factors | 8 Environmental Factors (Limited Number) | Unlimited Number of Environmental Recovery Factors (ERF) |
| Unexpected Factors | No Unexpected Factors are Considered | Unlimited Number of Unexpected Recovery Factors (URF) |
| Method of Weight Assignment | Based on the experience of IT Project Manager | Based On Standard Mathematical Approach (Rank Order Centroid) [9] |
| Adjusted Points Estimation | Adjusted Use Case Points (UPC) | Adjusted Business Function Recovery Points (ABFRP) ABFRP= UBFRP*TRF*ERF*URF |
| Effort Estimation | Effort = UCP * Hours/UCP | Recovery Time Effort (RTE) = (5000/ABFRP ²) - 3 |

The primary idea has been focusing on implementing criticality ranking of a business function by calculating its recovery time (RTE). However, during the empirical calculations regarding the business continuity points method the following conclusions have been inferred:

a) When the Unadjusted Points value is either very low or very high the corresponding criticality ranking of a BF can be determined/predicted without calculating its recovery time. This approach has been named as the *speedy classification of a given business function*, and is based on the *representative recovery scenario* (see Table 3).

b) The speedy classification cannot be applied with confidence for middle values of the Unadjusted Points. In such cases the recovery time is calculated and used as a criticality ranking reference value regarding the specific business function. Then the detailed criticality ranking is implemented.

Thus, a dual approach for BF criticality ranking is proposed. The equations which provide all the recovery complexity are all included in the summarized comparison between the Use Case Points and the BCPTs methods (see Table 1).

c) When different recovery scenarios are applied, different criticality ranking results can be inferred. Thus, one crucial issue in order to ensure the validity of the constructed classifier is to compare the criticality ranking inferred results between the speedy and detailed classification by applying the default [17] and the alternative recovery scenario. The latter validation is implemented in the current paper and the output is delineated in the results section. In any case, the classification of any individual business function follows the standard business continuity criticality ranking as recommended by domain experts.

The major advantages of the BCPTs approach are the following:

- The criticality ranking classifier has been based on simple calculations that can be easily analyzed to and rapidly interpreted by business managers.
- The method refers to a broad industrial sphere due to the fact that recovery time values can be foreseen for highly critical as well as less important activities.
- The estimated recovery timeframes can be easily compared by standard mathematical tools utilized towards the uninterrupted operation of critical business activities. A typical tool is the *system availability* provided by the following formula (5) [32]:

$$A = \frac{MTBF}{MTBF + MTTR} \tag{5}$$

Where A = Availability of the service, software application, network, business function, MTBF = The mean Time Between Failure and MTTR = Mean Time To Repair. We can thus, compare the proposed availability rates (%) by the domain experts, with the availability rates calculated by the BCPTs when replacing the MTTR with the calculated RTE value.

D. CONNECTING CRITICALITY RANKING OF BUSINESS FUNCTIONS WITH RECOVERY TIMEFRAMES

A thorough and practically used classification approach has been proposed by [5]. The specific approach determines 4 Impact Value Levels (IVLs) where each level includes the corresponding recovery timeframes, namely the Rational Time Objective (RTO) and Maximum Accepted Outage (MAO) [10]:

- IVL L1: Maximum Acceptable Outage (MAO) = 2 hours. Recovery Time Objective (RTO) < 2 hours.
- IVL L2: Maximum Acceptable Outage (MAO) = 24 hours (1 day). Recovery Time Objective (RTO) < 24 hours.

- IVL L3: Maximum Acceptable Outage (MAO) = 72 hours (3 days). Recovery Time Objective (RTO) < 72 hours.
- IVL L4: Maximum Acceptable Outage (MAO) = 168 hours (1 week). Recovery Time Objective (RTO) < 168 hours.

The above levels are used for classifying a business function as Critical (L1, L2) and Non-Critical (L3, L4). Similar criticality ranking methods have been also proposed[4].

E. CORE CALCULATIONS – DATASET PREPARATION

The recovery time effort (RTE) has been derived with the help of a created by the author initial data set. The initial data set assumed recovery parameters for 47 different business functions. The dataset was prepared in Microsoft Excel 2013 and the implemented calculations have been the following [2]:

UBFRP Value, derived from calculations with the equations described in the first table (see Table 1).

RTE Value is derived from calculations with the help of the equations included in the above mentioned table (see Table 1).

The most representative recovery scenario, regarding the involved human and application level actors as well as the number of the involved business activities is below depicted (see Table 3). The scenarios indicate their level of severity (difficulties during the recovery process, i.e. non-skilled workers, many distributed systems, network unavailability e.t.c.).

However, during the above delineated procedure, it had been realized that for specific business functions (i.e. UBFRP = 8.2 points), even if the predicted criticality ranking is IVL=L4 and BF=Non-Critical when applying the default recovery scenario (simple), a different IVL is inferred when we apply an alternative recovery scenario (i.e. average, IVL=L3). This occurs due to the different ABFRP and RTE for the same business function. The default and the alternative recovery scenario are comprised of the following cases (see Table 2):

TABLE 2. Summarized delineation of the default_RS and the alternative_RS.

| Recovery Case | UBFRP (Points) | TRF,ERF, URF | Selected Approach (Default Scenario) | Selected Approach (Alternative Scenario) |
|---------------|----------------|--------------|--------------------------------------|--|
| Case 1 | <= 9 points | 0.85 | SIMPLE | AVERAGE |
| Case 2 | >9 and <=21 | 1 | AVERAGE | COMPLEX/SIMPLE |
| Case 3 | >21 | 1.15 | COMPLEX | AVERAGE |

F. VALIDATION OF THE LEARNING DATASET WITH DATA MINING TOOLS

The constructed learning data set has been prepared by implementing the calculations of the business continuity points

TABLE 3. The representative recovery scenario [2].

| Human Actors | Application Level Actors | Business Activities | UBFRP | ABFRP | RS | RTE(Hours) |
|--------------|--------------------------|---------------------|-------|-------|----|------------|
| 3 | 3 | 3 | 9 | 5.5 | S | 160 |
| 5 | 5 | 6 | 15 | 15 | A | 21 |
| 7 | 7 | 9 | 21 | 31.9 | C | 1.9 |

TABLE 4. Part of the constructed learning dataset [2].

| BF_ID | UBFRPVALUE (points) | IVL_UBFRP | RS | RTE(hours) | IVL_RTE |
|-------|---------------------|-----------|----|------------|---------|
| 1 | 8.1 | L4(NO) | S | 199.06 | L4(NO) |
| 3 | 8.12 | L4(NO) | S | 198.06 | L4(NO) |
| 5 | 8.15 | L4(NO) | S | 196.59 | L4(NO) |
| 33 | 13 | L3(NO) | S | 75.44 | L4(NO) |
| 35 | 13.29 | L3(NO) | S | 72.05 | L4(NO) |
| 39 | 13.35 | L3(NO) | S | 71.38 | L3(NO) |
| 41 | 13.36 | L3(NO) | S | 71.27 | L3(NO) |
| 61 | 17 | L2(YES) | A | 14.3 | L2(YES) |
| 63 | 18 | L2(YES) | A | 12.43 | L2(YES) |
| 65 | 19 | L2(YES) | A | 10.85 | L2(YES) |
| 84 | 25 | L1(YES) | C | 0.45 | L1(YES) |
| 86 | 26 | L1(YES) | C | 0.19 | L1(YES) |

method in order to estimate the recovery complexity parameters as well as the corresponding recovery time effort (RTE) for every individual business function. The data set has been imported to the R-Studio software package where it has been divided to the training (70% of the records) and testing (30% of the records) subset. Part of the learning set is below illustrated (see Table 4).

From the illustrated data set it is realized that for the speedy classification, the IVL Classification (predicted value) is based on the UBFRPVALUE (predictor) which stems from the representative recovery scenario.

For the detailed classification the predicted value is IVL based on the RTE predictor. Table 4 illustrates the default recovery scenario for 47 different business functions.

III. RESULTS

A. DECISION TREE – DETAILED CRITICALITY RANKING (ALTERNATIVE_RS)

The present work includes a further evaluation of the business continuity points classifier by measuring the accuracy of the speedy criticality ranking when the Alternative Recovery Scenario is applied.

It is important to indicate that the speedy classification will lead again to the same classification results (see Fig.1) due to the fact that the specific approach to criticality ranking follows the representative recovery scenario. Thus, the next step for evaluating the proposed classifier is to compare the detailed criticality ranking based on the alternative recovery scenario with the results inferred by the initial classifier (see Fig. 1).

The inducted decision tree according to the detailed criticality ranking and the alternative scenario is depicted in Fig. 3.

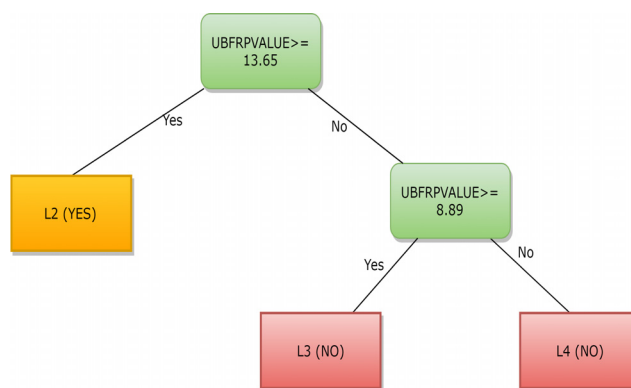


FIGURE 3. Decision tree for the detailed criticality ranking (Alternative_RS) (Source: Author).

The code in R-Studio for decision tree induction based on the CART algorithm is the following:

```

> install.packages("rpart") {\#}package
  for decision tree using CART
> install.packages("rpart.plot")
{\#}package for better visualization
> library(rpart)
> library(rpart.plot)
> datafile <- read.csv("C:/Users/
datafile.csv", sep=";")
> View(Case1)
> fit1a<-rpart(IVL\UBFRP~UBFRPVALUE,
Case1)
> plot(fit1a, uniform=TRUE,margin=0.2)
> text(fit1a, use.n=TRUE, all=TRUE,
cex=.8)
    
```

By applying the confusion matrix technique to measure the accuracy rate between the speedy and the detailed criticality

ranking when the alternative RS is performed, the following results have been obtained:

Confusion Matrix and Statistics

| | Reference | | | |
|------------|-----------|----------|---------|---------|
| Prediction | L1 (YES) | L2 (YES) | L3 (NO) | L4 (NO) |
| L1 (YES) | 2 | 1 | 0 | 0 |
| L2 (YES) | 8 | 11 | 3 | 0 |
| L3 (NO) | 0 | 0 | 9 | 9 |
| L4 (NO) | 0 | 0 | 0 | 4 |

Overall Statistics

Accuracy: 0.5532

B. COMPARING THE DETAILED DEFAULT_RS WITH THE DETAILED ALTERNATIVE_RS

Another important task regarding the measurement of the accuracy levels of the developed business continuity classifier, is the comparison of the detailed classification results between the Defaults_RS and the Alternative_RS. This task is also an extension to the initially published research results [17].

The specific process demonstrated a lower level of accuracy. Again the confusion matrix technique has been applied:

Confusion Matrix and Statistics

| | L1 (YES) | L2 (YES) | L3 (NO) | L4 (NO) |
|----------|----------|----------|---------|---------|
| L1 (YES) | 2 | 1 | 0 | 0 |
| L2 (YES) | 8 | 11 | 3 | 0 |
| L3 (NO) | 0 | 0 | 4 | 14 |
| L4 (NO) | 0 | 0 | 0 | 4 |

Overall Statistics

Accuracy: 0.4468

Nevertheless, the derived accuracy results regarding the performance of the speedy classifier are not that discouraging since they concern a 4-level classification of any individual business function. From the discussion section that follows it will be comprehended that a binary classification (Critical_BF/Non-Critical_BF) with the help of the initial classifier can be highly accurate. The R Code for applying the confusion matrix technique in order to control the accuracy of the predicted versus the actual values is the following:

```
> install.packages("caret")
> library(caret)
> require(caret)
> install.packages("e1071")
> library(e1071)
> require(e1071)
> K<-confusionMatrix(Case1$IVL_RTE,
  Case1$IVL_UBFRP)
```

IV. DISCUSSION

The primary critical issue to be analyzed is the consistency of the constructed data set. The data set has been based on the comparison of the inferred recovery timeframes with the proposed by the available literature recovery timeframes for

critical business functions which ensures reasonably derived recovery timeframes and criticality ranking for individual business functions.

The second point which requires further discussion is the initial validation of the constructed data set. The learning set was split into a training (70%) and a testing (30%) subsets. The derived results indicated minor differences. More precisely [17] “the critical UBFRP Value in the training data set was 12.65 while in the testing subset was 13.32. The criticality ranking in the first case was determined as L2 (YES) (critical BF) in the first case when $UBFRP \geq 12.65$, and L1 (YES) in the second case when $UBFRP > 13.32$ ”. Despite the different Impact Value Levels the differences are not of major importance since both subsets classify a BF as critical when UBFRP is approximately lower than 13. The study includes also the proposed decision trees.

Another point to be explained is the low level of accuracy of the speedy criticality ranking (55.32%) when compared with the detailed classification based on the alternative recovery scenario. According to the confusion matrix results, even if in some cases the IVL is different still the given business function is binary classified in the same way (critical, or non-critical). Thus, no significant influence regarding its binary classification as critical/non-critical had been observed. Only 3 out of the 47 tested business functions are classified differently, which means that the prediction accuracy level regarding the binary classification is almost 96%. The speedy classification predicted $IVL = L2$ (YES – Critical BF) while the real (Reference) $IVL = L3$ (NO – non critical) which is based on the detailed classification. The speedy criticality ranking, when compared with the detailed ranking default recovery scenario, has been proved more accurate in predicting precise IVL levels (89.36%) than in the case when the alternative scenario is utilized. Yet, it has been proved less accurate with respect to the binary classification (5 out of 47 different results, accuracy level 90%). Furthermore, no matter the 44.68% similarity rate between the detailed default and the detailed alternative recovery cases, only 3 out of 47 BFs were differently categorized when the binary classification is implemented (critical/non-critical).

Another issue to be discussed is the existence of a standard classifier for determining among the simple, average and complex recovery scenario. In some cases, when the RTE value is not calculated, it is quite demanding to determine the ideal recovery case (Simple, Average or Complex) for classifying a Business Function. The precise IVL classification is hard to be conducted since it varies between different recovery scenarios. The decision tree which has been inducted according the speedy (default) criticality ranking and the initial learning data set is the following (Fig. 4):

A final issue to be discussed is the possibility of misclassification of the presented classifiers. Based on the proposed recovery timeframes [5], [10], broad intervals permit the

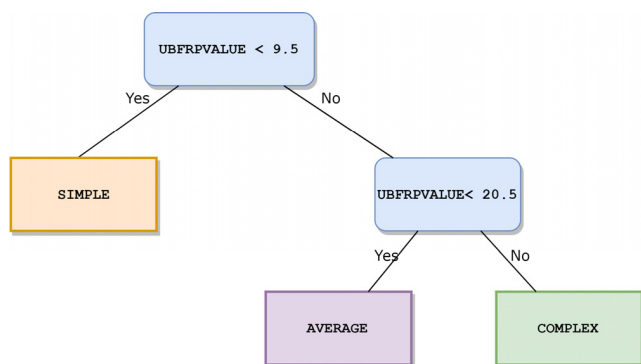


FIGURE 4. The induced decision tree for selecting simple, average or complex recovery case based on the speedy classifier (Source:Author).

classification of business functions in the same category, i.e. $IVL=L2$ for $2 \text{ hours} < RTO < 72 \text{ hours}$. This means that misclassification risk is extremely low.

A problem emerges when selecting the alternative recovery case instead of the standard recovery scenario. The accuracy level in this case of the initial speedy default classifier is only 55.32% or 44.68% when we follow the detailed default or the detailed alternative recovery case respectively. However, the selection of the precise recovery scenario does not affect the BFs which are highly critical, i.e. $UBFRPVALUE > 20.83$ points. Additionally, based on the derived decision trees (Fig. 1, Fig. 2 and Fig. 3) regarding business functions for which the $UBFRPVALUE \geq 14.45$ points, the predicted classification (IVL) is the same no matter the recovery scenario selected. The question remains for the BFs for which $13.29 \leq UBFRPVALUE < 14.45$. In this case the BF must be classified as critical ($IVL=L2$) in order to mitigate recovery risks in practice. Yet, the interval remains very short.

Based on these assumptions ensemble classification techniques like boosting which can improve the performance and the accuracy of the induced decision trees as well as resolve the weakness of high variance of decision trees are currently under investigation. Examined techniques are the ensemble classifiers like adaptive and gradient boosting which fit best and have been proved to be valuable for the CART algorithm [34]. Finally, for resolving overfitting in the future regression trees the pruning procedure will be considered for the accurate prediction of RTE values. In this part of the research the post-pruning algorithm entitled Cost Complexity Pruning can be applied, due to the fact that “it has been proposed by Breiman in the development of the CART system” [26]. For inducing decision trees that predict the RTE values again the CART algorithm will be utilized since ID3 and C4.5 require target attributes with discrete values [20].

V. CONCLUSIONS

The present article analyses the measurement of the accuracy of the speedy criticality ranking by comparing the predicted

impact value levels with those inferred by the detailed criticality ranking - alternative recovery scenario. The classifier, which is based on the supervised learning [18] approach to the induction of decision trees, is 89.36% accurate when applying the default recovery scenario according to [17] and 56% accurate when applying the alternative recovery scenario for an individual business function. The latter case is an extension to the initial work, additionally with a) the proposal of the detailed recovery classification algorithm for a business function when the alternative recovery case is applied (Fig. 3), and b) the similarity control between the detailed default classifier and the detailed alternative recovery classifier, which is estimated 44.68%. However only 3 out of 47 BFs, according to the inferred results, are different when the binary classification is applied, which is also a highly positive remark regarding the business continuity points classifier.

Additionally, the speedy classifier is 90% accurate, no matter the recovery scenario when we want to implement binary criticality ranking for the same business function (Critical/Non-Critical). Hence, the pruning technique is not considered a vital task in the currently proposed classifiers due to the short size of the derived decision trees as well as their high binary classification accuracy. The pruning technique is proposed for the regression trees which will predict RTE values based on the Unadjusted Points (UBFRP).

Ensemble classifiers, adaptive and gradient boosting are currently under investigation. The currently derived accuracy levels are satisfactory but only if we focus a) on the binary target values of the critical functions and b) not on middle values of Unadjusted Points. For these intervals regression trees have to be induced and ensemble classifiers will further improve the prediction of precise IVL Levels.

The core advantage of the criticality ranking classifier can be used in the early stages of the business impact analysis formulation in the industry. The specific algorithm is a crucial part of the Business Continuity Points method as it has been introduced by the author. No matter the successful performance of the speedy criticality ranking classifier, the detailed classification is strongly recommended due to the fact that precise recovery timeframes and impact value levels can be determined ignoring possible weaknesses of unstable classification based solely on the UBFRPVALUE. Finally, future research tasks are mainly oriented to the proposal of an integrated database tool which can perform accurate predictions. A conceptual model of the specific database solution is currently designed.

REFERENCES

- [1] E. Tucker, “Business continuity: A definition and a brief history,” in *Business Continuity from Preparedness to Recovery: A Standards-Based Approach*, 1st ed. Atlanta, GA, USA: Elsevier, 2015, p. 1.
- [2] A. Podaras, K. Antlova, and J. Motejlek, “Information management tools for implementing an effective enterprise business continuity strategy,” *E M, Ekon. Manage.*, vol. 19, no. 1, pp. 165–182, Jan. 2016, doi: [10.15240/tul/001/2016-1-012](https://doi.org/10.15240/tul/001/2016-1-012).

- [3] G. Karner, "Use case points—Resource estimation for objectory projects," 1993. Objective Systems SF AB. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.604.7842&rep=rep1&type=pdf>
- [4] E. Fasolis, V. Vassalos, and A. Kokkinaki, "Designing and developing a business continuity plan based on collective intelligence," presented at the *I3E*, Apr. 2013. [Online]. Available: https://link.springer.com/content/pdf/10.1007%2F978-3-642-37437-1_23.pdf
- [5] D. Gibson, "Mitigating risk with a business impact analysis," in *Managing Risks in Information Systems*, 2nd ed. Burlington, NJ, USA: Jones & Bartlett Learning, 2010, p. 320.
- [6] T. Rahlf, "Data for everybody," in *Data Visualisation With R: 100 Examples*. Cham, Switzerland: Springer, 2017, pp. 1–4.
- [7] Ö. F. Söylemez and B. Ergen, "Eye location and eye state detection in facial images using circular Hough transform," presented at the *CISIM*, Sep. 2013. [Online]. Available: https://link.springer.com/content/pdf/10.1007%2F978-3-642-40925-7_14.pdf
- [8] A. Podaras, "A non-arbitrary method for estimating IT business function recovery complexity via software complexity," presented at the *EEWC*, Jun. 2015. [Online]. Available: https://link.springer.com/content/pdf/10.1007%2F978-3-319-19297-0_10.pdf
- [9] E. Rorkowska, "Rank ordering criteria weighting methods—A comparative overview," *Optimum Stud. Ekon.*, vol. 65, no. 5, pp. 14–33, 2013. [Online]. Available: http://repozytorium.uwb.edu.pl/jspui/bitstream/11320/21891/02_Ewa%20ROSZKOWSKA.pdf
- [10] *Societal Security—Business Continuity Management Systems—Requirements*, Standard ISO 22301, 2012. [Online]. Available: <https://www.iso.org/obp/ui/#iso:std:iso:22301:ed-1:v2:en>
- [11] R. C. Barros, A. De Carvalho, and A. A. Freitas, *Automatic Design of Decision-Tree Induction Algorithms*. Heidelberg, Germany: Springer, 2015.
- [12] L. Breiman, *Classification and Regression Trees*. New York, NY, USA: Chapman & Hall, 1984.
- [13] L. De Micheaux, *The R Software: Fundamentals of Programming and Statistical Analysis*. New York, NY, USA: Springer, 2013.
- [14] K. GrDąbczewski, *Meta-Learning in Decision Tree Induction*. Cham, Switzerland: Springer, 2014.
- [15] P. D. C. de Almeida and J. Bernardino, "Big data open source platforms," in *Proc. IEEE Int. Congr. Big Data*, Jun. 2015, pp. 268–275.
- [16] *R: A Language and Environment for Statistical Computing*, R Develop. Core Team, Vienna, Austria, 2015.
- [17] A. Podaras, "A rules based decision making model for business impact analysis: The business function criticality classifier," presented at the *EOMAS*, Jun. 2017. [Online]. Available: https://link.springer.com/content/pdf/10.1007%2F978-3-319-68185-6_8.pdf
- [18] A. L'heureux, K. Grolinger, H. F. Elyamany, and M. A. M. Capretz, "Machine learning with big data: Challenges and approaches," *IEEE Access*, vol. 5, pp. 7776–7797, 2017.
- [19] L. Ponnambalam, L. Wenbin, X. Fu, X. F. Yin, Z. Wang, and R. S. M. Goh, "Decision trees to model the impact of disruption and recovery in supply chain networks," in *Proc. IEEE Int. Conf. Ind. Eng. Eng. Manage. (IEEM)*, Bangkok, Thailand, Dec. 2013, pp. 948–952.
- [20] L. Rokach and O. Maimon, "Top-down induction of decision trees classifiers—A survey," *IEEE Trans. Syst., Man, Cybern. C, Appl. Rev.*, vol. 35, no. 4, pp. 476–487, Nov. 2005.
- [21] A. Abdelhalim and I. Traore, "A new method for learning decision trees from rules," in *Proc. IEEE Int. Conf. Mach. Learn. Appl. (ICMLA)*, Dec. 2009, pp. 693–698.
- [22] N. D. Marom, L. Rokach, and A. Shmilovici, "Using the confusion matrix for improving ensemble classifiers," in *Proc. IEEE 26th Conv. Elect. Electron. Eng. Isr. (IEEEI)*, Eliat, Israel, Nov. 2010, pp. 555–559.
- [23] B. P. Salmon, W. Kleyhans, C. P. Schwegmann, and J. C. Olivier, "Proper comparison among methods using a confusion matrix," in *Proc. IEEE Int. Geosci. Remote Sens. Symp. (IGARSS)*, Milan, Italy, Jul. 2015, pp. 3057–3060.
- [24] P. P. Shinde, K. S. Oza, and R. K. Kamat, "Big data predictive analysis: Using R analytical tool," in *Proc. IEEE Int. Conf. I-SMAC (IoT Social, Mobile, Anal. Cloud) (I-SMAC)*, Palladam, India, Feb. 2017, pp. 839–842.
- [25] A. F. A. Pinem and E. B. Setiawan, "Implementation of classification and regression tree (CART) and fuzzy logic algorithm for intrusion detection system," in *Proc. IEEE 3rd Int. Conf. Inf. Commun. Technol. (ICOICT)*, Nusa Dua, Bali, May 2015, pp. 266–271.
- [26] W. Zhang and Y. Li, "A post-pruning decision tree algorithm based on Bayesian," in *Proc. IEEE Int. Conf. Comput. Inf. Sci.*, Shiyang, China, Jun. 2013, pp. 988–991.
- [27] Y. Maboudian and K. Rezaie, "Applying data mining to investigate business continuity in petrochemical companies," *Energy Sources B, Econ., Planning, Policy*, vol. 12, no. 2, pp. 126–131, 2017, doi: [10.1080/15567249.2015.1076907](https://doi.org/10.1080/15567249.2015.1076907).
- [28] L. Zheng *et al.*, "Data mining meets the needs of disaster information management," *IEEE Trans. Human-Mach. Syst.*, vol. 43, no. 5, pp. 451–464, Sep. 2013.
- [29] S. Zou, Y. Tang, Y. Sun, and K. Su, "An identification decision tree learning model for self-management in virtual radio access network: IDTLM," *IEEE Access*, vol. 6, pp. 504–518, Nov. 2017, doi: [10.1109/ACCESS.2017.2768402](https://doi.org/10.1109/ACCESS.2017.2768402).
- [30] J. R. Quinlan, "Induction of decision trees," *Mach. Learn.*, vol. 1, no. 1, pp. 81–106, 1986.
- [31] K. A. Shastry, H. A. Sanjay, and H. Kavya, "A novel data mining approach for soil classification," in *Proc. IEEE 9th Int. Conf. Comput. Sci. Educ.*, Vancouver, BC, Canada, Aug. 2014, pp. 93–98.
- [32] M. K. Rahmat and S. Jovanovic, "Reliability and availability estimation of DC uninterruptible power supply systems using Monte-Carlo simulation," in *Proc. IEEE 13th Int. Conf. Ind. Inform. (INDIN)*, Cambridge, U.K., Jul. 2015, pp. 76–81.
- [33] C. Wang, Z. Du, Z. Liu, and Y. Liu, "Study on decision tree land cover classification based on MODIS data," in *Proc. IEEE Int. Workshop Earth Observ. Remote Sens. Appl.*, Beijing, China, Jun./Jul. 2008, pp. 1–6.
- [34] H. Lee and S. Kim, "Decision tree ensemble classifiers for anomalous propagation echo detection," in *Proc. IEEE Joint 8th Int. Conf. Soft Comput. Intell. Syst. (SCIS) 17th Int. Symp. Adv. Intell. Syst. (ISIS)*, Sapporo, Japan, Aug. 2016, pp. 391–396.
- [35] S. B. Kotsianitis, "Decision trees: A recent overview," *Artif. Intell. Rev.*, vol. 39, no. 4, pp. 261–283, 2013, doi: [10.1007/s10462-011-9272-4](https://doi.org/10.1007/s10462-011-9272-4).



ATHANASIOS PODARAS received the M.Sc. Diploma degree in informatics from the Department of Information Engineering, Faculty of Economics and Management, Czech University of Life Sciences, Prague, Czech Republic, in 2005, and the Ph.D. degree in information management in 2010. In 2013, he was a Post-Doctoral Researcher with the Department of Informatics, Faculty of Economics, Technical University of Liberec, Czech Republic. He was an IT Specialist, a Programmer, IT Project Manager, and a Business Continuity and ISO Project Quality Assurance Team Member with the Software Applications Division, Alpha Bank, Greece, from 2007 to 2013. He completed his research in the ICT crisis management domain in 2015. Since 2015, he has been an Assistant Professor and a Researcher with the Department of Informatics, Faculty of Economics, Technical University of Liberec, Czech Republic. He has published multiple articles in peer-reviewed journals and international conferences. His research interests mainly focus on databases, data mining, business intelligence, information management, UML, business process requirement analysis, business continuity management and IT system recovery in crisis situations, and IT project management.

• • •