# CONTENTS

# Operational Research in the Emergency Medical System of Romania

Ionuț NICA
The Bucharest University of Economic Studies, Romania
ionut.nica@csie.ase.ro

*The explosive development of the human society in contrast to the limited character of resources determines the need for successful implementation of mathematic models in the decision-making process concerning the use of available resources. One of the critical areas where the need for rigorous criteria for resource allocation is strongly felt is the medical field. This issue appears to be currently affecting the great majority of nations in the world, being considered one of the most important challenges for modern states. The limited amount of resources allocated to the medical system brings forward the importance of optimizing the decision-making process concerning this field using models able to reflect the increasing complexity of the medical system, its interactions with the human society and its dynamics, therefore providing the perturbation control and adjustment instruments. From this point of view, the economical and mathematical modeling of the social phenomena provides strong, elegant and rigorous tools for the description of medical system that appears to be organized as a cybernetic system with a high level of complexity, focused on maximizing the social utility, and allowing the use of cybernetic methods designed for diagnosing, developing automatic medical archives, reducing time consumption and increasing overall efficiency.*

***Keywords:*** Markov network, Poisson distribution, cybernetics models, optimal decision, emergency medical system

# 1 Introduction

The resource allocation issue, especially for key sectors such as the medical one, seems to become more and more of a problem for the great majority of nations in the world. The limited character of resources, those allocated to the medical sector included, points out to the necessity of optimizing the decision-making processes involved in resource allocation based on rigorous methods capable of avoiding unnecessary expense and maximizing social utility.

The high complexity of the medical system and its multiple connections to human society determine the need for careful observation of its functioning and malfunctioning in order to produce the best suitable tools for perturbation control and adjustment.

The use of cybernetic and mathematical models for managing the functionalities of the medical system may lead at first to a mismatch between human expectations and the results achieved by following the path indicated through modeling. However, even these differences may provide useful information for further improvement, so that the medical system may eventually achieve the optimal resource allocation and maximize social utility.

In the pursuit of these goals (optimal utility/resource consumption ratio), cybernetic methods prove their utility for: designing diagnose systems for different clusters of diseases, developing automatic electronic medical archives (capable of minimizing searching times), as well as increasing overall efficiency by limiting waste, while keeping up with the evolution of biocybernetics.

## 2. The Romanian Medical System in the European Context

Due to Romania joining the European Union, the Romanian medical system started to use the similar structures in the EU as a reference point, which led to the need to increase efficiency by using mathematical methods applied in other European

countries. This would be the only way for the Romanian medical system to function as required by the new parameters implied by the efficiency standard implemented by the member states. According to data gathered in order to prove the need for implementing these changes, reforming the medical system using mathematical methods is more than critical, considering that the malfunctioning of the medical systems causes more than 60 000 deaths per year, which is the equivalent population of a small city.

Another challenge for the Romanian medical system is the shortage of medical staff (doctors, dentists, nurses, pharmacists), as compared to the other countries in the EU. Even though in Romania the financial efforts for sustaining almost all types of medical services have recently increased, there is still a general feel of system failure. By comparing the Romanian medical system to those of other European countries, and even by comparing the medical services provided in different regions in Romania, one can easily notice the massive differences in the access to medical services, as well as the gap between the values of most of the medical indicators, all picturing a worrying situation for the health of Romanians.

Romania has the highest mortality rate in the EU for both men and women, and this can be related to the difficult access to medical services of the general population, as well as to the fact that we have the smallest number of doctors, nurses and pharmacists reported to the size of the population. Moreover, in rural areas, where more than half of the population is located, there is even less medical personnel, and almost no functional hospital whatsoever.

According to a study by Ajay Tandon, Christopher JL Murray, Jeremy A. Lauer and David B. Evans in 2000 [0], Romania ranked 99 out of 191 countries in terms of medical system global performance.

However, the percentage of the GDP allocated to the health system is not definitory for its efficiency, considering that the USA ranks 37, even if they have the highest percentage of the GDP allocated to health. Still, compared to the European average, Romania allocates to the health system only a third of the amount spent on health by other countries. For instance, in 2010, Romania spent 600 euros per capita, as compared to the 1800 euros per capita which is the European average, and the government only directs 4% of the GDP towards the health system, while in France the percentage is 11% and the European average is around 8%. The difference can be explained by considering the low number of tax payers (only approx. 30% of the total population) and corruption, and it is a direct indicator of the struggle Romania has to put up in order to improve the efficiency of its medical system.

Since 2007, international mobility became even more accessible for Romanians, especially medical staff: almost 10% of the doctors decided to emigrate to countries such as France, Germany or Sweden, allured by the latest medical technologies available there and the high wages, and this percentage is still growing.

Even when it comes to medical equipment, Romania is one of the lowest ranked countries in the EU regarding the use of modern medical technology, being severely underequipped, which brings up even more the necessity to optimize the use of the existing resources, considering the low capacity for budgetary investment in medical technology.

However, almost every nation in the world has yet some challenges to deal with when it comes to the health system: no country has enough resources, money or medical personnel to cover all medical needs. More and more people have to live in fear of getting ill and not being able to access medical care. Therefore, there is a real necessity to improve efficiency of the medical system in order to optimize the use of the existing resources and to meet as

much as possible the demand for medical services.

Beyond the theory of providing public health services and the specific legal frame, a mathematical approach of the issue is also recommended, since the medical system is a good example of a cybernetic system [1], [15], [17], which includes not only complex components, but also dynamic and sensitive interactions that need careful planning.

**The Legal Frame For The Functioning Of Emergency Units**

*Law 95/2006* states that qualified first aid should be provided within:
a) 8 minutes for urban areas, in at least 90% of the cases;
b) 12 minutes for non-urban areas, in at least 75% of the cases.
The emergency medical care service should be organized in such manner that the maximal time for an intervention must not be longer than:
a) **15 minutes**, for emergency and intensive care units in urban areas, in at least 90% of the cases;
b) **20 minutes**, for emergency and intensive care units in rural areas, in at least 75% of the cases.
In order to implement an integrated emergency services management at regional level, all hospitals within the region should be included in a network, each network consisting in one regional first degree emergency hospital and several 2$^{nd}$ and 3$^{rd}$ degree local emergency hospitals. Moreover, the emergency and paramedic services department has to function around the clock in 12 hours shifts.
The mobile intensive care and emergency services provided by **SMURD**[1] have to abide a series of restrictions as well: emergency teams should consist of at least 4 people, including a driver/firefighter and doctor trained in intensive care and traumatology, the rest

---

[1] *Mobile Emergency Service for Resuscitation and Extrication*

of the team being supplied by other emergency structures, local authorities and local hospitals, or specially trained volunteers.

Also, according to law, the emergency and first aid structures in charge of the mobile units are responsible for providing functional medical equipment and drugs for the care of at least 20 critical patients.

*Order no. 1706/2007* states that:
• County capital cities with less than 500 000 inhabitants have to provide at least one Emergency Unit or Emergency Department within the county hospital. In case there is a regional or county children hospital, it must include a paediatric emergency department.
• Paramedics transporting critical patients have to report the emergency with at least 10 minutes before arriving at the emergency unit and provide all necessary information regarding the medical condition and treatment received by the patient in question.

According to the *National Triage Protocol*, the medical staff has to evaluate every patient presented to the emergency unit in order to determine the severity of the emergency and the urgency of accessing the medical services of each individual, and the average triage time should be 2 minutes or less. The triage procedure has to take into account two very important parameters:
• The time the patient was registered by the triage personnel;
• The time of the first medical consultation.

Since doctors are tempted to perform thorough examinations, and therefore become unavailable for patients who might need urgent interventions, the triage procedure is performed by other specialized personnel in order to optimize the use of doctors' time.

The National Triage Protocol states that patients registered to the emergency unit have to be sorted in order to be included into one of the following emergency levels:

Level 1 – CPR (code red): special room with life support equipment and defibrillator.

- The patient requires *immediate* life saving intervention.
- Time to be admitted in treatment area: 0 minutes.

Level 2 – Critical (code yellow): first degree emergency room.

- The patient is in severe pain or major disconfort, is of high risk or is in an altered mintal status.
- Time to be admitted in treatment area: 10 minutes.

Level 3 – Urgent (code green): 2nd degree emergency room.

- Stable patient requireing 2 or more of the resources defined in the Triage Protocol.
- Time to be admitted in treatment area: 30 minutes.

In case the time to be taken over by a doctor exceeds 15 minutes or there are changes in the patient's status the triage algorithm is repeated in order to update the procedures necessary for the patient in question.

Level 4 – Non-urgent (code blue)

- The patient is stable and requires the use of only one of the resources described in the Triage Protocol.
- Time to be admitted in treatment area: 60 minutes.

Level 5 – Consult (code white)

- The patient does not require emergency medical assistance and none of the resources described in the Triage Protocol.
- Includes people coming to the hospital for:
  - ✓ Getting vaccine shots;
  - ✓ Administrative resons such as medical permits, prescriptions etc.;
  - ✓ Social cases without medical complications;
  - ✓ Time to be admitted in treatment area: 120 minutes

In order to avoid overloading the Emergency Unit, the triage area can accommodate some of the medical procedures, so that the time to solve all cases is minimized.

Given all these restrictions and constraints, the need for a resource management algorithm becomes obvious, even more so considering the challenges the Romanian health system has still to deal with.

## 3. The mathematical model

Since emergency departments always focus on the quality of the medical services they provide, the 4 hour target set for a patient's waiting time is of critical importance, and this raises a problem that still has to be solved properly: how to allocate the human resources in order to meet this target.

The data collected from different emergency departments prove that most of them manage to meet the target and many other are close enough, but considering that lately the number of accidents increased, maintaining this target is still a priority. The so-called staffing algorithm appears to ease the managers' decision-making process regarding the efficient allocation of human resources in order to decrease the waiting intervals.

In all emergency units, the number of patients is a time variable, depending on the time of day, on the day of the week and even on the season (it is expected to have more patients presenting fall-related injuries and broken bones during winter, given the weather conditions). Therefore, the staff allocation differs on various intervals during a day.

The purpose of this research is determining the need for medical personnel (doctors, nurses, lab technicians, triage specialists etc.) for each interval during a day in order to reach the 4 hour target. [0] [0] [0] However, finding the optimal algorithm for such a complex system as the medical one is difficult, since the parameters taken into consideration (especially the patient's arrival time) are not constant. This characteristic was noticed by several specialists and therefore, there are multiple approaches to the matter. Unfortunately, all the approaches so far focused on single service systems, such as call-centers. The allocation of the human resource in an emergency department

is much more complex, because of the nature of the medical services: every patient arrived at the emergency unit is submitted to multiple tests, and therefore needs a number of different resources. Moreover, the severity of the case is also an important factor determining the access to certain resources, and resources can be used for more than one patient at a time, which means that, in order to be effective, the allocation algorithm has to take into account all these constraints. The suggested heuristic algorithm uses models based on waiting queues to estimate the amount of resources needed and the loading time for each resource in the system in order to optimize their allocation, while the quality of the medical services is measured through the probability for delays. The model includes a waiting queue M/M/1 with a single serving station [3] [4], with arrivals determined by a Poisson process and an exponential serving time. By using Wolfram Mathematica 9.0, the algorithm is implemented in order to optimize the allocation of doctors, nurses and triage specialists and maximize the number of patients treated.



**Fig. 1**. The possible flow for a patient
(Source: Author prelucration)

According to the Romanian emergency procedures, the patients arrived at the ER are sent to the triage room, where they are sorted (using the National Triage Protocol) by the severity of the case and distributed to the medical disciplines required by their specific situations. This is where the decision regarding further examinations is made, and the patients can be submitted to additional testing, EKG or ultrasound examination. The results of these investigations determines whether the patient will be admitted or discharged. The possible flow for a patient who arrived at the emergency room is represented in Figure 1.

**Table 1**. *The flow of patients to the emergency room*

|  | 00-04 | 04-08 | 08-12 | 12-16 | 16-20 | 20-24 | weekly | Percentage per ii | Total / hour |
|---|---|---|---|---|---|---|---|---|---|
| CPR | 1 | 2 | 3 | 0 | 5 | 4 | 23 | 0.02 | 3 | 0.5 |
| Major emergency 1 | 64 | 86 | 90 | 60 | 40 | 26 | 374 | 0.33 | 53 | 8.83 |
| Major emergency 2 | 70 | 42 | 156 | 184 | 164 | 120 | 736 | 0.65 | 105 | 17.5 |
| Total | 135 | 130 | 257 | 252 | 209 | 150 | 1133 | 1 | 161 | 26.33 |

(Source: Author computation)

The table above (table 1) contains data regarding the patient flow at the emergency room within the Bucharest Emergency University Hospital. During the observed interval, the emergency unit registered a total of 1133 patients per

week, 161 patients per day and 27 patients per hour, included in three different emergency categories: CPR, Major Emergency type 1 and Major Emergency type 2, in order of decreasing priority. The highest number of patients is registered in the 8-12 interval, followed by the 12-16 interval and the 16-20 interval.



**Fig. 2.** Variation of the flow of patients by hourly intervals
(Source: Authors computation)

The variation of the patient flow during the 24-hour interval



**Fig. 3**. Output of Performance Measures using Poisson distribution.
(Source: Author computation)

In Figure 3 it is represented the variation of the patient flow during the 24-hour interval. We calculated the network load by determining MeanSystemSize. In the analysis we also took into account the fact that some of the patients are also admitted, thus diminishing the outputs. Thus, the number of medical interventions is 106/4, for one hour. Also, upon leaving the emergency reception network it is found that 105/4 (for one hour) hours of people / hours are waiting to be consulted at different offices. The average waiting time in the network is about 35 minutes and the number of people waiting is 105/4 hours.Based on the statistic data gathered using weekly reports, the average number of people arriving at the emergency room is 28/hour. The interval considered for simulation is 4 hours, since in Romania, reports are filed every 4 hours.

Considering that the events requiring medical response follow a Poisson distribution with an average of 28 and that serving times are exponentially distributed, a network structure can be identified, leading to the processing of approx. 106 people per simulation cycle. Moreover, there are 105 people waiting in the network, the average time for using a node in the network being of 35 minutes/patient.

Aiming to minimize the probability of delays ($\alpha = [1 + \beta \frac{\Phi(\beta)}{\phi(\beta)}]^{-1}$, where $\phi$ and $\Phi$ are the density function, respectively the normal standard distribution function) and considering that in fact it is approx. 50% (50% of the cases are solved within the system), the resulting loading coefficient is 56%.

In order to determine the number of personnel necessary in an ideal situation, a generic network is considered, characterised by Poisson [7] [9] random entries with exponential serving times and variable number of serving stations for the considered intervals. Once registered, a patient can be transferred to any of the medical specialties or can be discharged (either admitted or case solved and discharged), which means that the Romanian medical system appears to be organized as a Markov chain structure with probabilistic flows between the different medical specialties. [13] [14]

Based on the data gathered at the Bucharest Emergency University Hospital, the state transition matrix is:

$$\begin{bmatrix} 0 & 0,5 & 0,2 & 0,3 & 0 \\ 0 & 0 & 0,57 & 0,13 & 0,3 \\ 0 & 0,4 & 0 & 0,4 & 0,1 \\ 0 & 0,2 & 0,3 & 0 & 0,5 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

The medical emergency network includes three major states: a patient demanding medical care is called a potential patient and is identified when entering the system. Anybody entering the system has no possibility of going through it unless it transitions to one of the following states, which explains the transition probability associated to the first column of the matrix (0). The state corresponding to a person transitioning the stations associated to the specialized medical disciplines is the consultation state. Therefore, a person has a 0,5 probability to access the services of the medical lab, 0,2 to enter the EKG room and 0,3 for getting an ultrasound. Anybody entering the system cannot exit without going through at least one of the specialized consultations.

A patient entering the consultation state leaves the medical lab to enter the EKG room with a probability of 0,57, the ultrasound room with a probability of 0,13 or may get discharged with a probability of 0,3, if the case is considered closed by the emergency system. The patient getting an EKG may return to the lab with a probability of 0,4, be transferred to ultrasound with a probability of 0,4 or be discharged with a probability of 0,1.

Because the path the patient follows within the emergency department is rather complex, this is considered the cause for the queues that affect the efficiency of the Romanian medical system.

A patient sent directly to ultrasound may leave the system with a probability of 0,5, but can return to the lab or EKG with a probability of 0,2, respectively 0,3. A patient reaching the final transition state is considered cured (case solved), with no possibility of going back to any of the prior states, situation that is reflected on the bottom line of the matrix, which only contains 0 valued elements.

**Fig. 4**. The distribution of Markov network.
(Source: Author computation)

In Figure 4 s represented the system with infinite number of service stations, the average number of personnel needed is 4. In the graphic representation we simulated on a time horizon of 100 weeks, the allocation of the resources used by the emergency system. The blue line represent the doctors, the purple line are the allocation of the nurses, the yellow line represent the laboratory technicians and the green line is represented by the EKG specialists.As evidenced by a large number of scientific papers and studies, the Markov network structure is generally non-stationary; therefore it can be decomposed in serving networks with deterministic paths, equivalent to the situation of a single serving station with random entries but with unlimited order processing capacity (unlimited number of serving stations).[11] [12] [16]

This is justified by the sorting of patients into groups, by the type of consultation required, which makes the network sequential, with deterministic paths, so that the number of untreated patients is easily calculated by adding up the average unused factor of the station in question.

The ideal allocation of human resources for the emergency system is determined by taking into consideration a queue with an average of 28 entries (patients) per hour, which is equivalent to 112 patients for a 4 hour interval. This way, the system is used to its full capacity, which means that all patients arrived are treated and the average number of specialized consultation rooms in use in the network is 4. Theoretically, within such structure, the number of treated patients is undetermined, this being the ideal situation.

Considering the number of necessary specialists as 4, the model determines the over- and under-load of medical resource for each hour, as a difference between the actual number of doctors existing in the system and the ideal number of doctors necessary at a given time. Over-loading the system leads to an increase of the budgetary costs associated to the system, so this occurrence is penalized with the factor $p^o = 1$. Still, this situation is not that bad, since theoretically a larger number of doctors in the system would consequently mean a larger number of patients treated and an increased efficiency of the medical emergency system. The under-load of the system, on the other hand, means reducing the budgetary costs by

reducing the main resource in the system (doctors), which may lead to a diminished efficiency of the system. Therefore, this situation is penalized with the factor $p^u = 2$.

The available human resource is allocated so that the total number of people necessary per shift equals the total number of people allocated on that specific shift. Also, the difference between the variable measuring the overload and the variable measuring the underload should match the difference between the actual situation and the ideal situation resulted from the model. In the specific case of the Romanian medical emergency system, conventionally, the work hours are grouped in three 8-hour shifts. In this particular situation, the differences described above are considered positive (there are more doctors than in the ideal situation), while the maximum value for each resource in the system is considered 100 (maximal number of people available in the ER).

The model considers 49 constraints, the objective function being the minimizing of penalties paid for over- and under-loading the emergency medical system. Considering that the decision variables are whole numbers, the result is a whole number programming case, which can be solved using the Mathematica software. The solution offered by the model is that the first shift should be covered by 23 people, the second by 38 people, while for the third shift, 37 will suffice.

**Table 2**: *The allocation of the medical staff*

| Time frame | 00-08 | 08-16 | 16-00 |
|---|---|---|---|
| Doctors | 6 | 9 | 8 |
| Nurses | 11 | 25 | 22 |
| Triage Specialists | 4 | 6 | 5 |
| Lab Doctors | 2 | 2 | 2 |
| Total | 23 | 38 | 37 |

(Source: Author computation)

Below, we evaluated the number of doctors needed in a system with an infinite number of service stations (equivalent to the ideal emergency system) is 4 doctors if it is considered that the number of patients is 26 per hour.



**Fig. 5.** Evolution of the flow of medical personnel. Simulation of the network in which the number of the patient in 26 and the number of doctors is 4 over a 100 – week horizon.
(Source: Author computation)

The algorithm for determining the necessary human resource consists of 5 steps:

- Step 1: Initialize a = 1;
- Step 2: Calculate parameter P, considering the constraint $= [1 + \beta \frac{\Phi(\beta)}{\phi(\beta)}]^{-1}$, where $\phi$ and $\Phi$ are the density function, respectively the repartition function of the standard normal distribution.
- Step 3: Determine the need for each resource for each specific interval, using the function:
  $s_k(t) = [x + \beta \sqrt{x}]$,
  where $x = m_\infty{}^k(t)$;
- Step 4: Estimate the percentage of discharged patients based on the resources estimated at step 3.
- Step 5: If the percentage calculated at step 4 is not 98%, a can be increased or decreased and we go back to step 2. Otherwise, STOP and save the solution.

$$\min[p^o \sum_{j=0}^{23} \Delta_j^+ + p^u \sum_{j=0}^{23} \Delta_j^-]$$

$$\sum_{i \in I} p_{ij} x_i = a_j, j = 0, \dots, 23$$
$$a_j - s_j = \Delta_j^+ - \Delta_j^-, j = 0, \dots, 23$$
$$\sum_{i \in 1} y_i \le k$$
$$x_i \le M y_i, i \in I$$
$$x_i \ge 0, i \in I$$
$$x_i \ge 0, x_i - integer \ and \ y_i = 0,1, i \in I$$
$$\Delta_j^+, \Delta_j^- \ge 0, j = 0, \dots, 23$$

$(s_0, s_1, s_2, \dots, s_{23})$ - *levels of resource allocation generated by the algorithm*;

$p^o$ - Penalties paid for each over allocation of the human resource per hour;

$p^u$ - Penalties paid for each underallocation of the human resource per hour;

$\Delta j+$ - overload at hour j, j=0...23;

$\Delta j-$ - underload at hour j, j=0...23;

I - the shift range allowed considering the legal constraints;

$xi$ - the decision variable expressing the number of employees scheduled to work on a shift, $i \varepsilon I$;

$$p_{ij} = \begin{cases} 1 & \text{if the shift } i \text{ includes the hour } j \text{ as working time} \\ 0 & \text{if it doesn't} \end{cases}$$

$$\sum_{i \varepsilon I} p_{ij} x_i$$
$-$ *the total number of employees working at hour j,j* $= 0, \dots, 23;$

M - *A very high number*;
k - The maximal number of shifts;
$y_i$ – Artificial variable - value 1 if at least one employee works on shift i, $i \varepsilon I$;
The suggested heuristic algorithm [8] [10] uses models based on waiting queues to estimate the amount of resources needed and the loading time for each resource in the system in order to optimize their allocation, while the quality of the medical services is measured through the probability for delays.[18]

## 4. Conclusions

To conclude, by comparing the Romanian emergency system to the British one, the first is found wanting, with a satisfaction level of only 56%, as opposed to 93% for the latter. According to the European standard, a medical emergency system is considered of high quality if it has a 98% quality indicator, which means that it can solve 98 cases out of 100. The Romanian emergency system is also affected by delays, the average time spent by a patient waiting for the different procedures necessary being 35 minutes, as opposed to only 7 minutes which is the European average. Since the main resource used in the emergency system is the human resource, the issue of staff allocation in the emergency departments is a major challenge for many medical systems. In the specific case of the Romanian medical system, it was found, using data gathered in a major Emergency Hospital, that a mathematical model is extremely helpful in determining the necessary human resource per shift, considering multiple factors such as the patient flow per hour and the budgetary constraint.

## References

[1] Ajay Tandon, Christopher JL Murray, Jeremy A. Lauer, David B. Evans - *The Comparative Efficiency Of National Health Systems In Producing Health: An Analysis Of 191 Countries*, GPE Discussion Paper Series: No. 29, EIP/GPE/EQC World Health Organization, 2000
[2] Coats, T.J., Michalis, S., 2001. Mathematical modelling of patient flow through an accident and emergency department. Emergency Medicine Journal 18 (3), 190-192;
[3] Eick, Stephen G., Massey, William A., Whitt, Ward, 1993. The physics of the Mt/G/1 queue. Operations Research 41 (4), 731-742;
[4] Feldman, Zohar., Mandelbaum, Avishai., Massey, William A., Whitt, Ward, 2008. Staffing of time-varying queues to achieve time-stable performance. Management Science 54 (2), 324-338;

[5] Fletcher, A., Halsall, D., Huxham, S., Worthington, D., 2006. The DH accident and emergency department model: A national generic model used locally. Journal of the Operational Research Society 58 (12), 1554-1562;

[6] Green, Linda V., Kolesar, Peter J., Whitt, Ward, 2007. Coping with time-varying demand when setting staffing requirements for a service system. Production and Operations Management 16 (1), 13-29;

[7] Green, Linda V., Soares, Jao., Giglio, James F., Green, Robert A., 2006. Using queueing theory to increase the effectiveness of emergency department provider staffing. Academic Emergency Medicine 13 (1), 61-68;

[8] Mayhew, L., Smith, D., 2008. Using queuing theory to analyze the Governments 4-h completion time target in accident and emergency departments;

[9] Gunal, M.M., Pidd, M., 2009. Understanding target-driven action in emergency department performance using simulation. Emergency Medicine Journal 26 (10), 724-727;

[10] Mortimore, Andy, Cooper, Simon, 2007. The ''4-hour target: Emergency nurses'' views. Emergency Medicine Journal 24 (6), 402-404;

[11] Munro, J., Mason, S., Nicholl, J., 2006. Effectiveness of measures to reduce emergency department waiting times: A natural experiment. Emergency Medicine Journal 23 (1), 35-39;

[12] Sinreich, David., Jabali, Ola., 2007. Staggered work shifts: A way to downsize and restructure an emergency department workforce yet maintain current operational performance. Health Care Management Science 10 (3), 293-308;

[13] Sinreich, David, Yariv, Marmor, 2005. Emergency department operations: The basis for developing a simulation tool. IIE Transactions 37 (3), 233-245;

[14] Whitt, Ward, 2007. What you should know about queueing models to set staffing requirements in service systems. Naval Research Logistics 54 (5), 476-484;

[15] Chiriță, N., Nica, I., 2019. Cibernetica Firmei. Aplicații și Studii de Caz. Ed. Economică;

[16] Ashour, O., Kremer G., 2013. A simulation analysis of the impact of FAHP-MAUT triage alghorithm on the emergency department performance measures. Expert Systems with Applications 40(1), 177-187;

[17] Nica, I., Chiriță, N, Fabian, C., 2018. Analysis of Financial Contagion in banking network, 32nd International Business Information Management Association Conference.

[18] Shih, C. L. and S. Su. 2003. Modeling an emergency medical services system using computer simulation. International Journal of Medical Informatics 72(3),57-72

**Ionuț NICA** (b. May 2, 1992) has graduated from the Faculty of Cybernetics, Statistics and Economic Informatics at the Bucharest University of Economic Studies in 2014. He followed a master's degree in Cybernetics and Quantitative Economics, within the same faculty. Currently, he is a PhD student, teaching assistant in the department of the Faculty of Cybernetics, Statistics and Economic Informatics and work in bank as Basel II Expert in the department of Retail Credit Risk Methodology and Validation. He has high interest in areas such as Cybernetics, Operational Research, Economic Dynamics, Applied Mathematics and Big Data.

# Business Analytics Applications for Consumer Credits

Claudia ANTAL-VAIDA
The Bucharest University of Economic Studies, Romania

*The fast-paced and dynamic economical background determines all the industries to quickly adapt to change and adopt emerging technologies to remain competitive on the market. This tendency led to high volumes of data generated each second and to a decreasing ability of the manpower to analyze it and use if for beneficial purposes. This paper reviews the impact of Digital Transformation on the Banking area and how financial institutions leverage the advantage created by this trend, especially in the credit risk management field. Multiple papers on consumer credit scoring models written after the financial crisis from 2007 were reviewed and their results were summarized in this article, to increase the accuracy of future analysis by leveraging the already known results.*
**Keywords:** *Business Analytics, Machine Learning, Banking, Credit Risk Assessment, Scoring Models, Consumer Credits*

## 1 Introduction

Considering the current economic background, public and private companies are facing lots of challenges and difficulties when setting up their growth strategy due to the highly competitive environment where they operate. In order to overcome them, these entities require changes in their approach when making decisions, relying more on historical facts. Therefore, these companies should change focus – instead of intuition, they should aim to become data-driven companies, but for this approach to be successful, they need to invest in data-processing and methods of generating insights to get valuable information.

Change is all around us and if we analyze how fast technology, economy, medicine and other fields are evolving, we will notice that data and data-analysis are key players in it. As most of the industries are highly developed nowadays, any organizational decision should be taken based on results generated by business analytics to reach an ascending trend and a sustainable growth. Fortunately, companies understood the importance of data and the business value add they can bring when used correctly and consistent.

Gartner defines the Analytics area as a bunch of methods and techniques used for building analytical models and simulations to understand reality and predict the future state of a system. [1] It is a must-have area due to its impact on financial growth and profitability by improving company's competitive advantage and minimizing the error when making a decision. Examining the growing interest in this area and the resources invested in research by multiple companies, we can deduce that the popularity of this field has exponentially increased and companies acknowledged the impact of more sophisticated analytical decision-making tools for creating new opportunities, choosing the timing and enhancing the know-how.

The Banking sector is facing similar changes and needs to invest in business analytics as well. Its focus should be on understanding customer's behavior, how to improve interactions with customers and what motivates him to carry through their obligations. [2]

Even though this industry seems to be complicated, its activity is quite straightforward: the banking sector handles cash, credits and secures funds aiming to make more money out of them. Its products mainly consist of credit and debit accounts, loans and mortgages, helping people and companies to invest in their future, especially when they do not have enough liquidities to do that.

Players in each industry have a key focus on improving profitability, looking into cost reduction, revenue increase and process optimization and the banking field does not make an exception. Having said that, Business Analytics can play a key role in achieving those goals and there are already examples on the market which exceeded expectations. Some applications worth mentioning would be the algorithms of predicting risks, that analyze transactions and identify uncommon behavior of a customer or the credit-scoring algorithms that predict whether a customer will pay his debts, or he will miss them.

## 2. Algorithms used in the Banking Area

The following part of this paper explains why and how the concept of big data appeared in the banking field and how companies decided to leverage the advantages created. All the data created each second increased the influence of Business Analytics in business applications, with lots of solutions already implemented and many others being explored.

### 2.1. Digital Transformation of Banking

The Digital Era is characterized by new technologies which increase the speed of knowledge turnover [3]. The emerging development of analytics, cloud, social media and mobility technologies caused overall disruption across industries, Banking being one of the most impacted fields by this trend. [4]

Companies understood the importance and impact of technology in their activity, therefore they started to invest in research and development labs that use not only social media analytics, machine learning algorithms and big data, but also research on possible innovative scenarios that leverage artificial intelligence, robotics, automation, advanced data visualization and not only.

American Express, one of the giants in the Banking and Financial services from US, has set up a new tech lab to focus on big data, cloud computing, analytics and mobile technologies, as well as on futuristic ones. Thru this initiative, they aim to analyze customer's behavior on the market and quickly respond to it with customized products. [5]

Another example of company acknowledging the importance of new technologies and investing in its development would be the Fidelity Investments, one of the largest asset management companies in the world. In 2014, they announced the opening of Financial Labs, a research unit that will partner with The Massachusetts Institute of Technology (MIT) and Stanford University to get the "outside view" and develop innovative applications. [6]

These were only few examples of Financial companies investing into the tech area. The trend in the industry is to digitalize as much as possible and behave almost like tech companies. Banks must be quick in converting an idea into a service in order to survive on the market but also stay relevant in the years to come. Moreover, technology offers new opportunities to address untouched markets, by simplifying the communication and removing the geographical barriers. [7]

From previous examples, we draw the conclusion that companies have become aware of the importance of technology and the role they play in business performance management. Moreover, they started to massively invest in developing new technical capabilities to optimize their processes and improve relationships with customers.

What's the result of these investments? How digitalization translates in the real world? It results in high volumes of unstructured data that are harder and harder to analyze manually and that's the perfect scenario for Big Data and Analytics to come into the picture. As data is not meaningful enough, companies had to identify ways of converting data into information and insights to monetize digitalization.

### 2.2. Machine learning at a glance

Business Analytics represents the use of data, technologies, statistics, mathematics and computer-based models to help

management understand the business, solve issues and make fact-based decisions. [8] This area has four stages that have different business impacts, depending on their complexity and level of knowledge required [9]:

- **Descriptive Analytics** is answering the question "What happened?", by analyzing and displaying historical data in reports and dashboards to simplify the decision-making process;

- **Diagnostic Analytics** researches the causes and effects of a certain event in order to avoid it or increase its frequency in the future, depending on the impact. It usually answers the question "Why did it happen?";

- **Predictive Analytics** provides an answer to the question "When can it happen?" and implies statistical methods and Machine Learning techniques;

- **Prescriptive analytics** in the area that recommends decisions based on simulations and process optimizations, trying to answer the question "How can it happen again?".

Going forward, this paper will focus more On Machine Learning and its applicability. **Machine Learning** is defined as a tool or mechanism that uses statistical models to facilitate the solutioning of a problem, by studying the past behavior, identifying patterns and constantly improving itself based on the data analyzed. Its main purpose is to develop an adaptable algorithm to solve an issue, despite the external variables that might influence it or its complexity. [10] There are different types of algorithms used by machine learning listed below [11]:

- **Supervised algorithms** know the outcome from the beginning and its learning is guided by human observations and feedback thru tags and labels inputted from the beginning;

- **Unsupervised algorithms** rely exclusively on clustering separate data based on similar features and modifies the calculation process to respond to initial inputs; this type does not involve

any external feedback, nor tags to be considered for data processing;

- **Reinforced algorithms** are about taking the most appropriate action to maximize the output, regardless the situation. If in supervised learning the expected outcome is known from the beginning, this type of algorithm has the liberty to decide what is best to do to perform a task and tends to learn from its own experience.

Even though Machine Learning sounds evolutionary and promises to revolutionize the way things work, having a positive impact on the areas where applied, it has both advantages and disadvantages that should be considered when deciding to use them for a real use case.

Thus, few of the advantages worth mentioning of Machine Learning are: [12]

- **Easily identified trends and patterns**. When given a large dataset to analyze, the Machine Learning algorithms can quickly identify specific trends that might not even be obvious to human beings.

- **Constant Improvement**. While exposed to computation of data, the ML algorithms gain experience and improve efficiency and accuracy, leading to more reliable decisions and results.

- **Various applications.** No matter the area you work on, you will find for sure an application that would involve Machine Learning algorithms and would be beneficial for your area.

On the other hand, Machine Learning has some limitations that should be known before deciding to use and invest in them:

- **Requires high volumes of data**. For excellent results, Machine Learning algorithms require high volumes of data to train on. Besides volumes, data quality of the train dataset plays an important role as well, outputs being highly dependent on inputs.

- **Time and Resources**. Depending on the complexity of analysis you want to perform, Machine Learning algorithms

may require time to learn and massive computer resources.

- **May produce biased results**. Machine Learning algorithms are highly susceptible to error, but this aspect mostly depends on the diversity of the dataset it trains on. If the train set is small and not inclusive enough, the results might be biased, leading to irrelevant interpretation.
- **Interpretability**. Unfortunately, it is hard to understand the reasoning behind a decision taken by an algorithm, reason why these are considered "black-boxes".
- **Scalability**. Once a model is proven to be efficient, companies implementing it need to overcome the challenge of scaling it. This can become expensive due to the resources required, the need of further optimizing it and integrating it with other systems.

All these challenges can be overcome if the company is willing to invest and few of them might not even appear, it always depends on the use case. Therefore, Machine Learning usage can surface the potential and value of unstructured data in each company. Its ability to form adaptive behavior in the process, without being programmed for this, makes them incredibly powerful when used for the right processes and analysis.

To sum up, the main benefits of the ML algorithms in this area is given by the complexity of the analysis performed (due to various parameters considered), reduction of approval time, less human resources involved, thus avoiding the human error, fraud of subjectivism

## 2.3. Algorithms for Risk Management

Risk management has become more important in the banking fields since the global financial crisis took place, moment since when banks started to research on how risks can be detected, measured and managed . [13]

There are different types of risks in the financial area that can be addressed by Machine Learning, but we will analyze the applicability only for three of them: Operational risk, market risk and credit risk.

**Operational risk management** assumes that a firm wants to foresee the direct or indirect risk of financial loss due to a host of potential operational breakdown. [14] The risk can be triggered by internal factors (people, system, deficient processes) or by external ones (global economic background, frauds, operational errors). Considering the increasing variety and complexity of risks, especially for financial institutions, machine learning and artificial intelligence applicability increased consistently and started to play a key role in predicting these events, assessing their impact and minimizing their effects. [15]

Banks pursuit the evaluation of the best ways to protect their data, systems and clients and machine learning can support that. Process automation can increase the execution of routine tasks, minimize human error, analyze data to outline the relevant content and increases the ability to evaluate risky clients and networks. Machine Learning can also generate and prioritize alerts for uncommon activities and asses the risk involved.

Another risks that is worth investigating is the **Market Risk** to which the banks (and not only) are exposed due to investing, trading and playing on the financial markets. Machine Learning is mainly proper for identifying inadvertent risk in trading behavior, for understanding the impact of the firms that trade on their market price, for establishing new patterns and connections between assets and how they influence each other or even for creating bots to constantly monitor the financial indicators and send alerts once a trade would be profitable.

Finally**, Credit risk** is one of the highest risks faced by banks and usually the one requiring the most capital, therefore its management is of high interest for the financial institutions.

The objective of credit risk management is to optimize the credit portfolio and reduce the risk of customers not meeting their obligations. The high and extensive complexity of credit risk assessment made

this area proper for machine learning applications.

## 3. Machine Learning for Consumer Credit

Utilizing Machine Learning techniques is not a new trend, but it is a growing one. Back in the '90s, a comparative analysis between traditional statistical models of distress and bankruptcy prediction and an alternative neural network algorithm proved to be an effective combinations, with a significant increase in accuracy. [16] And the research in this area just started at that point. Over years, there were multiple implementations of machine learning techniques supporting risk management, which proved to be very efficient, making the most optimal decision.

The following part of this paper mainly focuses on the applications that were developed and deployed for consumer credits in the risk management area after the financial crisis.

### 3.1. Scoring Models Overview

One of the tools most used in the credit risk management are the credit scoring models, defined as statistical methods that consider financial indicators to predict the default risk of individuals or companies. These indicators are given a relative importance and are considered when predicting the creditworthiness, pointing out the probability of default of the borrower. [13] In Table 3.1. are listed in a chronological order all the articles considered in this paper with the utilized algorithm(s). They are all tackling Credit Scoring models for Consumer Credit, reason why they were considered:

| Article | Author(s) | Year | Algorithm |
|---|---|---|---|
| *Credit scoring with a data mining approach based on support vector machines* | Huang, Chen, Wang | 2007 | Hybrid Support Vector Machines |
| *Support vector machines for credit scoring and discovery of significant features* | Bellotti, Crook | 2009 | Support Vector Machines |
| *Consumer Credit Risk Models via Machine-Learning Algorithms* | Khandani, Kim, Lo | 2010 | Classification and Regression Trees (CART), Linear Regression |
| *A Proposed Classification of Data Mining Techniques in Credit Scoring* | Keramati, Yousefi | 2011 | Artificial Neural Networks; Bayesian classifier; Discriminant Analysis; Logistic regression; K-Nearest Neighbor; Decision Tree; Survival Analysis; Fuzzy rule-based system; Support Vector Machine; Hybrid Models |
| *Loan Default Prediction on Large Imbalanced Data Using Random Forests* | Zhou, Wang | 2012 | Random forest |
| *Quantitative credit risk assessment using support vector machines: Broad versus Narrow default definitions* | Harris | 2013 | Support Vector Machines |
| *Large Unbalanced Credit Scoring Using Lasso-Logistic Regression Ensemble* | Wang, Xu, Zhou | 2015 | Lasso logic regression |
| *Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research* | Lessmann, Baesens, Seow, Thomas | 2015 | Artificial Neural Networks, Support Vector Machine, Ensemble Classifier, Selective Ensemble Classifier, Threshold metric, Area under receiver operating characteristics curve, |

| | | | H-measure, Statistical Hypothesis Testing. |
|---|---|---|---|
| *redit scoring with a feature selection approach based deep learning* | Van-Sang, Nguyen | 2016 | Deep Learning |
| *A novel multistage deep belief network based extreme learning machine ensemble learning paradigm for credit risk assessment* | Yu, Yang, Tang | 2016 | Deep briefing Network, Extreme Machine Learning |
| *Ensemble Learning or Deep Learning? Application to Default Risk Analysis* | Hamori, Kawai, Kume, Murakami, Watanabe | 2018 | Bagging, Random Forest, Boosting |

Table 3.1. Articles list on Scoring Models for Consumer Credit

### 3.2. Machine Learning Algorithms for Consumer Credit Models

First reference found is the study presented by Huang, Chen and Wang in 2007 that proposed a hybrid SVM-based credit scoring model that analyzes the credit applications based on the information provided by applicants. [17] This application would help the banks decide whether to grant credit to consumers who require it. The authors compared the SVM classifiers with neural networks, genetic programming and decision tree classifiers, and even though the number of imputed features was low, the result was comparable with the other techniques' outputs. Moreover, the SVM algorithm was proven to be very effective in the classification area, successfully grouping credit applications into accepted or not, hence minimizing the money lose due to underperforming credits. Lastly, the researched revealed that if the SVM classifiers are combined with genetic algorithms, they can serve multiple purposes: performing feature selection tasks but also optimize the model parameters.

Even though most of the conclusions were positive, the authors outlined some negative sides of the hybrid Support Vector Machine – Genetic Algorithms models too. One of their downsides would be the long time required for training and the high computational complexity demanded for a good classification accuracy. Another inconvenience would be the "black-box" nature of SVMs, but this could be overcome with the use of SVM extraction techniques or the use in combination with other interpretable models.

In 2009, another article focusing on SVMs for credit scoring was published by Belloti and Crook. They compared three traditional techniques (logistic regression, discriminant analysis and K-Nearest neighbors) against SVMs using a dataset of around 25000 customers. The results reinforced the outcomes of previous researches (that SVMs are successful in classifying credit card customers), but also revealed that they can be well applied selecting the features that have the highest impact on likelihood of default. They also discovered that a very important indicator in this analysis is the type of credit card, as this could influence the other variables to be examined in the model.

Though, one major disadvantage acknowledged is the high number of support vectors required for the best performance, mainly due to the broad indicative nature of credit data. [18]

In 2010, Khandani, Kim and Lo [19] published a study in which they were using machine-learning techniques for forecasting models of consumer credit risks. The framework used consisted of generalized classification and regression trees (CART)

By combining credit bureau data with customer transactions and applying linear regression $R^2$ with a delinquency rate of 85%, they reached a better accuracy of classification rates for the default of an

application. They also analyzed the patterns of the time-series of delinquency rates and concluded that aggregated consumer credit-risk analytics may have a high influence in forecasting systemic risks. Moreover, by applying machine learning forecast models on the decision to cut credit lines, they estimated a cost saving between 6% to 25%.

On the same topic, Keramati and Yousefi presented a study in 2011 in which they acknowledge the importance of analyzing the huge amount of data generated by credit scoring in a fast-growing credit industry, but also the human impossibility of manually reviewing and interpreting it, hence the need of data mining techniques to support this effort. In their paper, they analyzed ten different data mining approaches and their results outlined the following [20]:

- K-Nearest Neighbor (K-NN) proved to obtain the best results for the credit scoring purposes;
- The Employ Multi-Group Hierarchical Discrimination (M.H.DIS) resulted to have better classification abilities than the traditional models;
- Support Vector Machine – MARS (SVM MARS), logistic regressions and neural network are very good for classification, but LDA and CART are easy to use in building such a model,
- Integration of Self Organization Maps (SOMs) with supervised classification methods proved to bring more advantages.
- Kernel Based RBF neural network was the best choice in identifying the true positive.
- The comparison between discriminant analysis, logistic regression, neural network and regression trees for predictions and classification tasks outlined that CART and neural network are the best to apply for best results.
- When evaulating the accuracy of K-NN, SVMs and neural networks, it resulted that the integration of all these methods with effective feature selection improved the accuracy of the classification.

One of the main conclusions they draw was that calculating the probability of default for an applicant if more meaningful that classifying them into the binary classes.

One year later, in 2012, Zhou and Wang proposed a study in which they applied improved random forest algorithms in the binary classification field, by attributing weights to the decision trees in the forest. These weights were calculated based on the previous performance, namely errors in training. Their approach outperformed expectations, beating the result of benchmark algorithms like traditional random forest, SVMs, KNNs in terms of accuracy and proved that parallel random forests can considerably reduce the learning time [21].

In 2013, Terry Harris published an article on Support Vector Machines (SVMs) applied for credit-scoring models from two perspectives: a broad one, considering the credits that are less than 90 days past due, and a narrow perspective, analyzing the credits that are more than 90 days past due, reaching the conclusion that the last produced more accurate, mainly for severe cases of default. The main explanation for this conclusion could be the greater number of cases fed to the model, leading to a better learning of the pattern for un-creditworthiness. [22]

Wang, Xu and Zhou published a new article in 2015, in which they outlined a new mix of algorithms for credit scoring that exceeded expectations and previously known results. Their approach consisted of applying clustering and bagging algorithms to generate balanced training data and diversify data, applying Lasso-logistic regression ensemble to evaluate credit risks. [23]

During the same year, Lessman, Baesens, Seow and Thomas updated the study started by Baesens et al., including new classification algorithms used in the credit scoring area. They considered 41 classifiers for 8 credit scoring data sets and their results proved that there are more performant classifiers than the standard logistic regression. More than that, they outlined the

business value add of improving the prediction models, variable selection and data quality and suggested that focus should change into these areas. [24]

One year later, in 2016, Van-Sang and Nguyen published a study on Deep Learning, a powerful classification tool that provides training stability, generalization and scalability with big data. This method surpassed results previously obtained with baseline methods and showed competitive performance with other feature selection models extensively used in credit scoring area. The study also outlined that fewer features considerer for the evaluation procedure allow for collecting essential variables, hence reducing the resources allocated on performing the research. Moreover, parallel processing proved to decrease processing time, whilst obtaining the same results. [25]

During the same year, Yu, Yang and Tang proposed a novel multistage deep belief network based extreme learning machine (DBN – based ELM) ensemble learning methodology as a promising mix for credit scoring problems. These 2 techniques were already known for the time-saving characteristics and for the high-learning capacity thru hidden layers. The structure of multistage ensemble learning model, working in three stages, conducted to better results than typical single classifiers. The steps followed for this analysis are the following: in the first stage, the bagging sampling algorithms are applied to generate multiple and diverse training subsets of data; during the second stage, the ELM is utilized as classifier and different ensemble components can be properly defined with right subsets and different initial conditions. The last stage merges the individual results to form the final classification output thru the DBN model, which can effectively outline the relevant information hidden in ensemble members. [26]

In 2018, Hamori, Kawai, Kume, Murakami and Watanabe published an article in which they assessed payment data and compared the prediction accuracy and the classification ability of three ensemble-learning methods with neural-network methods. The three methods assessed were bagging, random forest and boosting. The study outlined that the boosting method has a superior classifying ability, even when compared to neural networks. The performance of the lastly mentioned proven to be highly dependent on the choice of activation function, dropout inclusion and number of hidden layers. [27]

### 3.3. Discussions

All these articles are approaching the same topic from different angles and thru different methods. The observations and conclusions obtained by the authors can be further leveraged in analysis, hence improving efficiency and accuracy by using the already known results.

Overall, the main idea that was outlined by each article is that Business Analytics plays a key role in the evolution of the financial institutions and in the optimization of their processes. When applied, it minimizes costs by reducing the number of human resources involved and increases the accuracy of the decisions.

Credit scoring algorithms assigns numerical values to the client outlining whether the entity is likely to default or not. Most of the studies were focused on this area, treating it as a classification problem in order to facilitate the credit decision, but also minimize the credit risk exposure.

Machine Learning algorithms performed better than the traditional techniques in classification steps and obtained increased prediction accuracy. The SVMs were widely tested and proved to be very effective in the classification area and in the feature selection process, especially when combined with genetic algorithms. Another algorithm that exceeded expectations was the Random Forest applied in the binary classification area, which showed outstanding results and a reduced learning time. Deep learning was outlined as one of the tools that provides training stability, generalization and scalability.

Beside all these, another valid point that should be considered is the dataset used for researches: it should be varied, integer, divers, to cover as many scenarios as possible and reduce biased results. Having said that, if the availability of real data would increase, if would encourage more researches on evaluating all the problems encountered in credit risk management, and not only.

## 4. Conclusions

Companies acknowledged the importance of adopting new technologies and the business value it creates and it is expected that financial institutions will increase the machine learning applications in the risk management fields to enhance their capabilities.

Even though these applications have some known limitations, a major one being the inability of understanding the mechanism used to reach a decision (mentioned in the literature as "operating like a black box"), the business value it creates it significantly higher, main benefit worth outlining being: high complexity of analysis performed, limited number of human resources involved, minimal error and reduction of performing time.

Machine Learning proved to be evolutionary and promises to revolutionize the way things work, hence it has the potential to transform the risk management area and enables the discovery of complex, nonlinear patterns in broad datasets.

This paper introduced an assessment of the researches around credit scoring algorithms for consumer credit within the banking industry, mostly because credit risks is considered the highest risk for a financial institution. However, the advantages and disadvantages of various machine learning tools for credit scoring can be further studied to refine them, improve results and maximize their values.

In conclusion, even though there have been different studies performed in this area, there's still room for research and improvement to extend the beneficial

applicability and impact of machine learning in the financial field and not only.

## References

[1] "Gartner Glossary," Gartner, Available: https://www.gartner.com/en/information-technology/glossary/business-analytics.

[2] M. Dwight, A framework for Applying Analytics in Healthcare – What can be Learned from the Best Practices in Retail, Banking, Politics and Sports, Pearson Education Inc, 2013.

[3] G. Doukidis, N. Mylonopoulos and N. Pouloudi, Social and Economic Transformation in the Digital Era, IGI Global, 2004.

[4] B. Raghynathan and R. V. Maiya, SMACing the Bank - How to Use Social Media, Mobility, Analytics, and Cloud Technologies to Transform the Business Processes of Banks and the banking Experience, CRC Press, 2018.

[5] T. Groenfeldt, "Forbes," 24 Dec 2014. [Online]. Available: https://www.forbes.com/sites/tomgroenfeldt/2014/12/24/american-express-opens-tech-lab-in-palo-alto/.

[6] I. Schmerken, "WallStreet and Technology," 10 Feb 2014. [Online]. Available: http://wallstreetandtech.com/asset-management/fidelity-labs-takes-innovation-to-the-next-level/d/d-id/1316288d41d.html?.

[7] R. Browne, "CNBC," 18 Nov 2019. Available: https://www.cnbc.com/2019/11/18/banks-must-behave-like-tech-companies-to-survive-amid-fintech-threat.html.

[8] J. R. Evans, Business Analytics - Methods, Models and Decisions, Second Edition, Pearson Education, Inc., 2016.

[9] H. Chahal, J. Jyoti and J. Wirtz, Understanding the Role of Business Analytics, Spinger, 2019.

[10] G. Sunila, Practical Machine Learning, Packt, 2016.

[11] O. Theobald, Machine Learning from Absolute Beginners, Scatterplot Press, 2017.

[12] D. TEAM, "Data Flair," 1 Jan 2019. [Online]. Available: https://data-

flair.training/blogs/advantages-and-disadvantages-of-machine-learning/.

[13] L. Martin, S. Suneel and K. Maddulety, "Machine Learning in Banking Risk Management: A Literature Review," Risks, 2019.

[14] I. A. Moosa, Operational Risk Management, Palgrave MACMILLAN, 2007.

[15] T.-M. Choi, H. K. Chan and X. Yue, "Recent Development in Big Data Analytics for Business Operations and Risk Management," IEEE Transactions on Cybernetics, 2007.

[16] E. I. Altman, M. Giancarlo and F. Varetto, "Corporate distress diagnosis: Comparison using linear discriminant analysis and neural networks (the Italian experience)," Journal of Banking & Finance, 1994.

[17] C. L. Huang, M. C. Chen and C. J. Wang, "Credit scoring with a data mining approach beased on support vector machines," ScienceDirect - Expert Systems with Applications, vol. 33, pp. 847-856, 2007.

[18] T. Bellotti and J. Crook, "Support vector machines for credit scoring and discovery of significant features," ScienceDirect - Expert Systems with Applications, vol. 36, pp. 3302-3308, 2009.

[19] A. E. Khandani, A. J. Kim and A. W. Lo, "Consumer credit-risk models via machine leraning algorithms," Journal of Banking & Finance, vol. 34, pp. 2767-2787, 2010.

[20] A. Keramati and N. Yousefi, "A Proposed Classification of Data Mining Techniques in Credit Scoring," in International Conference on Industrial Engineering and Operations Management, Kuala Lumbur, Malaysia, 2011.

[21] L. Zhou and H. Wang, "Loan Default Prediction on Large Imbalanced Data Using Random Forests," TELKOMNIKA Indonesian Journal of Electrical Engineering, vol. 10, pp. 1519-1525, 2012.

[22] T. Harris, "Quantitative credit risk assessment using support vector nachines: Broad versus Narrow deafult definitions," Elsevier - Expert Systems with Applications, vol. 40, 2013.

[23] H. Wang, Q. Xu and L. Zhou, "Large Unbalanced Credit Scoring Using Lasso-Logistic Regression Ensemble," PloS Ones, 2015.

[24] S. Lessman, B. Baesens, H.-V. Seow and L. C. Thomas, "Benchmarking state-of-the-art classiÞcation algorithms for credit scoring: An update of research," European Journal of Operational Research, 2015.

[25] V.-S. Ha and H.-N. Nguyen, "Credit scoring with a feature selection approach based deep learning," MATEC Web of Conference, vol. 54, 2016.

[26] L. Yu, Z. Yang and L. Tang, "A Novel Multistage Deep Belief Network Based Extreme Learning Machine Ensemble Learning Paradigms for Credit Risk Assessment," Flexible Services and Manufacturing Journal, vol. 28, 2016.

[27] S. Hamori, M. Kawai, T. Kume, Y. Murakami and C. Watanabe, "Ensemble Learning or Deep Learning? Application to Default Risk Analysis," Journal of Risk and Financial Management, 2019.

**Claudia ANTAL-VAIDA** is a graduate of the Faculty of Economic Cybernetics, Statistics and Information at the Bucharest Academy of Economic Studies, bachelor's degree in Cybernetics and master's degree in Business Analysis and Enterprise Performance Control. She is currently a PhD Student at the same University, mostly interested in business analytics, machine learning and data structures.

# Business Intelligence and Machine Learning. Integrated cloud solutions providing business insights for decision makers

Laura - Gabriela TĂNĂSESCU
The Bucharest University of Economic Studies, Romania
lauratanasescu@gmail.com

*The aim of this paper is to present the latest trends in business intelligence and ways in which nowadays organizations can implement cloud technologies. This work is going to present challenges of the market, providers of integrated cloud business intelligences tools, advantages and disadvantages of moving to the cloud. A real life use case will argue the importance of taking advantage of data, as well as the necessity and the obvious benefits of having the right tools of transforming data into correct business decisions.*

***Keywords****: business intelligence, cloud computing, artificial intelligence, analytics, machine learning, innovation, data*

## 1 Introduction

Technology has evolved a lot in the last decade and nowadays we can even talk about a new paradigm related to cloud and how cloud technologies are going to influence organizations and their development.

In the same time, there are plenty of talks about a revolution related to data. How big data has been developing in recent years, how is going to challenge artificial intelligence our everyday work and what is the way in which organizations adopt business intelligence in order to gain insights from this data.

Therefore, the aim of this paper is to talk about the recent trends in technology, offering clear but relevant information about the most important concepts. In addition, the paper is going to provide an example about business intelligence applications in real life use cases, using cloud technologies from one of the top providers of cloud.

In the following chapters, this work is going to talk about all the theoretical concepts mentioned before in order to provide a clear image of the domain. Afterwards, it is presented an analysis of Oracle Corporation and its analytics solution in cloud, with advantages and disadvantages, competitors and benefits. Finally, a use case is going to be realized with this technology.

## 2 Cloud technology

This chapter is going to present some theoretical aspects related to cloud, as well as types of it and what are the advantages and disadvantages of using it.

### 2.1 Cloud computing

We can refer to cloud computing as the possibility to provision computing services with the help of the internet, services where we can include networks, software, servers, analytics and databases. All these cloud capabilities are used to offer a faster innovation and a more flexible way to use resources. [1]

A concept (**Fig. 1**) that comes with cloud technology is that the locations of the service used, the hardware, all the operating systems and also many more other details remain irrelevant to the final user. [2]

Practically, cloud providers offer services that enable the users to access, store or transmit file or applications on different remote servers as well as the power to access all the data using the internet. This being said, it is not required for any user to be in a specific place in order to gain access to it. [3]

**Fig. 1**. Cloud computing concept [4]
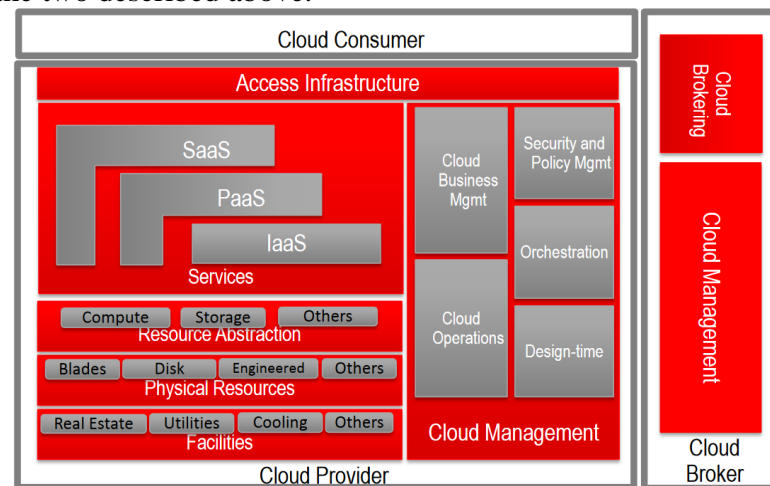
## 2.2 Cloud classification

Cloud computing is known as public or private cloud. The first one refers to those services that are offered to users for free over the internet. The second one provides access just for a number of people, offering services that are a system of networks. Here, we can also mention a third category, known as a hybrid cloud, that is a combination of the two described above.

Cloud computing cannot be seen as one piece of technology, but it is divided in three different services: software-as-a-service (SaaS), infrastructure-as-a-service (IaaS) and platform-as-a-service (PaaS). The first one is related to the part of the license for software applications that is offered to customers, license that is provided with a pay as you go model.

The second one provides customers the opportunity to practically rent infrastructure that includes servers, storage, operating systems and networks from any of the cloud providers of infrastructure.

The last one is especially designed in order to make it easier for developers to create web and mobile apps, without having the need to manage or set an environment and infrastructure for the development process. [5]

In **Fig. 2** can be seen a complete and detailed architecture of cloud.



**Fig. 2.** Cloud architecture

## 2.3 Advantages and disadvantages

Looking at this new technology that impacts our life nowadays, it is important to talk about the benefits that come with it. So, the first one to mention is related to costs. The fact that cloud computing eliminates all the expenses that were coming with the hardware and the software, as well as with setting up and running all the data centers, it is needless

to say that the overall costs are being reduced. The second one is effective for all the customers and also for the technical users and it is about speed. All the services are provided as self-service and on demand, as well as the fact that every service can be provisioned within minutes and without a great amount of knowledge.

The third benefit that deserves being mentioned is related to performance that comes from the fact that cloud services run

on a worldwide network of secure data centers. In addition, the performance comes also from the upgrades that are regularly being made to the systems, making them always faster and more efficient.

Last, but not least, we must talk about security. There are plenty of systems that cannot assure a good security due to the lack of knowledge or the missing budgets for improvements, so all the technologies and policies that are offered from the cloud providers offer a very important and needed secured system for customers.

Having mentioned all the benefits from moving to cloud and accept the innovations, it is equally fair to also talk about the downsides. The most important one that can be identified is again related to security. Moving and working with sensible data to a cloud that runs on a different country, for example, can cause concerns.

There are also several regulations that are unclear when talking about whether or not some critical national data can be stored in another country where the data center physically is. So, this is a risk that an organization should consider when trying to adopt cloud technologies.

Furthermore, the fact that just one portal is used by multiple employers at the same time, manipulating data and making changes too, can cause damage to the overall course of work.

## 3   Business Intelligence

This chapter is going to present some aspects of business intelligence as known today, as well as how this technology is used in cloud.

### 3.1 BI concept in nowadays technology

Business intelligence can be defined as a software application that is realized in order to analyze, report and offer visualizations of data. The entire procedure that includes reporting data, analyzing it and also accessing all the sources are achieved by a business intelligence software. This concept covers multiples directions like applications, technologies, processes and tools, as well as practices of translating relevant conclusions. [1]
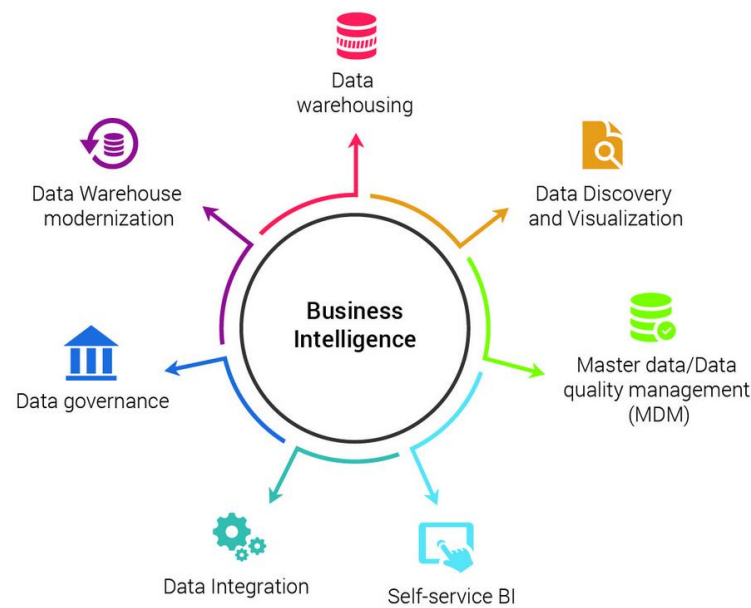
As mentioned above, business intelligence is, in fact, a process driven by and with data, where data storage and knowledge management make a combination that helps in the business decision process. [1]

BI technologies (**Fig. 3**) are used in order to help organizations achieve better decisions about the existent processes, requiring skills, relevant data and innovative technology. BI can be extended as a concept where it can include not only applications and tools, but also infrastructure and practices that enables those organizations to analyze information faster and better, optimize processes and formulate relevant conclusions and taking decisions. [1]

We can admit that a successful BI implementation should be focused on software development or hardware, but on the value that comes from information.

Taking this into considerations, it is important to understand the way data is created and used, what is the quality of that data, how it is constructed that system and the service levels.

So, coming to a point where a conclusion of business intelligence concept is needed for nowadays technology, we can best define it as a set of business data that is taken from multiple sources which are translated into information using different applications in order to support decision making and help organizations to achieve their needs. [6]

**Fig. 3.** Business Intelligence concept

### 3.2 BI in cloud

When using business intelligence solution in cloud computing environment, we should underline the great opportunities this combination can offer. Even though both of these technologies are at a starting point of their development, they are in trends for most of the organizations, having some difficulties to solve though. One of the main problems with these two is related to integration. Of course it is relevant to add here the costs that come from reorganization of processes and work, as well as from people trainings. Not only will these costs appear, but it is also possible that many employees will have a bad attitude towards change. Here it can be added the lack of resources to support these changes, the possibility of downsizing the targeted departments and also the uncertainty that comes with adopting new technology. Nevertheless, every risks and disadvantage that were mentioned about any cloud technologies can be translated to this combination of applications too.

### 3.3 Methods used to create a business intelligence system in cloud

In order to adopt such a system, there are several steps to follow so that the final result should be the one expected.

The first thing to consider is about data collections, where it is important to have the means of accessing and integrating all the places from where the data can be taken. In addition, an architectural model has to be proposed in order to have the best way to collect data.

Another step in this process is the validation part that also comes with reliability. So it is proposed to solve the reliability issues to ensure correct measurement accuracy and also the right measuring instrument used in the measuring process. Reliability is in fact used to phrase the measurement to which a metric provides correct results and no random errors.

The third part is about data preparation. In order to prepare data, we have to first collect all the data, combine it from all the sources where the data was found, structure it in order to be clear and organize it so that it can be easily analyzed. Analysis of the data comes with a process where statistical or logical techniques can be applied in order to illustrate, recap and evaluate data.

The last part and the one that also brings business value and insights is about data

28

*Business Intelligence and Machine Learning. Integrated
cloud solutions providing business insights for decision makers*

analytics. The most common way used was descriptive analysis, where reports were added. After big data has started entering in the biggest companies on the market, the traditional business intelligence has changed due to speed and ways of storing. Therefore, predictive and normative analysis has emerged lately, the firs ones being in the spotlight as well.

The evolution of big data and analytics has affected the overall way of business intelligence delivery. Information needs to be quickly extracted from data, organizations being more and more concerned about normative and predictive analytics that include machine learning capabilities and rapid ways of building relevant visualizations.
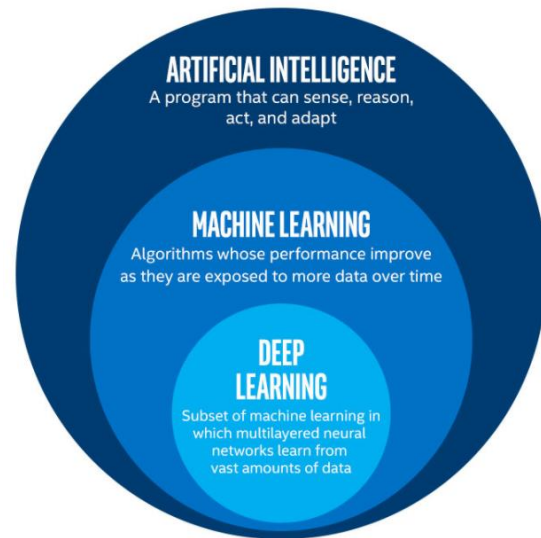
## 4 Artificial intelligence

In this chapter, the paper presents some theoretical information, as well as a brief introduction in artificial intelligence. Moreover, it is going to be made the obvious relationship between data and artificial intelligence in nowadays businesses, as well as the important part named machine learning.

### 4.1 AI definition and structure

While it is mentioned before that business intelligence works with the aim of collecting, reporting and analyzing data, artificial intelligence comes with another approach that impacts data.

In fact, artificial intelligence enables computers to make their own decisions. Thus, we can define artificial intelligence as the ability of a machine or computer to learn and think like human's brain.

Artificial intelligence contains subfields like machine learning, neural network, deep learning, compute vision and natural language processing (**Fig.4**). Explicitly, machine learning is working to automate analytical model building. This field uses different methods like neural network, operations research and statistics so that to find hidden insights from data. [7] [8]



**Fig. 4.** Artificial intelligence and its categories

### 4.2 Data and AI

There are multiple sectors of economy that deal with huge amounts of data which are available in different formats and sources. This enormous amount known as big data is becoming available and easily accessible due to the progress of technology. Multiples data applications of machine learning are formed through complex algorithms build into a machine or computer. The code used creates a model that identifies the data and, after data, it is building predictions around it. The model is going to use parameters built in the algorithm in order to form patterns that are going to help the decision making process. When new data is added to the process, the algorithm used will adjust those parameters mentioned before in order to check if the patter has changed. However, the entire model should remain the same.

AI along with Machine Learning and Deep Learning present multiple technologies that are utilizing Tensor Process Unit (TPU) and Graphics Processing Unit (GPU).

### 4.3 AI applications in cloud

Apart from the visualizations of data that are done using a business intelligence software, we can talk separately about

what is the value that cloud brings to the machine learning component.

Therefore, there are many reasons to talk about regarding using machine learning in cloud, along with business intelligence. First of all, it is about the leverage and speed provided by the power of the GPUs that are needed to train different algorithms, without investing a lot in hardware. Moreover, the scale up and down capability make it efficient and easier for users to improve the power depending on the needs and measure a project have.

In addition, the new picture offered by cloud providers in terms of business intelligence and machine learning does not require advanced skills and lots of knowledge in data science and programming.

## 5 Oracle as a cloud provider of analytics platform

In the following lines, the paper introduces Oracle as a cloud provider, as well as an interesting player in data and analytics market for cloud users. Finally, this chapter also propose a demo of Oracle analytics platform that is going to demonstrate the benefits of using cloud for analyzing data.

### 5.1 Oracle Analytics Cloud

One of the top cloud providers that also comes with an analytics platform is Oracle, which is proposing a comprehensive tool in a unified platform, including data preparation for enterprise reporting, self-service visualizations, advanced analytics, self-learning analytics and machine learning integration on top. (**Fig.5**)
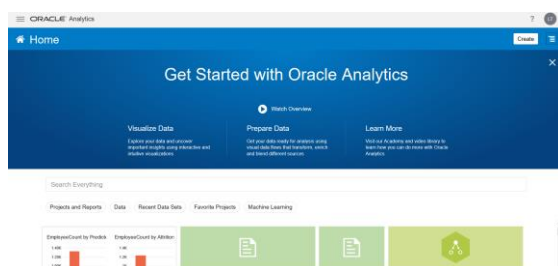
**Fig. 5.** Oracle Analytics Cloud interface

In the capabilities of this cloud platform we can enumerate data discovery, which helps users to easily collaborate with others, building intelligent analysis, machine learning models and statistical modeling.

Another thing to mention is related to the fact that developers can utilize interfaces that help them extend and customize all the analytics experiences in the flow.

It is very interesting the fact that in Oracle Analytics Cloud users can take data from any source, collaborate on project with others and explore real time data. Furthermore, unlike other providers that require the user to compromise between self-service, governed and centralized analytics, Oracle Analytics Cloud (OAC) solves this problems by offering a single solution that also incorporates Machine Learning and Artificial Intelligence.

Through the capabilities of OAC we find the data preparation enrichment that is built into the analytics platform. Another one is the business scenario modeling, a self-service engine for industry that helps in multidimensional and visual analyses. Moreover, we see here that proactive mobile that always learns from your work and offers contextual insights in daily activities. Last, but not least, is the enterprise reporting capability, the power of security and governance having a semantic layer which maps complex data into familiar business terms.

### 5.2 Augmented analytics – features of OAC

Keeping in mind the concepts mentioned above, we can converge business intelligence, artificial intelligence and more specific machine learning, into a term named augmented analytics.

We can see this concept as an evolution for the foundation build from analytics and business intelligence as well as big data, combining different and emerging technologies.

While business intelligence is about creating and finding data insights, AI and ML are about learning from different datasets in order to offer machine-driven decisions.

As it is known at the moment, a BI platform actually ingest a lot of data from multiples sources before anyone can prepare and reorder data.

An augmented analytics system is taking these latter steps and automates them using machine learning and artificial intelligence technologies. As an explanation, machine learning handles data preparation and artificial intelligence handles initial analysis.

Looking at the benefits of such a concept, we can tell that, in spite of those that are offered by multiples providers, there are some that offer a level of efficiency and accuracy that is possible due to computer processing. Thus, one of the most important aspect of augmented analytics includes accuracy. If the analysis is made by data scientists, there is likely that a mistake is going to occur. When using machine learning for that, these situations are eliminated from the beginning. [9]

Another thing to consider is speed. There are gaps that can appear when we first initiate a project using a BI platform like when we manually prepared data and also wait some time in order to receive an answer from different parties. Using augmented analytics, this process begins immediately, launching AI to cull the specific and needed data and also to begin the drilling down for the specific output needed for the project.

One more aspect to consider is the reduction of bias. Bias does not have to come as a personal shortcoming, but as a habit or a routine. Humans tend to revert to patterns so there can be a blind spot for data scientist that can lead to overlooked insights. In this case, computers and machines are going to work more efficient without inherent bias.

Last thing on this list is about the resources used. Augmented analytics can increase the resources by having them do more important things than some manual labor. So, for data scientist, it is going to mean more time to create different business problems and extract deeper insights form data. [9]

## 5.3 Oracle versus competitors

One of the advantages that Oracle Analytics Cloud has, as CEO of Red Pill Analytics said [8], is the fact that Oracle offers all in a single solutions. In fact, there is known that other providers offer multiple products in order to satisfy the same need and the problems is this process takes more time, resources and configurations before getting value from the investment.

Another thing to consider is the ability to scale up or down in order to adjust the resources, depending of the nature of workloads.

Also, its ability to offer not only visualizations, but a comprehensive view that helps the enterprise is considered important by another group director of Qubix International [8].

## 6   HR Attrition case study using OAC

Having all these concepts about business intelligence, artificial intelligence and cloud technology in mind, a small demo can be easily provided. An HR data set added in Oracle analytics cloud is going to be used so that to present the advantages and extended possibilities for data analysis.

## 6.3 Data   loading   and   hypothesis    formulation

We are going to use a public data set about employees and some details about them, as well as staff attrition.

These data set contains details about age, department, hours worked, over time, distance from home, daily rate, education, employee satisfaction, gender, job level, job role, marital status, relationship satisfaction and years since promotion.

These variables are considered suitable in order to make an analysis for the

organization's employees (in order to see their satisfaction and problems based on work life balance, benefits and capabilities) as well as realizing an algorithm in order to predict whether or not an employee that we do not know anything about is going to leave the company or not.
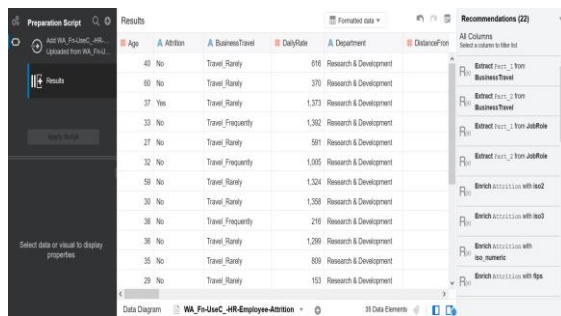


**Fig. 6.** Data loading menu for OAC

### 6.4 Data preparation

As mentioned before, Oracle Analytics Cloud, the tool used for the analysis, is offering intelligent recommendations and possibilities to arrange and filter data, as well as change a measure into an attribute and vice versa.
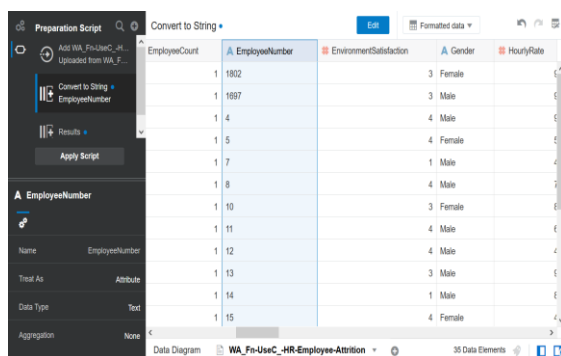


**Fig. 7.** Data preparation menu

There is also a part provided by some machine learning capabilities of the tool that is helping the users to enrich the actual data. Depending on this, we can add this recommendations or not.



**Fig. 8.** Data enrichment capability

In the same time, the explanation mode that can be seen below is part of the augmented analytics. Thus, for a numeric variable like monthly income, we see that the tool offers us different graphics that are relevant for the analysis.



**Fig. 9.** Explanation mode for monthly income – basic menu

Finally, the same thing can be seen in key driver tab that shows us which are the variables that best explain monthly income (**Fig. 10**).



**Fig. 10.** Explanation mode for monthly income – key driver menu

### 6.5 Data visualization

One type of dashboard we can make is a general one that provides us with general information about that data we use.

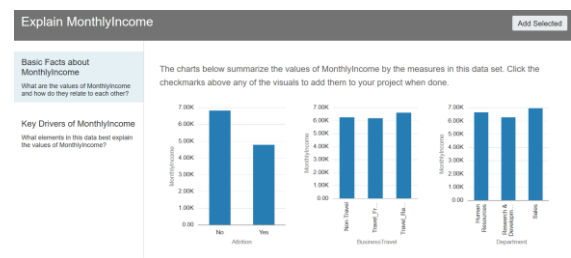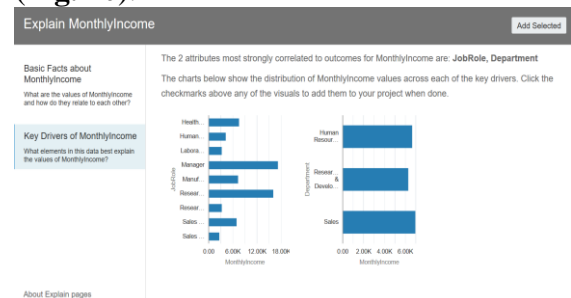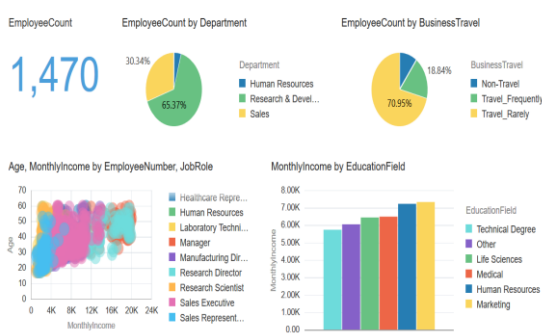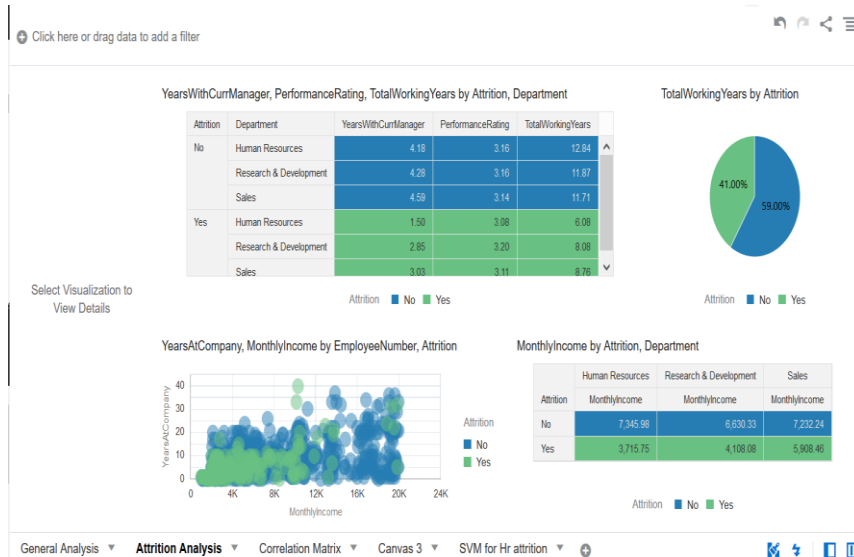**Fig. 11.** Data visualization for general overview

addition, we see that 70% of the people travel rarely, while 19% travel frequently.

In the same time, if we look at the monthly income of the people, as well as the role they have and the age, we see that the greatest income is for research scientists and research directors, with age between 35 and 60. On the other hand, we have the minimal income values for sales representatives and laboratory technicians, where the start age begins with 20.

Last, but no least, we see that marketing provides the greatest monthly income, while a technical degree provides the smallest value when analyzing this organization.

Therefore, we can see in **Fig. 11** that the total number of people analyzed are 1470. The majority are working in research and development and after that sales. In



**Fig. 12.** Data visualization for attrition

The second dashboard (**Fig. 12**) provides information about attrition. We can see from the charts that most of the people that are going to leave the organization are not for many years with the current manager, that they have lower performance ratings and lower total working hours. This shows us that people who will stay within the organization have a history in it, they dedicate a lot of time to it and they perform.

In the same time, most of the people that are going to leave have a smaller income and they have been with the organization for little time than the others.

Nevertheless, we see that the younger employees are those that will leave the organization.

Another interesting thing to consider is the fact that the tool provides an instant visualization of correlation between variables like in the picture below.



**Fig. 13.** Correlation matrix

Not only (**Fig. 13**) is this correlation easy to build because it does not request any technical knowledge, but it instantly provides useful information about our data like the fact that there are powerful and positive connections between age and total working years, monthly income and total working years, age and total working years.

On the other hand, we see some negative connections between years in current role and number of companies worked and number of companies worked and years with current manager.

### 6.4. Machine learning for HR use case

For the machine learning part of this project, we are going to build different machine learning algorithms of classification. After building them, an analysis is useful in order to decide which model is the best one for the use case and that model will be used for prediction.

For the HR attrition use case we are going to use support vector machine algorithm and Naïve Bayes algorithm, both useful for the binary type of classification.

For building this models, we are going to use the data flow that is available in Oracle Analytics Cloud. (**Fig. 14**)

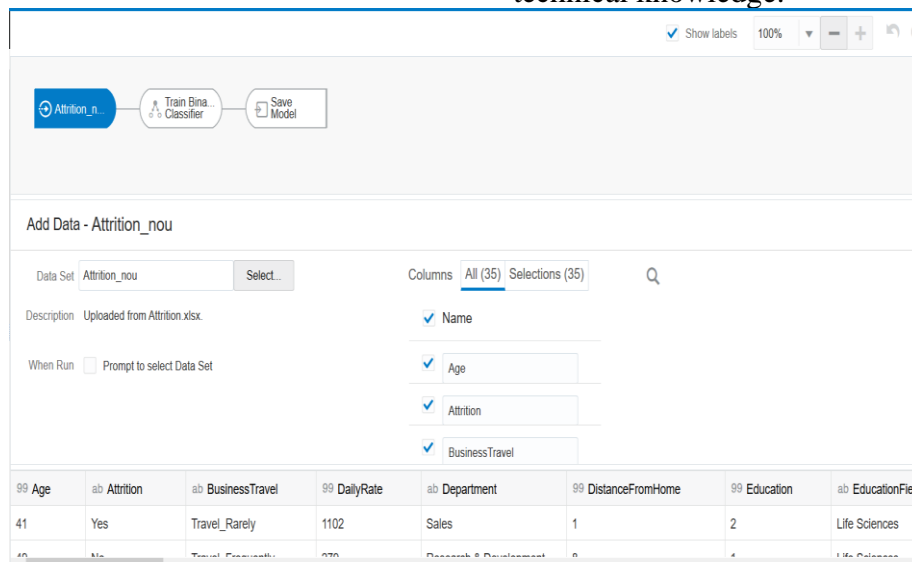As mentioned before, this is step does not require programming or very advanced technical knowledge.



**Fig. 14.** Data flow menu

After applying the chosen algorithm to the data set, we obtain a model that is going to be inspected in order to see how good or bad that model is.



**Fig. 15**. Model analysis for Naïve Bayes

The first one is built for Naïve Bayes and we have an accuracy of 87%, with a

precision of 65%, a recall 38% of and a false positive rate of 4%. (**Fig. 15**)



**Fig. 16.** Model analysis for SVM

For the second one we have an accuracy of 76%, a precision of 37%, a recall 72% and a false positive rate of 23%. (**Fig. 16**)

Therefore, in order to choose the best model, we are looking at the fact that the model is predicting attrition. In other words, we want to maximize the recall because we want to predict correctly all truly positive cases.

Lastly, a false negative for these models points out to the idea that we are wrongly going to conclude that a person is going to

leave the company, fact that might decrease the chances to prevent a person from leaving.

That being said, we can use support vector machine in the next steps so that to predict new values for possible leavers.

Using the same data flow where, to the initial set of data, the support vector machine model is applied to the data, we are going to obtain a prediction that is visualized in the pictures below.



**Fig. 17.** Visualization for SVM model prediction

So, from the total number of people, we can see that most of them were predicted for attrition correctly. 937 have no for

attrition and they were also predicted with no, 177 have yes and their prediction was yes too. The interesting part is that, like

remembered in the previous paragraph, 296 were mistaken, but for our analysis is better to think that 296 are going to leave, even though they are not. Finally, the smallest number is 60 for those that are going to leave and they were actually predicted as non-leavers. (**Fig. 17**)

In the next chart, we see the attrition and prediction for attrition grouped by department. (**Fig. 17**)

We see that most of the leavers are from research and development, as well as from sales. These two departments have also the biggest numbers for the false negatives, but the predictions are overall very good.

## 6.5 Use case results and proposal for improvements

We have seen through this use case the population that we analyzed. So, we have different people working in sales, research and human resources, that have experimented different levels of income, working years, type of managers, number of trainings, levels of job satisfaction, work life balance and more.

For these people, we have seen that those with lower levels of income, people that work for little time in the organization, that are little experience or that have been working for less time with a manager are going to be exposed to attrition.

So, in order to prevent this event, a machine learning algorithm of classification is being used so that to predict the possible employees that can leave. We have chosen the best one, more exactly the support vector machine that offered the minimal false negative rate.

With this algorithm we have predicted those employees that might leave and, the targeted organization can now address to them in order to find solution to the existing problems.

As seen from the data used, some of the solutions to propose might include a solid plan for development that includes levels of income, ways of promotion and trainings. Moreover, organizations should adapt to younger people that tend to leave

early when something is wrong, on contradiction to those that are older and that have spent many years in just one place.

## 7 Conclusions

First thing to mention in this final part of the paper is the fact that, using all the technologies presented before, a business problem was solved within days and with little technical knowledge.

All things being considered, a BI tool like Oracle Analytics provides us an integrated platform that is going to support the work from preparation till predicting future behavior, facts that are going to help business decision makers to act faster and better in their daily work.

Artificial intelligence and machine learning was useful not only for the suggestion area, the explanation mode or enrichment part, but also as providers of useful algorithms that can be applied right away.

Last, but not least, the cloud has offered multiple benefits in the entire work process. First of all, the permissions and roles part that helps more users to work on the same project or use it at the same time. One administrator can be responsible to create all these users and give them the right privileges. [4]

Second of all, the power offered by the cloud in order to run machine learning algorithms is very important. Considering the fact that a support vector machine classification can take a lot of resources, it is clearly an advantage when we can run in just one minute an entire algorithm in order to build a model and a prediction.

Lastly, the platform can be integrated with other solutions, it can take data from other applications, database or personal computer and everything is going to be in one place.

## References

[1] A. Pyae, "Cloud Computing in Business Intelligence," Asia Pacific

University of Technology and Innovation, November 2018.

[2] https://www.zdnet.com/article/what-is-cloud-computing-everything-you-need-to-know-from-public-and-private-cloud-to-software-as-a/, "What is cloud computing?," Steve Ranger, 2018.

[3] https://www.investopedia.com/terms/c/cloud-computing.asp, "Investopedia Cloud Computing".

[4] A. Banafa, "Ten Myths about Cloud Computing," https://www.experfy.com/blog/ten-myths-about-cloud-computing, 2019.

[5] https://azure.microsoft.com/en-gb/overview/what-is-cloud-computing/, "Azure Microsoft - Cloud computing".

[6] https://www.saksoft.com/information-management-services/business-intelligence/, "Information management services".

[7] https://www.sas.com/en_us/insights/articles/big-data/artificial-intelligence-machine-learning-deep-learning-and-beyond.html, "SAS Big Data and Artificial Intelligence".

[8] R. Clayton, "Oracle Analytics Cloud Succeeds Where Others Fall Short," Oracle Blogs, 2019.

[9] https://blogs.oracle.com/analyticscloud/what-is-augmented-analytics-v2, "What is Augmented Analytics?," Oracle Blogs, 2019.

# Testing Approaches for an Electricity Market Simulator

Anca Ioana ANDREESCU, Ana Ramona BOLOGA
The Bucharest University of Economic Studies, Romania
Anca.andreescu@ie.ase.ro, Ramona.bologa@ie.ase.ro

*Software testing methodologies represent different approaches and techniques to ensure that a particular software application is fully tested. This paper addresses the functional testing of a software system for simulating electricity markets. The tested functionalities are briefly presented and the main testing techniques compatible with this type of system are analyzed. The purpose is to recommend the best combination of techniques recommended in this case and to present results obtained by applying them to the simulator.*
*Keywords: functional testing technics, behavioral testing, electricity market, simulator*
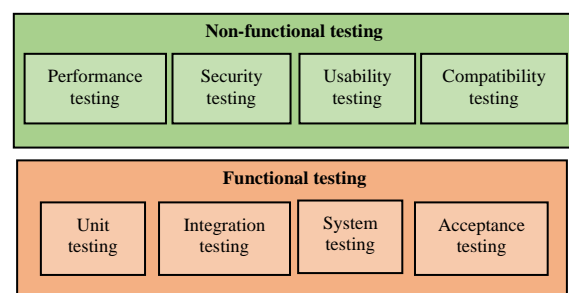
## 1 Introduction

Software testing is one of the main stages of the development cycle of an IT system, vital for software quality assurance. According to IEEE, software testing can be defined as: "the process of exercising or evaluating a system or system component by manual or automatic means to verify that it satisfies the requirements" [1]. Often, however, software testing involves considerable costs and time, hence the continuous concern for reducing costs and streamlining the testing activity.

The test cases can be generated based on the results obtained in the analysis and design stages of the computer system or they can be obtained based on the code developed in the implementation stage.

This work will focus on functional testing and the main techniques recommended to perform a more complete and efficient testing of a computer system developed to simulate participation in electricity markets. The paper is structured in four parts, as follows: section 2 - a brief presentation of software testing approaches for complex systems; section 3 - the main functionalities a wholesale electricity market simulator for which test cases have to be prepared; section 4 – presents examples of various testing techniques that may be applied to the simulator, including their advantages and disadvantages; section 4 is dedicated to results analysis, discussions and conclusions.

## 2 Testing in complex software systems

Software testing methodologies represent the different approaches and ways to ensure that a particular software application is fully tested. Software testing methodologies must include a broad spectrum of elements, from testing individual module units, testing the integration of an entire system, to specialized forms of testing such as security and performance [2]. According to the types of requirements of a computer system, two general levels can be identified at the level of testing: functional testing and non-functional testing, each having a series of components presented in Fig.1.



**Fig. 1.** Types of testing for software systems [2]

Functional testing (black box) is performed using the functional specifications provided by the client or by using the analysis and design specifications, such as use cases.

Non-functional testing (white box) involves testing the system in accordance with non-functional requirements, which usually involves measuring / testing the system

according to the defined technical qualities, for example vulnerability, scalability or usability. Traditionally, software testing considers only static views of the code and does not sufficiently address the dynamic behavior of the computer system. But, during this stage, the testing of an entire system should be performed, based on its specifications, considering the specific dynamics and interdependencies.

If in the first case we are talking about *structural testing*, based on the code resulting in the coding stage, in the second case we are talking about *behavioral testing* based on the requirements and design specifications. The use of the models obtained during the analysis and design stages can lead to a reduction of the costs through reuse and the improvement of the verification and validation. In addition, the test activity can be started in parallel with the writing of the code, as soon as the modeling results are available.

For example, in the context of object-oriented modeling using UML (Unified Modeling Language), use cases and corresponding sequence and collaboration diagrams, class diagrams, can be used as sources of information relevant to the test cases. The analysis of the use case scenarios offers a complete understanding of the system, but, being an informal description, it is difficult to generate automatic test cases. Therefore, an active research direction is the formalization of specifications to make it possible to automatically generate test cases.

Ryser proposed a method of creating test cases starting from use cases, scenarios and use case diagrams called SCENT approach [3]. Use cases and scenarios are transformed into semi-formal representation state charts that are further used for test case generation. But the sequence dependencies between use cases was first approached by Briand and Labiche [4], who used an activity diagram of use cases for each actor in the system. Later, Touseef and Zahid [5] used the same approach for representing sequential dependencies

between use cases, but added some execution contracts for use case scenarios in form of logical expressions.

Use cases and sequence diagrams were used by Swain et al. [6] to generate test cases, taking into consideration the dependency sequences between use cases in form of a graph.

## 3 Main functionalities of the Wholesale Electricity Market Simulator

The system for simulating the participation of a producer / supplier / trader / consumer of electricity in different types of markets aims to identify how decisions regarding the price or quantities offered can influence profit optimization. It is a complex system, with numerous restrictions and operating rules, but also with important interdependencies of the decisions made on the different markets.

In this section we will describe the simulator's functionalities regarding the transactions on the following markets ([7], [8], [9], [10] and [11]): Bilateral Wholesale Markets (BWM), Day Ahead Market (DAM), Intra-Day Market (IDM), Balancing Market (BM) and Ancillary Services Market (ASM). Also, before simulating the market trading, several basic configurations are needed to be made in a General Settings Module.

For the implementation of the simulator, a series of classes have been identified that contain the entities participating in its construction. The identified classes represent abstractions of the objects that interact within the system. A distinction can be made between two classes of classes: basic classes, which include objects characterized by dynamic behaviour and independent classes, which have been used to retrieve data from other external systems. The class diagram in Figure 2 includes the independent classes found at the simulator level.

**Fig. 2.** Class diagram for the identified independent classes

A brief description of the information included in the independent classes is presented below:

• BWM_LE_HISTORY- Average hourly prices for deficit and surplus for the Balancing Market, contract type LE;

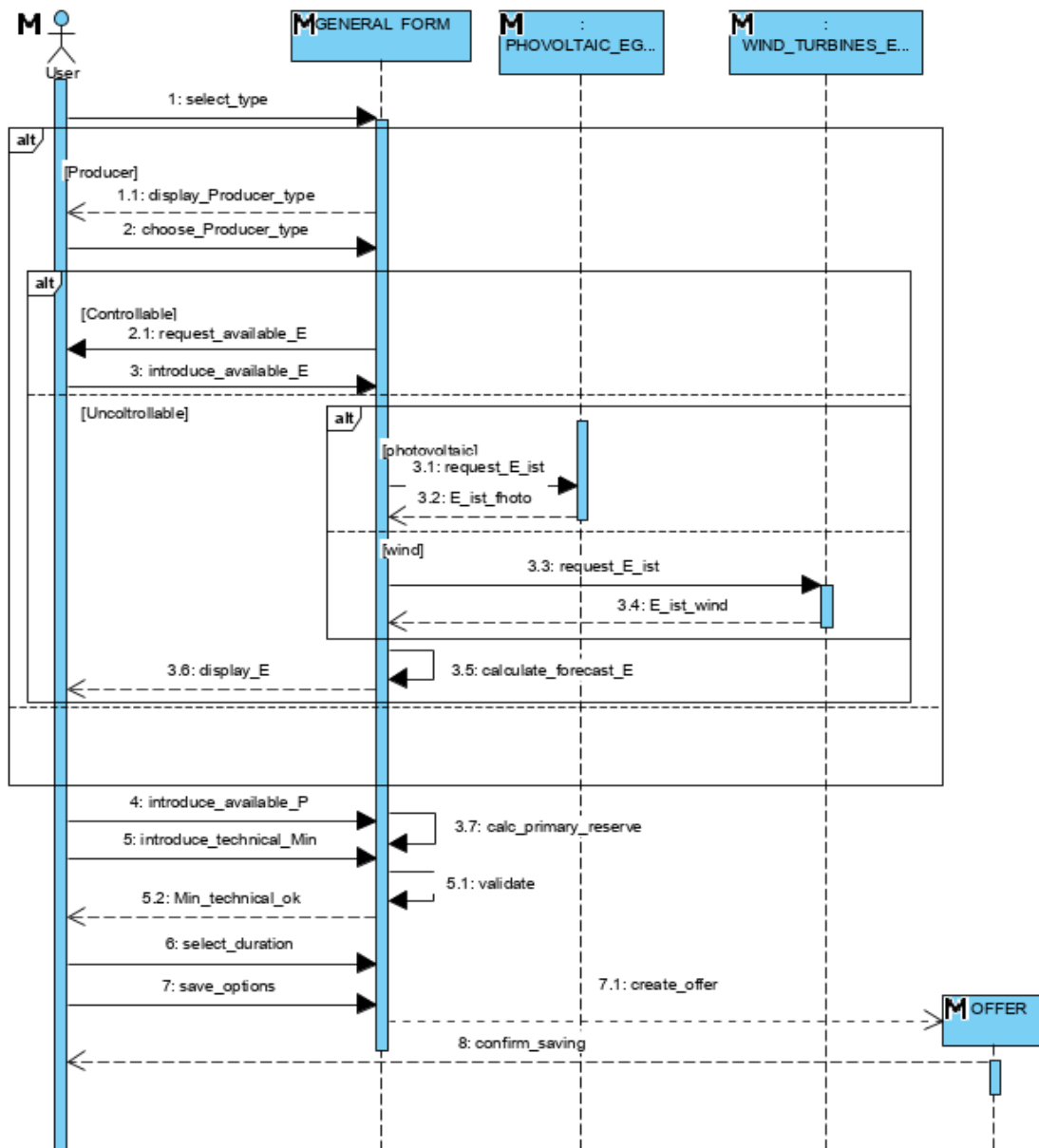• BWM _NC_HISTORY - Average hourly prices for deficit and surplus for the Balancing Market, NC contract type (non-compliant);

• DAM_HISTORY - Historical energy hourly prices for DAM;

• RESERVE_HISTORY - Historical hourly prices for ancillary services;

• BM_HISTORY - Average hourly prices for deficit and surplus for the Balancing Market;

• PHOTOVOLTAIC_EG_PRODUCTION - Data for estimating photovoltaic energy production;

• WIND_TURBINES_EG_PRODUCTION Data for estimating wind energy production.

In order to model the dynamic aspects of the trading simulator, UML interaction diagrams were constructed. These diagrams are made up of a set of objects and the relationships between them, including messages that the objects send from one to the other. There are two types of interaction diagrams: the sequence diagram and the communication diagram. The two diagrams are semantically equivalent and can be transformed from one another. The sequence diagrams illustrate the messages exchanged between the actors / users of the system and the computer system, through the interfaces that it offers. We have started with simple versions of the message sequences, and later we detailed the involved objects.

Figure 3 shows the detailed sequence diagram for configuring the basic options in the General Settings Module. This was built by detailing the messages exchanged between the objects participating in the creation of the scenario. The scenario that underlies this sequence diagram is the following: the user chooses one of the predefined options to select which type of business is simulated: producer, supplier, trader or consumer. If it is a producer, it will choose whether it is type or controllable (thermo, hydro, nuclear) or uncontrollable (photovoltaic, wind). Subsequently, the user will enter the volume of electricity available for trading. The total available power per hour will be introduced, meaning the installed power, if it is a producer or the power required for trading or buying if it is a supplier, consumer or trader. Then, a technological minimum must be entered. The user must choose from a predefined list the periods for the simulation. Then, the month of the year with which the simulation begins will be chosen. The options will be saved, and the offer will be created.

**Fig. 3.** Detailed sequence diagram for General Settings configurations

The sequence diagram in the figure 4 describe the simulation of trading on BWM and it is based on the following scenario: the user selects one of the options for sale or purchase on the BWM. Then, the user chooses from a predefined list of standard products offered by the Romanian Gas and Electricity Market Operator (OPCOM) the one for which the simulation of the transaction is desired, in accordance with the time-related settings in the general module. Once the selection has been made, a validation will be performed and will block the products in the list that have overlaps with the selected product. Subsequently, the user selects information specific to the chosen product, such as semester, quarter, month or year. Will be available for access only the products that have an equal or a shorter period than the chosen period for simulation in the general module. Then, the user enters the desired hourly power to be contracted. This is validated in accordance with the settings in the general module or with other transactions performed on other markets. The cost of imbalances must be introduced, as a percentage of revenues, implicitly with

the value of 10%. This default percentage can be changed. The total hourly power offered on the BWM is calculated and displayed for the chosen option (sale / purchase). Data is saved.



**Fig. 4.** Sequence diagram for trading simulation on BWM

DAM transactions can be simulated using the following scenario: the user selects one of the options for sale or purchase on the DAM and the desired month for simulation. The number of months for which the simulation is performed must be less than the period specified in the General Settings Module and correspond to the periods of a possible simulation performed on the BWM. The available hourly power is displayed by correlating the settings in the general module with any sales

on the BWM and with the primary regulation reserve (for Producers). For information purposes, a historical monthly average price for the chosen month is displayed. Also, a correction coefficient for the selected month will be displayed. It has a default value of 25% and can be modified by the user. Average monthly per hour prices will be displayed for the period chosen by the user. These can be modified. The hourly power that is to be traded must be introduced. For selling, it must be checked that the maximum available power is not exceeded. The data is saved, and the total offer is displayed.

Trading simulation on ASM is based on the process described in the following. The producer selects the desired month for simulation. It is considered that the number of months for which the simulation is performed must be less than the period specified in the General Settings Module and correspond to the periods of a possible simulation performed on BWM or DAM. For this scenario, a historical price is displayed and used. It can be modified by the producer performing the simulation. The available power for reserves is completed and validated based on the total power offered in the general settings module and of the power contracted on BWM and DAM. If desired, the correction coefficient can be updated. Following data is saved. The simulator will calculate and display the total offer on ASM.

The next scenario was considered for trading simulation on BM: the producer selects the desired month for simulation. The number of months for which the simulation is performed must be less than the period specified in the General Settings Module and correspond to the periods of a possible simulation performed on the BWM or DAM. Historical monthly average prices per hour for deficit and surplus will be displayed. Historical prices can be modified by the user. The power offered

by the producer on this market is obtained as the difference between the total power completed in the general settings module and the quantities contracted by participating on other types of markets. Automatically a correction percentage is applied to the estimated revenues or expenses for the simulated period. By default, the risk level is 30% and can be changed. Data for the current month and the offer are saved. Total updated offer is displayed.

## 4 Recommended testing techniques

The main components underlying the testing of computer systems are the techniques, activities, instruments and the controlled environment. Of the mentioned components, we will focus on the techniques used in the testing process. The techniques of designing the test cases can be divided into three categories [12]:

a) Test techniques based on specifications or black box;

b) Testing techniques based on structure or white box;

c) Testing techniques based on experience.

In the following, several types of specification-based techniques (black box) are presented in detail. These tests derive from specifying the desired behavior of the system. It starts from the basic idea that specifications should define what a system should do, not how behavioral specifications should be implemented.

## 4.1 Partitioning into equivalence classes

The documents containing the system requirements normally indicate the rules that the system must follow. There may be rules that imply belonging to a certain range of values or there may be rules of the type if ... then. For example, Rule A might be "if n is less than one, then this is executed." Rule B could be "if n is greater than or equal to one and less than or equal to twelve, execute this".

We consider the rules for the values that a field can store for one month of the year. Thus, we have the following situations:

a) if the month is less than one, the value is not valid;

b) if the month is between one and twelve, it is admissible and the value can be accepted;

c) if the month is greater than twelve, the value is invalid and an error is displayed.

The entire infinite range of integers that could be introduced into the system for the value of the month must fit into one of these three criteria, one of these categories or classes: a) smaller than one; b) between one and twelve; c) greater than twelve. If a value is selected from each number class and we use it as a test case, it can be said that all three rules have been covered. Also, we can say that each value belongs to an equivalence class, hence the name of the technique.

There are three criteria that must be met when creating equivalence classes [13]:

• *Coverage*: each possible input value must belong to one of the equivalence classes.

• *Disjoint character*: the same input value cannot belong to more than one equivalence class.

• *Representation*: If, by using as input value a particular member of an equivalence class, the execution results in an error state, then the same state can be detected using as input value any other member of the class.

Partitioning into equivalence classes applies to all types of values, whether they are continuous or discrete. Using the scenario with the months of the year presented above, the following table can be defined for the partitions of inputs and outputs:

**Table 1**: Partitioning into equivalence classes for inputs and outputs

| Rule | Input partition | Output partition |
|---|---|---|
| Month less than 1, display "Invalid, value too small" | Month <1 | Invalid, too small |
| Month between 1 and 12, display "Valid" | (Month>=1, Month <=12) | Valid |
| Month greater than 12, display "Invalid, value is too big" | Month >12 | Invalid, too big |

### 4.2 Analysis of limit values

The errors tend to be grouped around the limit values. The assumption underlying this type of testing is that developers often omit special cases, represented by the "borders" of equivalence classes.

Using the previous example, we can see that in order to be valid, the month will have to fall within the "valid limit values" from 1 to 12. Also, for the month to be invalid, it must be outside these valid limits. The principle of testing with the limit value analysis is to use the limit value itself and another value as close to it as possible, to be able to reach any part of the limit, using the precision that was applied to the partition. In the example presented, 0,1,2 are valid values for lower limit testing, while 11,12,13 are valid values for upper limit testing.

We can say that the limit values analysis is a particular case of the test method based on equivalence classes, which is centered on exploring the limit cases. Instead of choosing an arbitrary representative of the equivalence class for testing, the method involves choosing a "boundary" element, which represents a particular case. A good way to represent valid and invalid partitions and boundaries is in a table with the following structure:

**Table 2**. Example of partitions and limit values

| Testing conditions | Valid partitions | Invalid partitions | Valid limits | Invalid limits |
|---|---|---|---|---|
| Average hourly price | 10 RON –100 RON | <10 RON >100 RON | 10 RON 100 RON 10,001 RON 99,999 RON | 9,999 RON 100,001 RON |

## 4.3 Testing based on decision tables

Equivalence class partitioning techniques and boundary value analysis are often applied to specific situations or inputs. However, if different combinations of inputs lead to certain actions, this may be more difficult to show using partitioning into equivalence classes and analyzing limit values, techniques that tend to be more focused on the user interface. The other two test techniques based on specifications, decision tables and transitions between states are more focused on business logic or business rules.

A decision table is an effective way of dealing with combinations of elements (for example, input data). This technique is sometimes referred to as the *cause-and-effect table*. Decision tables provide a systematic way of specifying complex business rules, which is useful for both developers and test engineers. Decision tables can be used in designing test cases, whether or not they are described in the specifications, as it helps test engineers explore the effects of different input combinations as well as other system states that need to correctly implement business rules. Decision tables help systematically select efficient test cases and can have the beneficial side effect of finding problems and ambiguities in specifications. It is a technique that works well in combination with partitioning into equivalence classes. The combination of the conditions explored can be a combination of equivalence partitions.

In using the decision tables for designing the tests, the first step is to identify a function or subsystem that has a behavior that reacts according to a combination of inputs or events. After identifying the elements to be combined, they can be placed in a table that shows all the combinations of True and False for each of the conditions. The next step is to identify the correct result for each combination. Each combination of inputs

and outputs is found in the specialized literature as a rule. From a testing perspective, a rule is a test case.

The outputs can be represented either in the form of an expected result (a single output line), or by listing each type of action that can be considered a result and specifying whether or not it will be performed according to the values of the input conditions. We illustrate below the two variants for specifying the results for the decision tables. For the first case, we create test cases for the purpose of verifying a login form, having the structure in the figure below.



**Fig. 5**. Login form

The constrained decision table for testing the functionality of the login form is presented in Table 4. The entry conditions verify that the user and password data have been entered.

**Table 3**. Decision table for testing the login form

| Input | CT1 | CT2 | CT3 | CT4 |
|---|---|---|---|---|
| **User** | Yes | Yes | No | No |
| **Pass-word** | Yes | No | Yes | No |
| **Login** | Yes | Yes | Nu | No |
| **Expected output** | System respon | Error message | Error message | Error message |

The following decision table was constrained to test the functionality of the electricity market simulator, having as input conditions the type of market that is simulated trading and the type of user. The expected results model the possibility of selling / increasing or buying / decreasing in one of the markets depending on the type of user.

**Table 4**. Decision table for testing the trading possibilities in the market simulator

| Input | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| **Producer** | Yes | No | No | No | Yes | No |
| **Supplier** | No | Yes | No | No | No | Yes |
| **Trader** | No | No | Yes | No | No | No |
| **Consume** | No | No | No | Yes | No | No |
| **BWM** | Yes | Yes | Yes | Yes | No | No |
| **DAM** | No | No | No | No | Yes | Yes |
| **BM** | No | No | No | No | No | No |
| **ASM** | No | No | No | No | No | No |
| **Selling/ Growth** | Yes | Yes | Yes | No | Yes | Yes |
| **Buying/ Drop** | Yes | Yes | Yes | Yes | Yes | Yes |

## 4.4 Testing based on the transition between states

Testing based on the transition between states is used when a certain aspect of the system can be described in the form of what is called a "finite state machine". This shows that the system can be in a (finite) number of different states, and the transitions from one state to another are determined by the rules of the "machine". This model is based on both the system and the test cases.

Any system where a different output can be obtained for the same input, depending on what happened before, is a finite state system. Such a system can be described in the form of UML diagrams of state machines. A model that highlights the transition between states has four basic parts:
• The states in which the system can be located (for example, closed or open)
• Transitions from one state to another (transactions are allowed only between certain states)
• The events that generate the transitions
• Actions resulting from transactions.

It is noteworthy that, in any state, an event may cause a single action, while the same event, from a different state, may cause another action or the transition to another state.

The test conditions can be derived from the diagram that models the machine with states in different ways. Each state can be considered a testing condition, as is each transition. The following table describes a general model for a state table that can be used to test a system in which finite states can be identified and modeled.

**Table 5**. General model for a table of states

| | | 1st test case | 2nd test case | 3rd test case | 4th test case |
|---|---|---|---|---|---|
| | | 1st transition | 2nd transition | 3rd transition | 4th transition |
| **STATES** | A state | | | | |
| | B state | | | | |
| | C state | | | | |
| | D state | | | | |

One of the advantages of the technique based on the transitions between states is that models can be built at the desired levels of detail and abstraction to verify the different critical or less critical aspects of the system.

## 4.5 Use case-based testing

Use case-based testing is a technique that helps identify test cases for verifying the entire system from a transactional perspective. A use case is a description of a particular way of using the system by an actor. Each use case describes the interactions that the actor has with the system to perform a specific task (or at least a measurable result for the user). In general, actors are people, but there may be other systems. Use cases are a sequence of steps that describe the interactions between the actor and the system.

The use cases are defined in terms of the actor and from his perspective, not of the computer system. I often use business language and terms, rather than technical

terms. They serve as a basis for developing test cases for system testing and acceptance testing. Use cases may reveal integration defects, i.e. defects caused by the incorrect interaction between different components.

Each use case usually has a main scenario, as well as other additional alternatives (covering, for example, special cases or exceptional conditions). Each use case must specify any preconditions that must be met to operate the use case. The use cases must also specify the postconditions that can be observed and a description of the final state of the system, after the use case has been successfully executed.

The system requirements can be specified as a set of use cases. This approach can facilitate the involvement of users in the process of collecting and defining the requirements, but also in the testing process.The following table shows the partial variant of a use case template that documents the simulation of the sale on the BWM market for the manufacturer. The basic and alternative flows are highlighted, with the highlighting of the actors involved. These will serve as input into the process of testing the scenarios involved in this use case.

**Table 6**. Partial description of a use case highlighting the actors involved

| Name | Simulation sale on the BWM market for the manufacturer |
|---|---|
| **Preconditions** | The manufacturer specific activity settings from the general module have been introduced |
| **Postconditions** | The manufacturer has simulated its offer on the market of bilateral contracts for the period specified in the general module. The simulation data on the BWM were saved. |
| **Basic flow** <br> **P: Producer** <br> **S: System** | 1. P: Select the option to sell on the BWM market. <br> 2. P: Choose from a predefined list of products (OPCOM) the one for which you want to simulate the transaction, according to the time settings in the general module. <br> 3. S: Check overlaps with the selected product. <br> 4. P: Selects information specific to the chosen product, such as semester, quarter, month or year. (Only available products that have a shorter or a shorter period with the chosen period for simulation in the general module can be accessed.) <br> 5. P: Enter the desired hourly power to be contracted. <br> 6. S: The hourly power desired to be contracted will be validated in order not to exceed the declarant availability in the general module. <br> 7. S: A historical monthly average price for the selected product will be displayed <br> 8. S: The cost of imbalances is completed, as a percentage of the receipts, implicitly with the value of 10%. <br> 9. S: Calculate and display the total hourly power offered on BWM for sale. <br> 10. S: Saves data for BWM simulation. |
| **Alternative flows** | 3A: There are overlaps with the selected product - S: The products in the list that have overlaps with the selected product will be blocked. |

6A: The hourly power exceeds the declarant available in the general module- S: An error message is displayed and the hourly power is re-entered

7A: The manufacturer wants to change the price - P: Changes the historical price

7B: There is a negotiated hourly price for a product - S: it will automatically fill in and use the negotiated price

8A: Manufacturer wants to change the cost of imbalances - P: Change percentage of imbalances

As an example of creating and using dependency charts, we selected only the participation of the energy producer on two electricity markets: BWM and DAM. The scenarios in which the Producer is the actor are:

(1) Basic configurations in the General module;
(2) Select product for BWM transaction;
(3) Fill in the hourly power to be contracted for the listed products;
(4) Fill in the price of the imbalances as a percentage of the receipts;
(5) Estimate the correction coefficient according to the weather forecast;
(6) Fill in the offer for each time interval of each month.
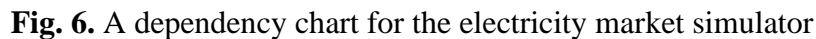
The scenarios in which the System is the actor are:

(10) Import historical data about products;
(11) Retrieve the predefined list of products and details;
(12) Calculate and display Average historical price per hour for the selected product;
(13) Hourly power validation;
(14) Block those products from the list that have overlaps with the selected product;
(15) Estimate the total contracted quantity and update the availability;
(16) Retrieve the historical DAM data;
(17) Calculate monthly average of the historical price;
(18) Validate offer per interval and per month.
(19) Total offer amount update

Dependency charts are very useful in testing, as they support derivation of additional test cases and ensure that dependencies between scenarios are tested. Dependency chart in Figure 6 uses the notations presented in Table 7.

**Table 7.** The notations for dependency charts [3].

| Symbol | Explanation |
|---|---|
|  | Unrestricted scenario |
|  | General dependency |
|  | Sequence |
|  | Alternative |
|  | Iteration |
|  | Real time dependencies |
|  | Structuring constructs |

**Fig. 6.** A dependency chart for the electricity market simulator

**Table 8.** Examples of test case identification

| Test preparation | | Producer has entered basic configurations in the General Module | |
|---|---|---|---|
| **ID** | **Scenarios** | **Expected result** | **Description** |
| 1.1 | (11), (2), (14) | User can not be able to add the selected product to his offer | User select products on BWM that overlaps the previous selected product(s) |
| 1.2 | (11), (2), (14), (12), (3), (13) | An error message is displayed and the hourly power has to be re-entered | Hourly power exceeds the declarant available in the general module |
| 1.3 | (16), (17), (5), (6), (18) | An error message is displayed and per interval and per has to be re-entered | Offer per interval and per moth on DAM exceeds the declarant available in the general module |

Use cases present their scenarios in the form of all possible paths for the specific functionality of the computer system. Therefore, all these scenarios must be verified for the actual implementation of the testing process. By crossing the dependency chart, the necessary test paths can be identified and the use cases efficiently generated. Table 8 presents some examples of test cases based on scenario dependencies, where scenarios are indicated using their identification numbers between brackets.

**5 Conclusions and future work**

As our first development phases used UML for modeling, we find best fitted UML-Based Approach to System Testing, with its obvious advantages.

Each of the functional testing techniques presented above are frequently applied, each

with its own recommendations. Often these are used in a complementary manner to benefit from their specific advantages. In the case of a complex system, such as the system for simulating the participation in the electricity markets presented in this paper, where there are numerous and interdependent restrictions between the different functionalities, the test based on use cases will be the most appropriate technique.

## 6 Acknowledgment

**References:**
[1] N. Kosindrdecha, and J. Daengdej. *A test case generation process and technique* Journal of Software Engineering, 2010;
[2] Inflecta, *Software Testing Methodologies - Learn The Methods & Tools*, March 2018, https://www.inflectra.com/ideas/topic/testing-methodologies.aspx, last accessed on 13.12.2019;
[3] J. Ryser, and M. Glinz, *A scenario-based approach to validating and testing software systems using statecharts.,* Proc. 12th International Conference on Software and Systems Engineering and their Applications, 1999;

[4] L. Briand and Y. Labiche. *A UML-based approach to system testing*, Software Quality Engineering Laboratory, Systems and Computer Engineering, Innovations in Systems and Software Engineering, Springer. pp. 12-24, 2002;
[5] M. Touseef and Z. H. Qaisar. *A use case driven approach for system level testing,* International Journal of Computer Science Issue. Vol. 9, 2012;
[6] S.K. Swain, D. P. Mohapatra, and R. Mall, *Test case generation based on use case and sequence diagram,* International Journal of Software Engineering", 2010;
[7] Guideline on Electricity Balancing. Retrieved January 27, 2020, from https://www.entsoe.eu/network_codes/eb/ ;
[8] Guideline on Electricity Transmission System Operation. Retrieved January 27, 2020, from https://www.entsoe.eu/network_codes/sys-ops/ ;
[9] Regulation regarding intra-day market operation of Romanian Energy Regulatory Authority (RERA). Retrieved January 27, 2020, from https://www.anre.ro/ro/legislatie/documente-de-discutie-ee1/proceduri-oper-regl-comerciale/regulamentul-de-organizare-si-functionare-a-pietei-intrazilnice-de-energie-electrica1387366406 ;
[10] Regulation regarding day ahead market operation of Romanian Energy Regulatory Authority (RERA). Retrieved January 27, 2020, from https://www.anre.ro/ro/energie-electrica/legislatie/documente-de-discutie-ee/proceduri-oper-regl-comerciale/regulament-de-organizare-si-functionare-a-pietei-pentru-ziua-urmatoare-de-energie-electrica-cu-respectarea-mecanismului-de-cuplare-prin-pret-a-pietelor&page=1 ;
[11] Regulation regarding negocitated contracts of Romanian Energy Regulatory Authority (RERA). Retrieved January 27, 2020, from

https://www.anre.ro/ro/legislatie/documente-de-discutie-ee1/proceduri-oper-regl-comerciale/regulament-privind-modalitatile-de-incheiere-a-contractelor-bilaterale-de-energie-electrica-prin-licitatie-extinsa-si-negociere-continua-si-prin-contracte-de-procesare ;

[12] C. Damodar, *Manual Testing Help*, 2012, Retrieved January 27, 2020, from

https://www.softwaretestinghelp.com/manual-testing-help-ebook-free-download/comment-page-1/;

[13] D. Graham, E. van Veenendaal, I. Evans, R. Black, *Foundations of Software Testing: ISTQB Certification 1st Edition*, Cengage Learning Business Press, 2006.

**Anca Ioana Andreescu** graduated from the Faculty of Cybernetics, Statistics and Economic Informatics of the Academy of Economic Studies in 2001. She got the title of doctor in economy in the specialty economic informatics in 2009. At present she is an associate professor in the Department of Economic Informatics and Cybernetics of the Bucharest University of Economic Studies. Her domains of work are: informatics systems and business analytics programming languages.

**Bologa Ana Ramona** graduated from the Faculty of Cybernetics, Statistics and Economic Informatics of the Academy of Economic Studies in 1999. She got the title of doctor in economy in the specialty economic informatics in 2007. At present she is a professor in the Department of Economic Informatics and Cybernetics of the Bucharest University of Economic Studies. Her domains of work are: informatics

# A Big Data Modeling Methodology for NoSQL Document Databases

Gerardo ROSSEL, Andrea MANNA
Universidad de Buenos Aires
Facultad de Ciencias Exactas y Naturales
Departamento de Computación. Buenos Aires, Argentina
grossel@dc.uba.ar, amanna@dc.uba.ar

*In recent years, there has been an increasing interest in the field of non-relational databases. However, far too little attention has been paid to design methodology. Key-value datastores are an important component of a class of non-relational technologies that are grouped under the name of NoSQL databases. The aim of this paper is to propose a design methodology for this type of database that allows overcoming the limitations of the traditional techniques. The proposed methodology leads to a clean design that also allows for better data management and consistency*

**Keywords:** *NoSQL, Document Databases, Conceptual Modeling, Data Modeling, NoSQL Database developing.*

## 1 Introduction

The need for analysis, processing, and storage of large amounts of data has led to what is now called Big Data. The rise of Big Data has had strong impact on data storage technology. The challenges in this regard include: the need to scale horizontally, have access to different data sources, data with no scheme or structure, etc. These demands, coupled with the need for global reach and permanent availability, gave ground to a family of databases, with no reference in the relational model, known as NoSQL or "Not Only SQL".

The NoSQL databases can be classified by the way they store and retrieve the information [1][2]:

- Key-Value databases.
- Document databases.
- Column Families databases.
- Graph Databases.

The development of conceptual modeling and general design methodology associated with the construction of NoSQL databases is at an early stage [SS17]. of data modeling is to highlight in [3]: "*Data modelling has an impact on querying performance, consistency, usability, software debugging and maintainability, and many other aspects*"

There are previous works on development methodologies we can cite, like the BigData Apache Cassandra methodology, proposed by Artem Chebotko [4][13]. It uses the Entity Relationship Diagram as a conceptual model, but it is oriented to a specific engine, Apache Cassandra. Thus, it is not generic and does not adapt to a design of other NoSQL Databases. Another proposal using a conceptual model for the design of NoSQL is described in [5]. It suggests the use of the various NoSQL databases common features to obtain a general methodology, in which an abstract data model called NOAM is used for conceptual data modeling. Such data model is intended to serve all types of NoSQL databases using a general notation.

Recently, an attempt to generate a universal modeling methodology adapted to both relational and non-relational database management systems was also presented, on the grounds of overcoming the constraints that the entity relationship model has, according to the author [6].

The use of conceptual modeling is also proposed in [7], although the background is not sufficiently studied, such as our work on interrelation of documents and the relationship between them and the conceptual model [8]. They use UML as a tool for the realization of the conceptual model and simple rules to transform it into a

logical model using UML stereotypes.

These efforts show that traditional methodologies and techniques of data modeling are insufficient for new generations of non-relational databases.

It is therefore necessary to develop modeling techniques that adapt to these new ways of storing information. In this sense, this paper will provide the tools to solve these limitations for document database design. As indicated in [14] the methodology should allow: "*describe the data-model precisely*"

The rest of the paper is organized as follows: Section 2 outlines the definition of document database; Section 3 describes the main elements of the methodology and phases of document database develop; Section 4 presents the logical design using the document interaction diagram or DID by extending our previous work: moving from logical to physical model using *JsonSCHEMA* is presented in section 5 and finally Section 6 presents conclusions and future work.

## 2 Document Databases

The proposed methodology is oriented to the design of databases based on documents. A document is a collection of field name and value pairs. The values can be a simple atomic value or a complex structure such as lists of values, another document or lists of child documents.

NoSQL documents are generally referred to as schema-less, which seems to suggest that it is not necessary to make a model before the development starts. The fact that the structure of the data does not need to be defined in advance has many advantages for prototyping or exploratory development, but as data expands and the applications make use of them, the necessity to have them organized in some way arises. In that sense it is more appropriate to say that they are agnostic with respect to the internal structure of the data. It is, therefore, necessary to make a design of the data organization.

to as *schema-less*, which seems to suggest that it is not necessary to make a model before the development starts. The fact that the structure of the data does not need to be defined has a priori many advantages for prototyping or exploratory development, but as data expands and the applications make use of them, the necessity to have them organized in some way arises. In that sense it is more appropriate to say that they are *agnostic* with respect to the internal structure of the data. It is, therefore, necessary to make a design of the data organization.
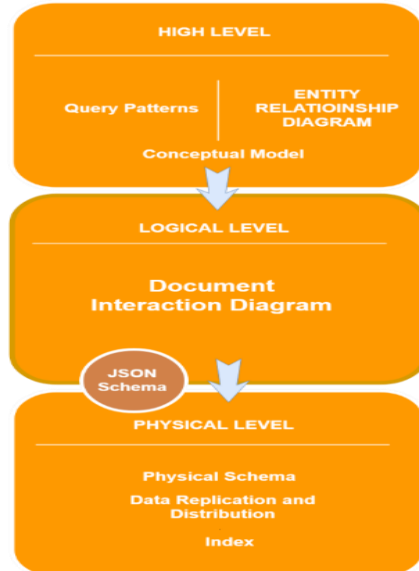
## 3 Methodology

The proposed design methodology has as its starting point the conceptual model, that can be considered as a high-level description of data requirements. Conceptual modeling is usually performed using some form of entity-relationship diagram ([9]) for conceptual class diagram in UML. Conceptual modeling is intended to describe the semantics of software applications.

In traditional relational database design methodologies, conceptual modeling gave way to a logical design that was later transformed into a physical design. It operates by transforming models from higher levels of abstraction to a model that maps directly into the structures of the database.

Phases of proposed NoSQL document database develop consists of high or conceptual level (conceptual model and access patterns), logical level (types of documents, interrelations and specifications), and physical design in steps like phases of traditional relational database.

In the high-level phase, a conceptual data model is developed in a similar way to the design of relational databases. In the current era, with the emergence of Big Data, the need for conceptual modelling is even more important than before.

As a tool of specification and communication with the other phases, the entity relationship diagram is used (ERD) [9]. In this phase, it is also necessary to

specify the query patterns that have been obtained in the analysis requirements. Query patterns can be specified in natural language or in a more formal language like ERQL [10].



**Fig. 1.** Phases

The Logical Level is the heart of the proposed methodology, in which the types of documents and their interrelationships are established. To represent the logical design, we use a new type of diagram that extends the ERD and that we call document interrelation diagram (DID)[8]. Each type of document is later specified using *JSONSchema*.

There are two ways of relating documents: referencing or embedding. The ability to embed documents allows the designer to store related data as a simple document.

In this way, what is called impedance

mismatch can be solved (that is, the difference between the structures of data in memory and the way in which they are stored) [2]. The decision whether to embed or reference is a design decision that is guided by query patterns.

The last phase of our methodology is the analysis and optimization of a logical model to produce a physical data model. In this phase, topics such as index creation, sharding, data distribution, and adapting the data types to the software of the database are considered. The utilization of JSONSchemes is essential in this regard.

**4 Logical Design**

The more important task in this phase is the development of the document interrelation diagram. The DID represents the logical model for a document-based database that captures the classes or types of documents, their structure and interrelation. The documents can be grouped into different classes. Each database uses its own terminology as collections in MongoDB or tables in RethinkDB. we use classes or document types as terminology to indicate a group of documents with similar characteristics.

In the DID each entity of the ERD corresponds to a class or document type, unless it is specifically indicated that this entity will have an independent existence as document type.

In order to exemplify, the entity relationship diagram of Fig.2 will be used. This ERD represents, in a simplified way, the conceptual model of a database that stores orders, products and customers.
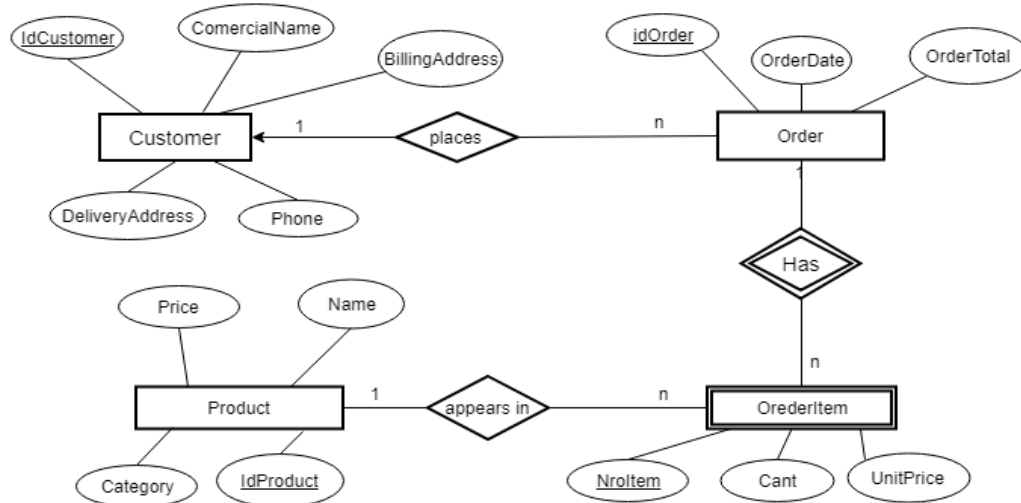
**Fig. 2.** ERD

The entities *Customer*, *Order* and *Product* becomes three document types. *OrderItem* is a weak entity so it is a special case.

To complete the document interaction diagram, it is necessary to decide how the interrelationships will be solved. For this it is necessary to consider the query patterns.

Let's start with the relationship "places". Many design decisions are possible:

• Reference from both sides
• Embed on both sides
• Reference from Order and embed from Customer (or vice versa)
• Embed partially from one side and reference from the other.
• Embed partially from both sides
• Embed total / partial or reference from one side and do nothing from the other

Fig. 3 shows how the reference of both sides is specified while Fig. 4 does the same with embedding of both sides. The arrow indicates reference and curly brackets indicates embedding [8].
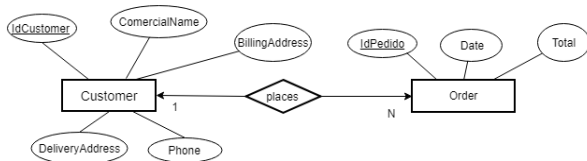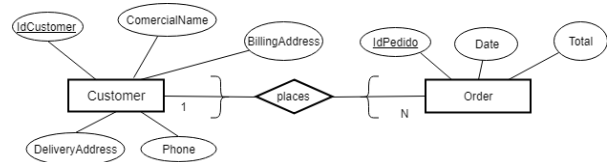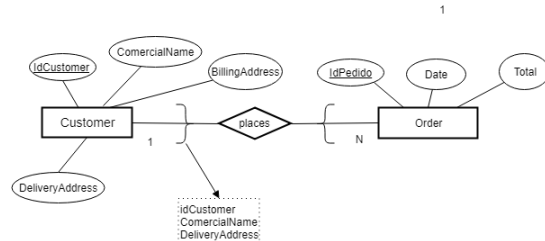


**Fig. 3.** DID: reference



**Fig. 4.** DID: embedding

Embedding simplifies access by minimizing the number of times it should be read from persistent storage. The goal is to keep data that is frequently used together in one document. Although it might be better for a document not to incorporate all the information of the document with which it is interrelated, but only the necessary information that arises from the query patterns.

Suppose that the query patterns indicate that a common way of access to the data is the printing of the order for which the customer's commercial name and shipping address are needed, in addition to all the associated order items. Also suppose that you want to get the dates of the orders made and the total amounts of the same. If the interrelation is solved using only references, the applications are being forced to make several roundtrips to server for to obtain the necessary data. In these cases, a partial embedding can be a better solution.

Fig.5 shows how partially embed is represented. It is necessary to indicate which
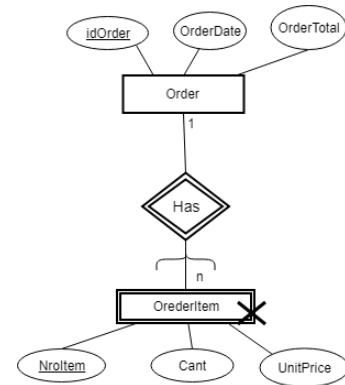
fields of the other entity that will be embedded.



**Fig 5** Embed Partially

Weak entities generally form an aggregate with the strong entity that determines them. It is the case of *ItemOrder* and *Order* in which *Order* can be considered as an aggregate or "a collection of related objects that we wish to treat as a unit" [1].
The simplest way to deal with this is to embed the weak entity in the type of document generated by the strong entity. It is also necessary to indicate that the weak entity will only have an embedded existence, which is done by placing a cross on it as shown in Fig. 6
The cross over any entity indicates that it is not generating a type of document that will be stored independently.



**Fig. 6.** DID: week entity

Although *ItemOrder* entity does not generate a document type, it has an interrelation with the *Product* entity that must be resolved in the logical model. The product information needed in the *ItemOrder* will depend on the domain over which the model was made and what are the access patterns. In this case, it can be assumed that only the name of the product is needed, for which we partially embed the name of the product in the item. When embedding the *ItemOrder* in Order it is embedded with everything it contains including references and embedded fields of other types of documents, in this case the name of the product. The final diagram is as in Fig. 7.



**Fig. 7.** DID

In some cases, it is not enough with the types of documents generated from the ERD to resolve all interrelationships. Assume the case of a database that must save user access to different modules and that a large number of daily accesses are made by each user. The most important query is to know on a given date which modules a user accessed. The ERD in Fig.8 is the conceptual data model.



**Fig. 8**. ERD Users and Modules

How to resolve the interrelation between *User* and *Access*? At first glance it seems to be a case like that of the previous order and item. But there are two important
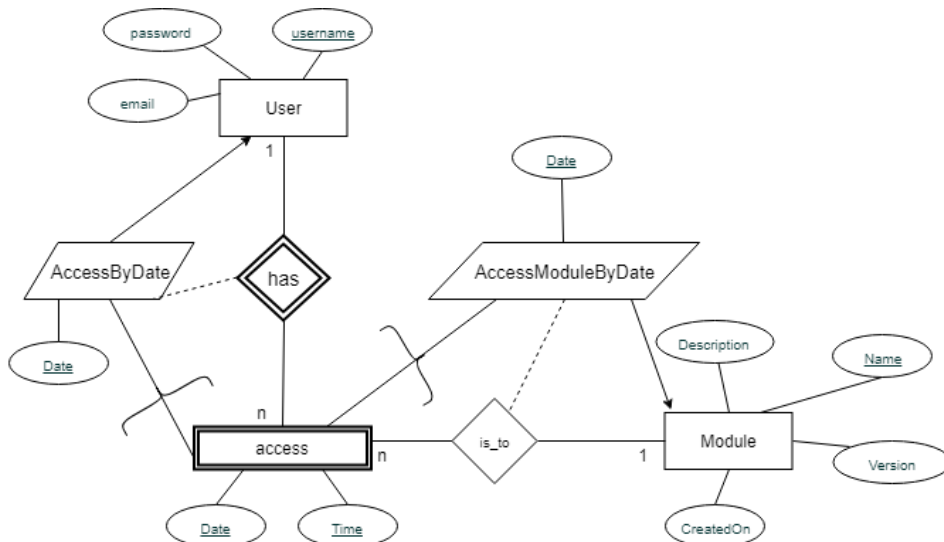
differences that change the design decision:

1. The immutability or not of the data: In the previous case, once the order has been sent to the client, the items can no longer be modified. However, in this new domain accesses are added frequently.
2. The volume of data: The items in an order have a limited amount of data. On the other hand, user accesses grow permanently and frequently.

In a document-based database the document is the unit of access, changes in their sizes may generate the need to reorganize the physical space where they are stored, if this is done very often there may be a degradation of performance.

The query patterns in the example indicate that, in general, accesses for a given date are consulted, so it would be a good design decision to divide the accesses by date. Also, once the date is finished, the accesses of the same are immutable. To have a document by date it is necessary to create an auxiliary document type. Fig. 9 shows how that document is specified.



**Fig. 9.** DID: Partition

The new document that does not correspond to any entity of the ERD is drawn as a parallelogram with two inclined sides. It is also necessary to indicate which interrelation that document is representing that is achieved
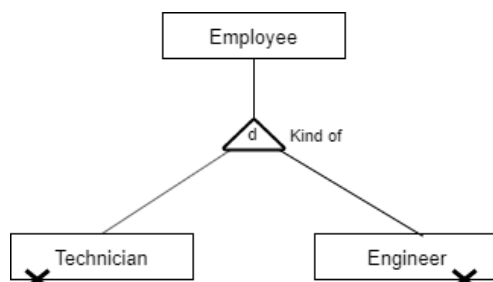
with a dotted line from the interrelation to the symbol of the intermediate document.

The auxiliary document has on one hand a reference to the user and on the other it embeds the accesses. The key will be the date and user id. We must explicitly mark as

a key the *Date* taken from the accesses to indicate that it is the partition key and therefore there is a single date per document, the user identifier does not need to indicate it since the arrow indicates reference to the key of the user and also the cardinality of the user-measurement relationship indicates that the measurements are of a single user. It is not necessary to keep the measurements as an independent document, so the cross is placed on that entity.

The extended entity relationship diagram also supports hierarchies between entities.

The hierarchies in the ERD can be with full or partial coverage, with overlapping or without overlapping. The possibility that documents of the same type have different schemes facilitates the design. We can generate a single type of document corresponding to the super-entity that also has the attributes of the sub-entities. For this, it is enough to indicate that the sub-entities do not generate a type of document as seen in the Fig. 10.
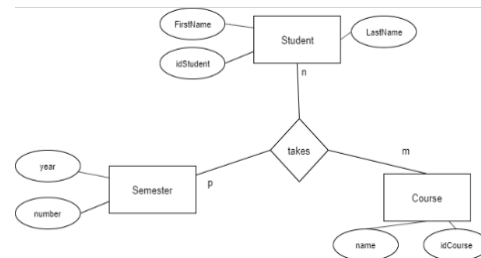


**Fig. 10.** DID Herarchies

Depending on the pattern of consultations, other decisions may be made:

- Mark the super-entity as not generating a document type and then generate one for each sub-entity. This is possible if the hierarchy has no overlap.
- Specify that both super-entity and sub-entities generate one document type each. Indicating which attributes would be placed in super-entity.

Another type of relationship that is necessary to model is ternary relationship.



**Fig. 11.** ERD: ternary relationship

Suppose a ternary relationship between *Student*, *Semester* and *Course* entities. The cardinality in this case is n:m:p, for a student and a semester there are many courses he takes, a semester and a course has many students enrolled, for a course and a student can be many semesters where he takes it. The DER of Fig. 11 shows this relationship.

The most complex part is deciding how to model the relationship takes. The decision on how to model will, as always, depend on the query patterns. The basic case is to generate a type of document that simply contains the information of the relationship with the identifiers of each of the entities involved. To do this, an auxiliary document is drawn with the name of the new document type and a dotted line that binds it to the entity as seen in Fig.12.



**Fig. 12.** DID: ternary relationship

That is the simplest model, but suppose that a very common query is to know which students are enrolled in a course in a semester, in fact you want to know first name, last name of them for a given course and semester. While the previous model

allows you to answer this query, you might decide to have a document type that stores the complete information to optimize access to it.

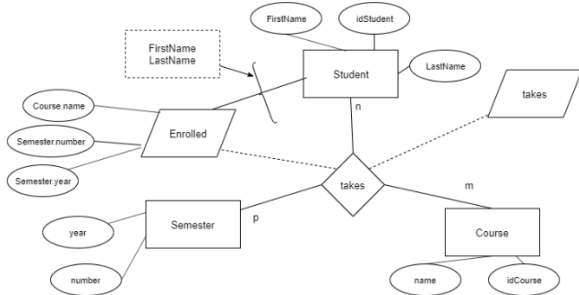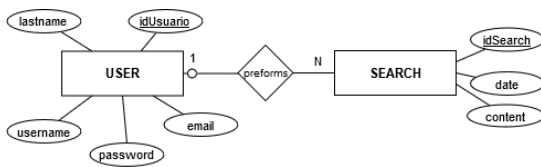The semantics of this diagram (Fig. 13) are that only key attributes are added or that allow you to group data from another entity from those participating entities in the interrelationship that are not related by any link to the new document type.



**Fig. 13.** DID: complex ternary relationship

One case to consider is when it becomes necessary to group multiple instances of an entity, by one or more attributes, into a single document. To exemplify let's assume a part of a DER where users and their searches are modeled.



**Fig. 14.** DER: User/ search

The relationship between user and search can be modeled in various ways, either by embedding searches in the user or by referencing. The relationship could also be resolved by partitioning by user and date in the same way as shown in Figure 9 for user and access. Let's say that a very frequent query is to know the searches performed on a given date. The solution of partitioning by user and date is not efficient for this because access should be made for each user who has a search on that date. In this case, the ideal is to have a single document with all the searches for a date. This would involve grouping by the date attribute, i.e. generating a

document for each date that has all searches. An auxiliary document should be created to save all searches with the date as key. The notation is similar to that seen before, although in this case the auxiliary document refers only to the entity on which it is being grouped.



**Fig. 15**. DID: User/Search

Figure 15 shows the corresponding DID. Note that the reference from *Search* to *User* is important, because marking the entity as not generating a type of document would lose the relationship.

It is also possible to generate an intermediate document to resolve the relationship between *User* and *Search*. There would be data redundancy in favor of access speed. The complete DID is shown in Figure 16.



**Fig. 16.** DID: Complete User-Search

## 5 From logical to physical level

Upon completion of the development model interrelationship of documents, which is equivalent to logic design relational database, it continues with the physical design.

The physical design implies making decisions about specific aspects of implementation such as: data distribution, index generation, use of engine facilities of the selected database, etc.

Many document databases support indexes.

Index creation must be based on query patterns. It's about doing a trade-off so you don't have a few indexes that could lead to poor read performance, but not so many that affect the write performance.

The use of *JSONSchema* for a more detailed specification of each type of document facilitates decision making process and implementation. A *JSONSchema* is a JSON document which describes the structure of another document.

The steps to follow are as follows.

1.  For each document type in the DID:
    a.  Define the appropriate data types for each attribute
    b.  Write the specification using *JSONSchema*.
2.  For each query analyze the ease of documents to respond to it. Ideally a single access should be enough for the most used queries.

From the DID each type of document is mapped to a *JSONSchema* which allows to specify in detail the structure of each document. For example, the document type *AccessByDate* in Fig. 9 is mapped to the the following scheme:

```
{"title":"AccessByDate",
 "type":"object",
 "properties":{
        "userId":{"type":"integer"},
        "date":{"type":"string","format":"date"},
        "accesses":{ "type":"array",
        "items":{ "type":"object",
                    "properties":{
                            "modulename":{"type":"string"},
                            "timestamp":{"type":"string",
                            "format":"date-time"}
                            }
                    }
                }
            }
        }
```

From the DID in Figure 16 JSONSchema will be generated for each of the following document types:

*User*: With the attributes in the diagram, specifying the type of each.

*UserSearchByDate*: having the *userid* and *date* as keys and a vector with that user's searches on that date.

*SearchByDate:* The key is the *date* and has a vector with the searches and in each the corresponding *userid*.

No other document types are generated. By indicating that an attribute is key we are claiming that it is unique and that it identifies each document, even though the database always generates an identifier attribute.

The flexibility of the *JSONSechema* to establish optional properties makes it an ideal tool for specifying document types of variable structure. In the case of hierarchies this facility is extremely useful because you can specify conditions for which an attribute exists or not. Looking at *JSONSchemas* it is possible to realize that in some case it is convenient to reserve space the same in such a way that the document does not resize it during its lifetime. If the document grows larger than the size allocated for it, the document may be moved to another location with the consequent input/output cost [12].

Some document-based databases have tools to validate if a document complies with a *JSONSchema.*

## 6 Conclusions

A methodology that allows obtaining a detailed design from a conceptual model has been presented. This work extends and completes previous work on document modeling in the design process.

The proposal presented allows flexibility to establish detailed design decisions. There is not currently, to the best of our knowledge, complete methodology such as that presented for document-based databases that have the same level of flexibility and specification capability.

The presented methodology was used successfully in several developments using different database engines. In future work we plan to report in detail the cases of success in the use of this methodology.

## References

[1] Adam Flowler, "The State of NoSQL", 1st edition, 2016

[2] Pramod J. Sadalage, Martin Fowler, "NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot Persistence", *Addison-Wesley*

[3] Gómez, P., Casallas, R., & Roncancio, C. (2016). "Data schema does matter,

even in NoSQL systems!" 2016 IEEE Tenth International Conference on Research Challenges in Information Science (RCIS), 1-6.

[4] Artem Chebotko, Andrey Kashlev, Shiyong Lu, "A Big Data Modeling Methodology for Apache Cassandra", IEEE International Congress on Big Data (BigData'15), pp. 238-245, New York, USA, 2015.

[5] Francesca Bugiotti, Luca Cabibbo, Paolo Atzeni, Riccardo Torlone. "Database Design for NoSQL Systems". *International Conference on Conceptual Modeling*, pp. 223 - 231 Atlanta, USA, Oct 2014.

[6] Ted Hills, "NoSQL and SQL Data Modeling", *Basking Ridge, NJ: Technics Publications*, 2016

[7] Shin, K & Hwang, C & Jung, H. (2017). "NoSQL database design using UML conceptual data model based on peter chen's framework". International Journal of Applied Engineering Research. 12. 632-636

[8] Gerardo Rossel, Andrea Manna, "Diseño de Bases de Datos Basadas en Documento: Modelo de Interrelación de Documentos" *XIII Workshop Bases de Datos y Minería*

de Datos. Congreso Argentino de Ciencias de la Computación CACIC 2016 San Luis Argentina..

[9] Peter P. S. Chen, "The entity-relationship model: toward a unified view of data", *Proceedings of the 1st International Conference on Very Large Data Bases,* ACM, New York, NY, USA, 1975.

[10] M. Lawley and R. W. Topor, "A query language for EER schemas," in *Proceedings of the 5th Australasian Database Conference*, 1994, pp.292–304.

[11] Storey, Veda & Song, Il-Yeol. (2017). Big data technologies and management: What conceptual modeling can do. Data & Knowledge Engineering. 10.1016/j.datak.2017.01.001.

[12] Dan Sullivan. 2015. NoSQL for Mere Mortals (1st. ed.). Addison-Wesley Professional

[13] Jeff Carpenter & Eben Hewitt. (2020). Cassandra: The Definitive Guide (3st. ed.). O'Reilly Media, Inc.

[14] Pivert, Olivier. NoSQL Data Models: Trends and Challenges. 2018. Wiley-ISTE

**Gerardo ROSSEL** graduated as Ms. Sc. in Computer Science from the Faculty of Exact and Natural Sciences of the University of Buenos Aires. He has a Doctor's degree from the National University of Tres de Febrero. At present, he is an assistant lecturer at Computer Department of FCEyN UBA. He has more than 20 years of experience in software industry and is Chief Scientist of UpperSoft software company. Her specific area of competences is in Databases, NoSQL, Data Science, Machine Learning, Patterns and Software Architectures, Epistemology and Philosophy of Computer Science. He is co-author of book "*Algoritmos, Objetos y Estructuras de Datos"*. He has published several papers in national and international conferences and journals. He was a member of the International Program Committee of several international conferences.

**Andrea MANNA**, graduated from the Faculty of Exact and Natural Sciences of the University of Buenos Aires in 2000. She got the title of Ms. Sc. of Computer Science. At present, she is assistant lecturer in the Faculty of Exact and Natural Sciences of the University of Buenos Aires. She has been working in the software industry since 1995. She is Chief Software Architect of UpperSoft Sofware Company and work in software development for more than twenty years. She is co-author of book "*Algoritmos, Objetos y Estructuras de Datos".*

# Learning view over the implementation of business process optimizations

Radu SAMOILA
Bucharest University of Economic Studies, Romania
radusamoila2@gmail.com

*Current context of great development and changes in the technological matters, national and global economies have to keep pace and undergone major changes. The ultimate aim of the companies and organization is to improve or optimize its business processes to cope with increased competitiveness in order to deliver more efficiently better products or services. The success of current businesses is linked to the efficiency but also the effectiveness of their core processes. Most part of the latest researches have recognized this importance leading to efforts concentrated on analyzing and optimizing the business processes. There are many techniques which are considered but also others which are not causing potentially significant opportunities of improvement not being addressed. However, currently, there is a scarce of an universal technique or methodology that can be used by the organizations to address business optimizations. This paper addresses different topics on how companies are analyzing the decision to initiate a business process optimization and how optimizations landed within organizations.*

***Keywords:*** *business process, optimization, efficiency and effectiveness of internal processes, business automation, optimization methodology.*

## 1 Introduction

Business processes have received ample attention during the last years. Many approaches and techniques have been discussed and proposed, there were many promises made, but the spectacular results that the reengineering and optimization revolution vowed were never fully realized. This made more and more people hesitant about the whole concepts.

As defined by the main relevant literature publications, Business Process Optimization ("BPO") is the concept of redesigning the internal processes to promote efficiency and effectiveness in order to strengthen the alignment of individual processes with the overall strategy of the company. While the optimization of a singular process or of the processes in a particular business function could trigger real business improvement, organizations that gather their efforts across the entire organization can see significant competitive advantage, better customer service (internal and external), and much more efficient operation.

## 2 Business processes' optimization approaches

Zhou and Chen [1] suggest that business process optimization should aim at reducing lead time and cost, improving quality of product and services, and enhancing the satisfaction of customer and personnel so that the competitive advantage of an organization can be maintained. Reijers [2] suggests that the goals of business process optimization are often the reduction of cost and flow time. However, Hofacker and Vetschera [3] underline that the concept of "optimality" of process designs is not trivial, and the quality of processes is defined by many, often conflicting criteria.

Zhou and Chen [4] remark that there is still no structured optimization methodology or technique for business processes. Optimization is not an option for diagrammatic process models. This is because optimization requires quantitative measures of process performance that cannot be offered by diagrammatic models. However, there are many qualitative improvement approaches applied to diagrammatic process models such as that by Zakarian [5] and Phalp and Shepperd [6]

to name a few. **Table 1** summarizes the main business process optimization approaches identified in literature, mostly related to Petri nets and mathematical process models. Taking into consideration the emphasis that has been put on Petri nets for their analysis capabilities, one would expect that they would also fit for optimization purposes.

But, according to Lee [40], Petri nets are not appropriate to solve optimization problems except using graph reduction techniques. Although they can capture system dynamics and physical constraints, they are not suitable for optimization problems with combinatorial characteristics and complex precedence relations.

*Table 1:*

| MODEL of business process | modelling SET(S) | TYPES of business process optimisation | APPROACHES to business process optimisation |
|---|---|---|---|
| –Petri–nets (and workflows) | –Diagrammatic models<br>–Mathematical/formal models | –Graph reduction techniques | – (Sadiq and Orlowska, 2000)<br>– (van der Aalst *et al.*, 2002)<br>– (Lin *et al.*, 2002) |
| –Mathematical models | –Mathematical/formal models | –Algorithmic approaches | – (Han, 2003)<br>– (Gutjahr *et al.*, 2000)<br>– (Jaeger *et al.*, 1995)<br>– (Hofacker and Vetschera, 2001)<br>– (Soliman, 1998)<br>– (Tiwari *et al.*, 2006)<br>– (Vergidis *et al.*, 2006)<br>– (Volkner and Werners, 2000)<br>– (Zhou and Chen, 2003a)<br>– (Zhou and Chen, 2002)<br>– (Zhou and Chen, 2003b) |
| | | –Activity/Task consolidation | – (Dewan *et al.*, 1998)<br>– (Rummel *et al.*, 2005) |

Zhou and Chen [1] developed a structured design methodology for business process optimization from strategic, tactical, and operational perspectives using quantitative methods that support the design. This optimization optimally assigns resource capabilities, organizational responsibilities and authorities, and organizational decision structure. Another approach to optimization is the consolidation of the activities (or tasks) of a business process. Rummel [9] proposes a model that focuses on decreasing the cycle time of an internal process by consolidating activities—assigning multiple activities to one actor—thereby eliminating the coordination and handoff delay between different activities that occurred when assigned to different actors. As this approach is activity focused, it ignores interactivity delay that may contribute significantly to overall process cycle time. Dewan [10] claims that there is no

structured methodology to determine the optimal re-bundling of information-intensive tasks. They present an approach to optimally consolidate tasks in order to reduce the overall process cycle time. The authors present a mathematical model to optimally redesign complex process networks but a limitation of the paper is that it refers to business processes with information flows only. Its main contribution is the effective business process restructuring and the reduction of the overall task time using handoff delay reduction or elimination as a result of a unified methodology applicable to multiple task-based business processes. Although formal languages have associated analysis techniques that can be used for investigating properties of processes, an optimization approach based on executable process languages was not observed in the literature. Since most of the optimization approaches—as discussed before—are based on algorithmic methods, these could be easily translated to executable software

programs. Analysis and optimization of business processes can be done best using an approach based on explicit and executable process models. Such models would allow evaluating performance in terms of flows, calculating costs against objectives, recognizing constraints, and evaluating the impact of internal and external events. Therefore, by being able to assess the process execution quality and costs, it is possible to take actions to improve and optimize process execution.

## 3 Key factors in financial and operational optimizations

There are various warning factors that signal the degradation of business processes. These factors are triggered either from the internal or external environment of the business. The signals are usually presented before the company enters into a crisis. There are cases where managers do not observe these signals or consider them to be some one-off or periodic difficulties. This approach does not only aggravate the outcome of the business, but, on a long term, threaten the existence of the organization itself [7]. Most companies agree to implement a change management plan. The overall review of the process can be split into five stages, as follows [8]:
- Stage 1: Analysis, usually takes from one week to one month.
- Stage 2: Planning – takes from one to three months.
- Stage 3: Implementation - six months to one or more than one year.
- Stage 4: Monitoring - six months to one year.
- Step 5: Return to business growth - from one year to two years.

### 3.1 Corporate approvals for business optimizations

The preparation of an effective restructuring plan (Slatter, Lovett, Barlow, 2006) is based on the following elements: crisis ending and business stabilization, appointment of a new director, stakeholders' management, strategic orientation, critical improvement process, implementation of organizational changes and financial optimization.

Six major milestones must be attended by an organization: (according to Downey, 2009):

**Step 1:** Changes in the management structure. It refers to bringing a new CEO or an external specialist. Involves the board of directors or senior management to recognize that a change is required and initiate a corporate review program.

**Stage 2:** Business Review. Rapid identification of the problems faced by society and assessment of business survival chances: strategy, operations, finance, infrastructure, people, commitment and ability to change.

**Stage 3:** Business Restructuring Plan. Establishing appropriate strategies and a well-structured recovery plan to deliver lasting results.

**Stage 4:** Implementation. Organizations can use sharp actions to save the company's performance: layoffs, department dismantling, and drastic cuts in all nonessential costs. Positive cash flow is critical and needs to be set quickly. In addition, the cash will be needed to implement the review strategy and must come in a timely fashion.

**Stage 5:** Stabilization. In this stage, the main aim is to increase the efficiency and the effectiveness of the business operations the focus is on. It is necessary to improve the profitability, but also to ensure the good functioning of the existing technologies.

**Stage 6:** Implement the change(s). The final stage is to implement the planned change, the organization / corporation restoring its financial loneliness. At this stage motivates staff and employees to achieve profitability and return on investment.

### 3.2 General corporate model for business restructuring through optimizations

Strategic analysis, monitoring, and strategic planning of the organization's/ corporations' activities can be monitored and controlled on the basis of several methods, models and

process diagrams. One such method is the **Critical Success Factors** method provided by Rockart (1979), which is based on the 80/20 Pareto rule for strategic management needs. The method involves identifying the critical success factors of the company as the results in important areas to ensure corporate success. This method identifies the most important business processes and the performance indicators (or KPIs) of these processes. It allows the plan to be compared with the results obtained, as well. The effectiveness of strategy implementation should be measured continuously in order to ensure continuous improvements of the operations. Key areas of corporate restructuring include: sales, finance, production and supply chain, management activities, services, business development, organization and human resources. Supply and logistics are usually embedded under the production processes or activities.

## 4. Managing Processes versus Projects

One important aspect is to make the difference between processes and the projects. Business processes consist of providing value to a customer through value-added activities, moving work across functional area boundaries, and controlling process performance indicators and standards and measuring process execution. Business processes are usually driven by facts or events, such as the maintenance of a factory, printing a product catalogue, the close of a billing cycle, or solving customers' issues in reconciling a checking account. These are activities that are typically replicated and repeated with specific resources allocated to an individual steering group such as factory line workers or customer service employees, to give some examples. Business processes are looking to the following core features: efficiency, agility and meeting customers' demands. While efficiency seeks to cut operational

costs and cost of capital, agility strives to cut the time required to develop products and services, and to respond to customer and market demands (thus through improving the effectiveness). Customer demands focuses on retaining the customers and their overall level of satisfaction.

A project, as defined by the Guide to the Project Management Body of Knowledge (PMI, 2004) is represented by a series of activities, and related tasks with a dedicated objective, bounded under a starting and an end date, and resources. A project consumes cash, people, time and equipment for the specified time period and defines what is planning to be done, when to do it, and ensures that the planned results are reached. The tasks of a project are unique and usually not repeated, and once the project is planned, changes to the plan are avoided to ensure meeting the schedule. While business processes look for efficiency, agility (effectiveness), and meeting customer demands, projects look to deliver the related objectives within established budget and time boundaries.

Thus, a processes optimization project is a short to medium term effort an organization puts upon to identify all necessary inefficiencies/ redundancies, decide the action plans to resolve them and ensure the optimization enhancements will meet the desired results without any negative outcome for the current operations.

## 4.1 Keys to Optimizing Processes

The executive management have to be equipped with the necessary tools to make the right decisions to realize the organizational course corrections with agility. The keys to optimizing process performance and execution capability are tied to the organizations' commitment to define and continuously assess and update the documentation of internal business processes. These documentations, including process maps, inputs and outputs, resource allocations, cycle times, etc., formally define the scope of the process from initiation to delivery and serve as the "process

roadmap".

Once the decisions are made, the Business Process Optimization ("BPO") projects concentrate the organization to identify and scrutinize opportunities to reduce costs, eliminate waste and reduce cycle times, while increasing products and services quality. There are cases when the vast majority of organizations continue to operate in silos. This is where the core business processes' activities must cross the traditional functional view and have neither an owner assigned nor measures of process execution success. Many successful organizations have mature business process management structures where the core and support business processes are well-defined, documented and assigned to an owner, measured for execution success, and scrutinized for efficiency, agility, and quality. Process management focuses on the management of the cross functional processes. This involves continuous monitoring, evaluation and measurement in terms of costs, quality, cycle time, etc. (Figure 1 below).



Fig. 1

Processes are well documented and monitored

KPIs are established at the level of each process

Process 1
Process 2
Process 3
Process 4
Process 5

1. The management of processes is clearly and structurally defined.
2. The performance of each process is monitored and looked after continuous improvements (response time, gaps, delivery time, etc.)

Formally defined processes and related documentation bring an organization with visualization of high-level and detailed core, supporting and management processes. Processes maps and supporting detail documents define how the

processes look like, their interactions across business functions, what each the process step delivers and how it produces its intended deliverables. Process owners own the core or supporting business process, not the individuals assigned to work the tasks within the process. The process owner is responsible to ensure that the execution of the process is successful and that it works to identify process issues, root causes and training needs. In order to manage a process that will bring continuous successful execution, process communications must ensure that process the information flow is established vertically and horizontally across the organization. A process owner must be able to monitor the external and internal flows of information. The purpose of process communication is to make sure that all employees are informed of process performance information and to control the company's progress toward its objectives and goals.

## 4.2    The    Efficiency    should    never compromise the Quality

For many consumers, the quality of the products and services is as important, if not more important than the cost of the product or service. While the focus of many business process optimization projects is centered on optimizing efficiencies and reducing the cycle times, businesses must continue to ensure that the business process optimizations does not compromise the quality of the product or services that are delivered by the internal processes. Six-Sigma methodologies are usually used to embed quality into process optimization projects. Six Sigma originated at Motorola in the early 1980's and is a methodology for disciplined issues solving and quality improvement. Six Sigma's goal is the near elimination of defects from any process, product, or service, limiting defects to just 3,4 defects per million opportunities. To ensure organizational alignment, Six Sigma methodology requires all improvement projects must be integrated with the goals of an organization. The DMAIC methodology

Six-Sigma (The Black Belt Memory Jogger, 2002) employs the following activities:

- Define: the phase whereby the customer needs are established and the processes and products to be improved are identified.
- Measure – determine the baseline and target performance of the process, defines input and output variables of process steps, and validates the measurement methods.
- Analyze – analysis of data to identify critical factors required for process execution.
- Improve – identification of necessary improvements/ optimizations (process, procedural, systemic, etc.) to optimize the outputs and eliminate and or reduce defects and variation. Statistically validates the new process operating conditions.
- Control – establishes the development of documents, monitors, and assigns overall responsibility for sustaining gains made by the implementation of process improvements.

# 5 BPO Project Management Methodology

Business Process Optimization projects follow the same method as defined by the Project Management Institute. Project management is obtained through the use of the process group such as: initiating, planning, executing, controlling, and closing. As stated in the Project Management Institute's Guide to the Project Management Body of Knowledge 2000 edition, "these process groups are linked by the results they produce- the result of one often becomes an input to another." Specifically:

- Initiating- approving the project or phase is part of the scope management;
- Planning- defining and refining objectives and selecting the best of the alternative courses of actions to reach the objectives that the project

was undertaken to address.
- Executing- coordinating people and other resources to carry out the plan.
- Controlling- ensuring that the project objectives are met by monitoring and measuring the KPIs periodically to identify variances from plan to identify corrective actions that can be taken, when necessary.
- Closing- formalizing acceptance of the project and bringing it to an orderly end.

## 5.1 Combining Six Sigma and Project Management Best Practices in the Initiating and Planning Process Groups

One of the single, most critical activities to ensure the success of a project, whether it be in the development of a software application, drug compound or optimizing a key business process, is the clear and concise definition of project objectives, goals and milestones in the projects planning phase. The purpose of the project should support the vision and mission statements of the organization and it requires the support and commitment of the management. Business process optimization projects should contain a section in its charter that defines the specific business process to improve. This formalized definition of the process optimization scope eliminates any confusion and formally defines the subject boundaries. Additionally, it assists in the identification of the established deliverables. For example, a fulfillment organization receives customer complaints on low order fill level. The customer places an order for a quantity of 100 for a particular item and receives only a quantity of 90. The project objective could be "to optimize the warehouse picking process to ensure an increase in the fill rate on customer orders from 90% to 99% by 4th Quarter 202X". The process scope has been narrowed specifically to the picking process and provides the basis for the process goal.

Six Sigma is a data driven problem solving methodology that requires the formal definition of performance key indicators. When planning for a process optimization

project, specific Six-Sigma tools and activities are used to characterize customer needs, and processes to be improved. These tools include the mapping of the high-level process in its current state, identification of the processes existing performance measures (i.e., pick time, product staging time) and a process financial analysis (i.e., resource cost, overhead). Specifically, Six-Sigma seeks to identify the Costs of Poor Quality (COPQ). COPQ includes costs of repairs, rework, rejections, inspection, testing and in the case of our fulfillment example, the cost of customer complaints. While a process optimization project's benefit can be measured financially (hard) or non-financially (soft) most business cases are based on the hard benefits. In the above example, the soft benefits of "improved customer satisfaction" should be considered as well.

While discussion of the "customers", Six-Sigma projects take the time to understand the needs of the customers. The project team must understand how the process issues link to the eventual customers. Six-Sigma mentions about the "*voice of the clients*" research to gain this important insight. There are many different methods to researching the customer's voice. These include, but are not limited to the following:

- Customer Complaint database- this is an acceptable place to start if the organization formally tracks issues;

- Direct Contact- if allowed, considers phone call surveys, focus groups, interviews at the point of provision.

- In-Direct Contact- includes mail surveys, feedback cards, market research and competitor analysis.

- Becoming the customer- order from your own distribution center, buy your own brand products, set up a new account with your own financial institution.

Another effective tool to use in a process optimization project is the SIPOC High Level Process Mapping tool. The acronym SIPOC stands for Suppliers, Inputs, Process, Outputs, Customer. It is a simple, but effective tool to align the project team and all stakeholders as to the core process within the scope of the project. It is important to mention that it is too early in the project to mention the existing process (that comes later in the Measure Phase).

The general approach to the SIPOC process identification includes the following steps:

- Begin with a simple definition of the in-scope process;
- Identify key steps of the process (expand these at the bottom of the SIPOC diagram);
- Have the project team identify the major inputs and outputs of the process;
- Have the team identify key suppliers of the inputs, and customer for each output.

Accordingly to any other project, a business process optimization project requires the formal identification of a project team with clear structure, roles and responsibilities. It can be used the SIPOC High Level Process Map to ensure all process stakeholders are represented on the main project team.

The Initiating and Planning Phase of a business process optimization project starts by formally identifying the process problem, not with the identification of the solutions to the problems/ issues. Six-Sigma tools such as the SIPOC, COPQ, and VOC help the project team identify the potential issues, process scope and essential process representatives, before the organization invests substantial time and money in the initiatives.

## 5.2 Measuring and Analyzing Current Process Performance

During the execution phase of a BPO project, the project manager is concentrated on executing the process optimization plan. These integral activities include the development of individual and team skills through the use of various team building exercises, reward and recognition systems and locating team member in the same physical area. The project manager is also focusing efforts to ensure the process

optimization plan is being carried out through regularly scheduled status meetings to exchange information about the specific project. During the execution phase, team efforts are focused in the identification of measurements to determine the effectiveness and efficiency of the process. There is necessary to develop the process measures which are critical for the process optimization project. It must identify and capture data on key performance indicators to determine process effectiveness and efficiency. Process effectiveness measures a customer's quantifiable service or product specifications. In addition, a process optimization project must track key performance indicators that reflect the internal efficiency of the process. In general, the following main steps are completed to measure the performance of a business process:
- Develop a data collection plan for the process;
- Identify process efficiency data collection sources;
- Identify process effectiveness data collection sources;
- Collect efficiency and effectiveness data to determine process performance baseline measurements.

**5.3 Controlling Key Business Process**
Following the development and testing of systemic, procedural, or responsibility enhancements, the BPO project team efforts should focus on ensuring the solutions are implemented and measured for their effectiveness. The project team must identify measures to be monitored after the desired state process is landed. This activity includes the identification of the persons responsible for collecting and analyzing the process data and reporting process efficiencies and effectiveness to the entire organization in the form of dashboards or status reports. Six-Sigma projects typically employ Statistical Process Control charts that monitor the

stability and variation of a particular process. A typical Statistical Process Control chart tracks the performance of a process over time and shows control boundaries which the results will lie between if the process is "in-control". Use of any Statistical Process Control charts require regular updating and review to ensure their feasibility. This ensures that process performance doesn't decline again.

Process change control is another key that ensures ongoing alignment with an organization's strategic goals. Processes are enabled by technological change, not hindered and that the appropriate organizational structure is in place to provide resources to support the business process. As documented in the Six Sigma Black Belt Guide (2001), a classical model for managing the change process has three phases: (1) unfreezing, (2) movement and (3) refreezing. Once a process change is identified and ready for deployment, the "unfreezing" of existing behavior patterns must be addressed. Typically, most work groups are resistant to changes and this must be solved. People or practices must then be moved (movement) to the process change by training or through technology adoption. Once, process resources have acquired the necessary skills and technology is in place, the process is then *refrozen* to ensure the process or function is aligned for organizational effectiveness. One effective technique used to facilitate the transition from existing processes to the new process is the use of a formal "White Paper Fair" where all functional areas impacted by the process changes have an opportunity to visualize the process enhancements.

**6. Implementation of optimization processes- market study results.**
Through a questionnaire developed according to the basic conceptual model of the project optimization implementation, process includes the following major specific modules, namely:
- Changes in business processes;

- Optimizations targets (planned and achieved);
- Optimizations implementation issues;
- Obtained benefits;
- Impact of the optimizations on corporate performance;
- Success factors of the optimizations.

The market study run over Romanian market (mostly energy sector) shows which are the main items of the above modules impacted by the optimizations.

## 6.1 Changes in business processes
The study shows that on average product design/ development, costs reductions aims, inventory management and production planning have led to a change in business processes to the largest extent. The sales and ordering, product design/ development, distribution, inventory management and production planning have determined the change in business processes measure.

On average, at the level of the study, only advertising/ promotion, billing/ payments and business planning received less attention.

## 6.2 Optimizations targets (planned and achieved)
The study shows that, as far as the objectives were included in the optimization plans, on average, at Romania level, the improvements based on automation, reduction of the costs and production costs, increase of competitiveness through costs reductions and the utilization of the novel technologies represented the main goals for which the objectives and targets were included in project plans. On the other hand, the increase in competitiveness through increased quality, concentration on the main results and the establishment of aggressive objectives received the lowest scores with regards to the optimization targets.

## 6.3 Optimizations implementation issues

The study shows that, having in mind the main implementation issues list at the planning level, on average, at the level of Romania, the available IT infrastructure does not support the planned optimizations, management reticence to allocate the funds and the business mistakes under the pressure of delivering the expected results were have been the major issues for the implementation of the optimizations.

## 6.4 Obtained benefits
The list of possible issues considered to be under the list of benefits of the optimizations' implementation found at the company level, as resulted from the market study, at the level of Romania, shows that the customer satisfaction level (improved response to clients' requests), concentrating the resources towards the selling aspects, increased flexibility through the adoption of new IT technologies and more efficient marketing and selling processes were the major benefits of the optimizations' implementation.

## 6.5 Impact of the optimizations on corporate performance
Based on the market study performed, the main items that were considered to be of impact of optimizations' implementation found at the company level, in Romania, shows that in terms of impact on corporate performance, on average, improving the development of new products, improving the costs' reductions and improving the investments return level, had a greater impact on corporate performance. On the other hand, the improvement of the sales rate, increasing the market rate and the improvement of the operational profits had the lowest impacts on corporate performance.

## 6.6 Success factors of the optimizations.
Based on the market study, the list of possible matters considered to be success factors at the optimizations level, found at the company level, show that the utilization of experts or external support, optimizations

driven by customer requirements and competitive pressures and involving all important employees represented to a greater extent impact on success in optimizations' implementation. On the other hand, the process mapping approach, the development of a well-defined project structure and using the internal surveys scored the lowest in terms of the main success factors for the optimizations done.

## 6.7 Summary of the study

The approach to optimizations is a broad one and aims at a sharp change in the quality of services offered, costs and production, including the analyses of the current state of scientific research in this field, based on the most recent and representative references in the relevant literature and interpretations and own contributions. This study aimed to provide certain contributions in the BPO environment, namely:

- to highlight factors and conditions necessary for the optimizations of financial and business companies;
- establish the optimization methodologies, but also how to apply them;
- Identify the common traits of business optimizations and/ or related methodologies, along with the actions and measures used;
- Developing a general optimizations model for companies, through which managers will be able to determine the reasons for the changes they want, as well as changes in the environment, all of which stimulate the need for improvements.

The results conclude that organizations do not focus on some of the most important tasks and actions recommended in the literature as a basis for optimizations, such as the use of time as a competitive advantage, changes in customer/ market business/ processes, the value-added item of each business activity and the application of the right

innovative technology. Therefore, one can assume that there is a major reason why many of the optimizations project objectives were only modestly achieved.

On average, the most common problems encountered in optimizations' implementation appear to be basic and difficult to solve in practice: implementation difficulties due to communication barriers between the organization/ functional sub-units, unexpected amount of optimization efforts required, interruption of operations, failure to achieve the expected benefits, pressure business mistakes to produce quick and overestimated results, and reluctance of top managers to commit the funds needed for the project.

Given that most optimizations benefit from innovative uses of information technology, an organizational problem that could condemn optimization projects to the failure of a particular company is the lack of communication between CEOs / top executives and CIO / IS managers.

## 7. Conclusions

Organizations can achieve sustainable and effective process improvement by combining project management best practices with certain Six Sigma methodologies and automation solutions. The ability to combine these proven methodologies provides the structure and discipline required to identify process improvement and optimization opportunities, develop sustainable solutions and lead the organization through the strategic change process. Use of these integrated techniques allows business processes to be efficient, agile, and meet the organization's customer demands. In today's challenging, global economy it is essential for organizations to combine the disciplines of Project Management, Six-Sigma and business process optimization to realize process gains that ensure "faster", "better", "cheaper" for their products or services, while maintaining a high level of quality in the marketplace.

## References

[1] Y. Zhou and Y. Chen, "Project-oriented business process performance optimization,", 2003, vol. 5.

[2] H. A. Reijers, "Product-based design of business processes applied within the financial services," vol. 34, 2002.

[3] I. Hofacker and R. Vetschera, "Algorithmical approaches to business process design," 2001.

[4] Y. Zhou and Y. Chen, "The methodology for business process optimized design," 2003

[5] A. Zakarian, "Analysis of process models: A fuzzy logic approach," vol. 17, 2001.

[6] K. Phalp and M. Shepperd, "Quantitative analysis of static models of processes,", 2000.

[7] Soininen, J., Puumalainen, K., Sjögrén, H., Syrjä, P., The impact of global economic crisis on SMEs, Management Research Review, 2012.

[8] Scherrer, P. S., Management turnarounds: diagnosing business ailments, Corporate Governance: The international journal of business in society, Vol. 3, 2003

[9] J. L. Rummel, Z.Walter, R.Dewan, and A. Seidmann, "Activity consolidation to improve responsiveness," vol. 161, 2005.

[10] R. Dewan, A. Seidmann, and Z. Walter, "Workflow optimization through task redesign in business information processes," 1998, vol. 1.

[11] Tristan Boutros și Jennifer Cardella, "the Basics of process improvement", Eng., 2016.

**Radu SAMOILA** has graduated the Master of Economy and Information Technology in 2011 at Bucharest University of Economy. Currently he is a PhD Student at this university, since 2019. Main fields of interest are business process optimization, automation of business processes and the continuous improvements concept.