

NOVA

IMS

Information
Management
School

MDSAA

Master Degree Program in
Data Science and Advanced Analytics

Business Cases with Data Science

Case 4: Millennium bcp – Business Process Conclusion Prediction

Bernardo, Pinto Leite, number: 20230978

Emília, Santos, number: 20230446

Nicolás, Zerené, number: 20230779

Ricardo, Kayseller, number: 20230450

Stepan, Kuznetsov, number: 20231002

Group F

NOVA Information Management School
Instituto Superior de Estatística e Gestão de Informação

Universidade Nova de Lisboa

May, 2024

INDEX

1. EXECUTIVE SUMMARY	2
2. BUSINESS NEEDS AND REQUIRED OUTCOME	2
2.1. Market Overview	2
2.2. SWOT Analysis	3
2.3. BCG Matrix.....	3
2.4. PESTEL Analysis.....	3
2.5. Competitive Landscape	4
3. METHODOLOGY.....	4
3.1. Data understanding	4
3.2. Data preparation	6
3.3. Modeling.....	9
3.4. Evaluation	10
4. RESULTS EVALUATION	10
5. DEPLOYMENT	11
6. CONCLUSIONS	12
6.1. Considerations for model improvement.....	12
7. REFERENCES.....	12
8. Annexes	13

INDEX OF GRAPHS

Graph 1 - Counts of Different Task Types	5
Graph 2 - Rejections by Activity ID.....	5

INDEX OF FIGURES

Figure 1 - Target definition.....	8
Figure 2 - Model's Evaluation.....	10

1. EXECUTIVE SUMMARY

Millennium bcp, the largest private Portuguese bank, embarked on a strategic project to enhance Business Process Management (BPM) capabilities with a focus on significantly improving the predictability of process outcomes. This initiative is important for maintaining the bank's competitive edge and operational efficiency in a highly dynamic financial sector. This project leverages extensive data analysis across four key datasets—task execution data, user information, specific request data, and rejection records—to enhance automation, advanced data management, and ensure high levels of accuracy and performance. By doing so, it aims to significantly improve customer satisfaction by providing more reliable and timely services.

The strategic approach of this project is to enhance predictive capabilities for determining whether requests will be rejected or successfully executed. To achieve these goals, comprehensive data management and advanced modeling techniques were implemented to enhance predictive accuracy. A range of predictive models was thoroughly evaluated, and their performance was assessed using critical business metrics to ensure robust and reliable results.

The Millennium bcp project is expected to deliver several key outcomes that will significantly enhance the bank's operational capabilities and customer service. Overall, it aims to position Millennium bcp as one of the leaders in digital transformation within the banking sector, improving customer loyalty through increased service reliability and speed.

2. BUSINESS NEEDS AND REQUIRED OUTCOME

2.1. MARKET OVERVIEW

The banking sector in Portugal is well-developed, with a significant presence of both local and international banks. Key players include Millennium bcp, Caixa Geral de Depósitos, Banco BPI, and Novo Banco. The market is characterized by a high degree of digitalization, with increasing investments in digital banking and fintech solutions. Portugal's banking market is dynamic, reflecting the country's economic growth and stability.

2.1.1. Customer Preferences

Digital Banking: Portuguese customers are increasingly favoring digital banking services such as online banking and mobile apps. This shift is driven by the convenience and accessibility of managing finances digitally. **Personalized Services:** There is a growing demand for personalized banking services that cater to individual needs and preferences.

2.1.2. Market Trends

Consolidation: The banking market in Portugal is seeing a trend towards consolidation. Banks are merging to improve efficiency and competitiveness, driven by changing customer preferences and technological advancements. **Sustainability:** There is a strong focus on sustainability and ethical banking practices. Banks are incorporating environmental and social considerations into their business strategies.

2.1.3. Local Special Circumstances

Regulatory Scrutiny: Following the global financial crisis, there has been increased regulatory scrutiny and a focus on risk management in Portugal's banking sector. **Demographic Changes:** The aging population and emigration trends in Portugal have impacted the demand for banking services, with a growing emphasis on retirement planning and international money transfers.

2.1.4. Underlying Macroeconomic Factors

Interest Rates: Low interest rates set by the European Central Bank have affected the profitability of banks, leading to a focus on cost-cutting and diversification of revenue streams. **Inflation and Economic Growth:** Economic stability and growth are crucial for shaping the demand for banking services. Consumer confidence and investment levels significantly impact the overall performance of the banking sector.

2.2. SWOT ANALYSIS

- **Strengths:** **Market Leadership:** Millennium bcp is the largest private bank in Portugal, offering a wide range of banking services including retail, commercial, and private banking. **Innovative Solutions:** Strong focus on digital transformation and automation, enhancing customer experience and operational efficiency. **Brand Recognition:** Well-established brand with a strong reputation in the Portuguese market.
- **Weaknesses:** **Regulatory Challenges:** Navigating complex regulatory environments can be resource-intensive. **Operational Costs:** High operational costs due to extensive branch networks and legacy systems.
- **Opportunities:** **Digital Banking Growth:** Increasing demand for digital banking services provides opportunities for expansion and innovation. **Market Expansion:** Potential to expand services into underserved segments or regions within Portugal and internationally.
- **Threats:** **Economic Volatility:** Economic fluctuations in Portugal and the broader Eurozone can impact banking operations. **Competitive Pressure:** Intense competition from both local banks and international players entering the market.

2.3. BCG MATRIX

- **Stars:** Digital banking services and mobile banking applications, which are high-growth areas with significant market share.
- **Cash Cows:** Traditional retail banking services that generate stable revenues.
- **Question Marks:** Emerging fintech solutions and blockchain technologies that require further investment to capture market share.
- **Dogs:** Legacy banking services that are costly to maintain and offer low returns.

2.4. PESTEL ANALYSIS

- **Political:** **Regulatory Environment:** Changes in banking regulations and policies by the Bank of Portugal can impact operations. **EU Membership:** Portugal's membership in the EU facilitates easier cross-border banking activities within the SEPA region.

- **Economic:** Economic Growth: The Portuguese economy is expected to grow, albeit modestly, which can positively influence banking activities. Interest Rates: Low-interest rates in the Eurozone impact net interest margins for banks.
- **Social:** Customer Preferences: Shift towards digital banking and the need for personalized financial services. Demographic Changes: Aging population may affect the types of banking services in demand.
- **Technological:** Digital Transformation: Adoption of AI, machine learning, and blockchain to improve banking services and operational efficiency. Cybersecurity: Increasing importance of robust cybersecurity measures to protect customer data.
- **Environmental:** Sustainability: Growing emphasis on sustainable banking practices and green financing. Climate Risks: Need to manage financial risks associated with climate change.
- **Legal:** Compliance: Adherence to stringent regulatory requirements and anti-money laundering (AML) laws. Data Protection: Ensuring compliance with GDPR and other data protection regulations.

2.5. COMPETITIVE LANDSCAPE

Millennium bcp competes with major banks like Santander, Caixa Geral de Depósitos, Banco BPI, and Novo Banco. Each bank is investing heavily in digital transformation to enhance customer experience and streamline operations. Additionally, international banks like Santander Totta and Deutsche Bank have a significant presence in Portugal, intensifying competition. Strategic Recommendations:

- Enhance Digital Services: Invest in advanced digital banking platforms and fintech collaborations to attract tech-savvy customers.
- Optimize Operations: Streamline operations by leveraging automation and AI to reduce costs and improve efficiency.
- Expand Market Reach: Explore new market segments and geographical regions to drive growth.
- Sustainability Initiatives: Implement sustainable banking practices to appeal to environmentally-conscious customers and comply with regulatory expectations.

By analyzing these aspects, Millennium bcp can strengthen its market position, enhance customer satisfaction, and drive sustainable growth in the competitive Portuguese banking sector.

3. METHODOLOGY

3.1. DATA UNDERSTANDING

3.1.1. Data collection and description

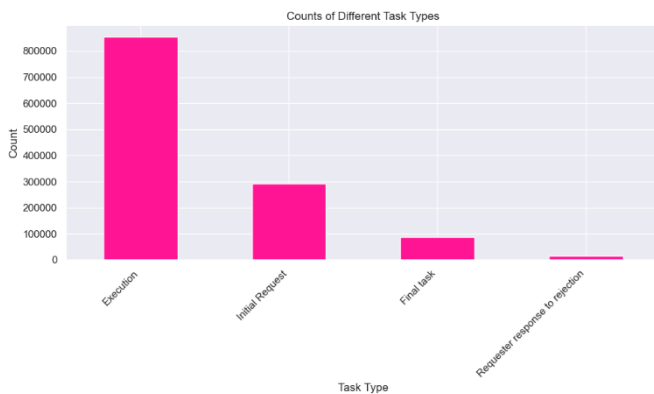
A CSV file from Millennium bcp was received, composed by 4 sheets, relative to information about task execution, user information, specific request data and rejections. About the first sheet, related to the tasks and their executions, it has 209017 observations and 12 variables. Some insights taken were that there might exist outliers on 'Request Identifier', 'Task executer department' and 'Task Executer', based on their minimum value, the 25% quartile and the mean, and on 'idBPMAApplicationAction' based on the maximum value, the 75% quartile and the mean. The second

sheet, 'Q2 – User Information', has 11370 observations and 7 variables. There are more males than females, managers than not managers and a higher number of not outsourcers in the dataset. Variable 'Sex' has 4 unique values, which will be assumed to be 'F' - Female, 'M' - Male, ' ' - robots (as mentioned by Millennium, do not have gender), 'U' - unknown. There might exist outliers on 'Task Executer', based on the minimum value, the 25% quartile and the mean, and a value of '2050' was found on 'BirthYear', which is not possible. Regarding the third sheet, due to confidentiality matters, there is not much information on the meaning of its three features and 297556 observations, but it was noticed that there might exist outliers on 'Request Identifier' and 'idField', based on the minimum value, the 25% quartile and the mean. Finally, the rejections sheet, where there are two variables, 4099 rejected tasks, and possible outliers on 'idBPMRequirement', comparing the maximum value (very high) with the third quartile (75%) and looking at the mean.

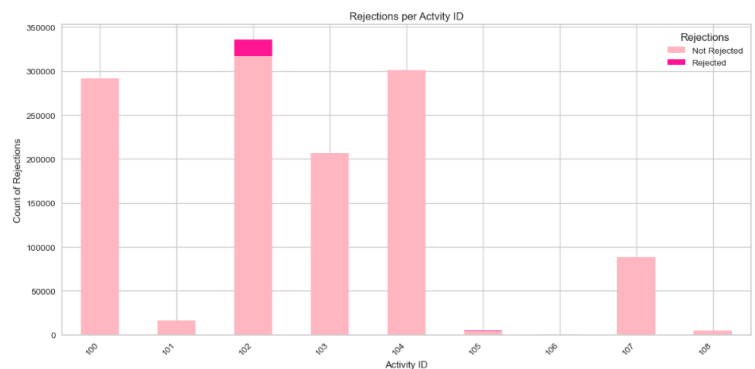
3.1.1. Data Exploration

Before continuing with data exploration and the upcoming steps, the 4 sheets were merged into one using the columns that were common among the sheets, in order to connect better their information, have more insightful analysis and only have the data about the business process in this case. Some pre-processing was done (for instance, the creation of the boolean variable 'Rejections', fitting in the information given by the fourth sheet) and important insights were extracted from the analysis.

The most occurring task is 'Execution', followed by 'Initial Request', then 'Final task' and, lastly, 'Requester response to rejection' (Graph 1). It was also concluded that most of the users are not outsourcers and are considered managers. Interestingly, the ones that are not managers and not outsourcers are related to the cases that were rejected (information present on the fourth sheet). The activity that occurs that most is 102, associated with different actions, for instance, administrative closure and tasks that were returned to the team. This activity is also the one related to the existence of rejections (Graph 2).



Graph 1 - Counts of Different Task Types



Graph 2 - Rejections by Activity ID

Upon the development of the two following Directly-Follows Graphs (DFG) (Annexes 1 and 2), it can be seen that there are three final actions, from which there are no arrows coming out, just coming in (except the ones looping in that same action): "Task terminated – administrative closure", "Task automatically terminated – SLA time reached" and "Request accepted by requester". Tasks that either leave the stage of being returned to the team or that were executed successfully can end up in any of those three final actions mentioned. An action that can lead also to administrative closure or to an automatically terminated task (SLA time reached) is the submission of an initial request.

Regarding the DFG relative to the activities, the most clear insight is that activity 100 is always a starting point, since there are arrows only pointing outside (except the loop). The numbers on the paths represent the frequency of transitions from one activity to another, highlighting the most common paths and bottlenecks within the process. Activities like ID 100 and 104 have high frequencies, indicating they are common steps, while the transitions to activities like ID 107 show critical junctures where decisions or significant actions occur.

3.1.3. Data Quality

Regarding **missing values**, only some were encountered in the first sheet (task execution data), on 'Action', 'Task predicted end date', 'Task Executer' and 'Task executer department' variables. Missing values on 'Action', are not actually missing values are empty values (checked on Excel) and correspond to specific values of 'idBPMAApplicationAction', such as 273 and 271, meaning that these ID's represent always empty actions. The missing values on 'Task predicted end date' are based on the actions that don't need a time interval to have them done, for example, it takes about three minutes to fill in a form, there is no need to measure this action. These are not mistakes or typos, have a meaning. Regarding 'Task Executer' and 'Task executer department', these were assumed to be individuals that probably left the company.

About **duplicates**, were just identified 24 in the second sheet, which were revealed to be indeed duplicated rows, that were consequently dropped.

3.2. DATA PREPARATION

In the initial stages of our analysis, a thorough examination of the dataset was conducted to ensure the quality and usability of the data for further analysis. Key steps included:

3.2.1. Correcting data types:

Dates were converted to date type variables (previously, objects), 'Task Executer' and 'Task executer department' to integer (previously, floats) and Rejections to object (previously, integer).

3.2.2. Missing Values:

After merging the datasets several strategies were employed to address missing values. For the 'Task predicted end date' column, missing values were filled using the corresponding 'Task execution end date' values, as these tasks did not require explicit time intervals. Missing values in the 'Action' column were confirmed as intentional empty values based on domain knowledge, so these values were filled with 'No Action'. The task executer departments that were missing corresponded to specific task executers, knowing this a symbolic value to the nan values, '0', in this case, according to the data type.

3.2.3. Duplicates:

The data was checked for duplicates, particularly focusing on the 'Task Id' and 'Request Identifier' columns. The analysis involved grouping by these identifiers and counting unique occurrences to identify any duplications. Duplicate entries were then removed to maintain the uniqueness and integrity of the dataset. The process ensured that no 'Task Id' corresponded to more than one 'Request Identifier' and the other way round.

3.2.4. Outliers:

Outliers were detected using a combination of statistical methods and domain-specific knowledge.

- **Z-Score Method:** Outliers were identified by calculating the z-scores for each numeric column. Data points with z-scores greater than 3 or less than -3 were flagged as outliers.
- **Interquartile Range (IQR) Method:** Outliers were detected by computing the IQR for each numeric column. Values falling below $Q1 - 1.5IQR$ or above $Q3 + 1.5IQR$ were considered outliers. This method helped identify outliers by focusing on the spread of the central 50% of the data.
- **Isolation Forest:** This unsupervised learning method was used to detect outliers. By training an Isolation Forest model on each numeric column, data points (i.e. outliers) that were isolated were identified.

Comparing the results given by these methods, some actions were taken. For the 'OrgUnitSince' variable, values below 2005 were removed to maintain a consistent and relevant timeframe for the data. Variables like 'BirthYear' were transformed into more interpretable forms (e.g., age) after outlier detection. This transformation helped in maintaining the practical relevance of the data. Overall, these methods ensured that outliers, which could bias the analysis, were appropriately managed. Later on, after analyzing the target values, the group decided to classify as outliers the instances with value 'Unknown' on this variable (further discussed).

3.2.5. Incoherencies

This section addresses the identification and resolution of data incoherencies or inconsistencies in the dataset. Key issues included invalid 'BirthYear' entries, anomalies in 'Task Executer' ID's, and possible inconsistent values in the 'Sex' field (robots), which did not occur after the merge of the 4 sheets.

3.2.6. Feature Engineering

Extract Month and Year from Task Arrival Date:

The task arrival date, initially in a raw date-time format, is decomposed into separate month and year variables. This transformation allows for the capture of temporal patterns and trends, enhancing the model's ability to learn seasonal behaviors or yearly cycles in task arrivals.

Create Variable Age:

An age variable is derived to quantify the duration or age of entities within the dataset.

Create Variable Time Difference:

The time difference between significant events or timestamps is calculated to understand the time lags within processes. This variable reveals delays or expedite patterns, contributing to a more granular understanding of the temporal sequences and their impact on outcomes.

Replacing Categorical Values with Numerical:

Categorical variables are converted into numerical values, facilitating their use in future engineering, for instance, 'F' and 'M' in 'Sex', converted to '0' and '1' (done in 'Sex', 'Is Manager', 'IsOutsourcer', 'Task Type' and 'Action', which were then turned into objects).

Define Target Variable:

The target variable is determined based on specific conditions applied to the dataset. The logic checks various activity IDs and action codes to classify the target variable into 'Closed administratively', 'Request canceled', 'Request finished' and 'Closed Administratively (requester rejects accounting impact)'. By defining the target variable through a structured and logical approach (Figure 1), we transform raw execution data into a critical component for machine learning models. This process ensures that the target variable accurately reflects the final states of tasks. Some of the values of the target were classified as 'Unknown', when they didn't meet the criteria for any predefined categories. An analysis of the 'Unknown' targets was performed, during which the corresponding 'Activity ID' values were examined. Graphs were generated to illustrate these findings, revealing that the final activities 100, 102, 105, 101, 103, and 107 were redundant. Consequently, it was decided to remove these activities from the dataset.

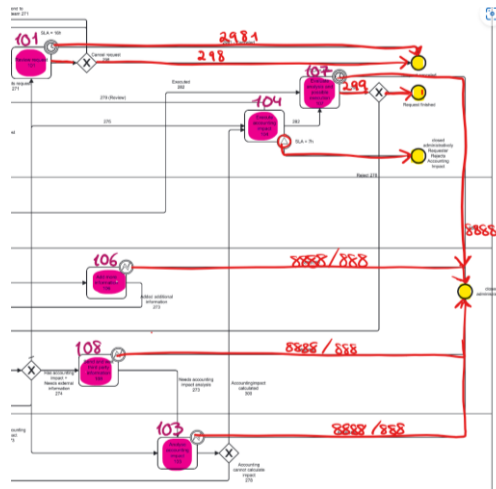


Figure 1 - Target definition

3.2.8. Feature Selection

Since the multi-classification predictive models would be trained based on the extracted paths and their prefixes, there was no need to develop this stage in depth. So, some of these features were removed looking at the correlation matrix and comments existent on the Profiling Report, others were selected based on our understanding of the problem.

3.2.7. Paths Extraction

In order to extract the paths of each 'Requester Identifier', an understanding of the historical traces and the sequences of events was made. These traces are crucial as they provide the foundational data from which we derive prefixes of varying lengths. The prefix extraction step involves segmenting these traces into subsequences, 'trace prefixes'. We obtained a minimum of 1 path and a maximum of 25.

3.2.8. Data frames creation based on prefixes

Each trace prefix was systematically categorized into buckets or data frames. This categorization is based on the encoded characteristics of the trace prefixes, paths with the same length are in the same data frame. The primary purpose of this bucketing step is to segment the data in a manner that allows for more tailored and effective predictive modeling and better enable the prediction of the outcome based on the stage in which someone is at during the process.

In practice, to facilitate targeted predictive analysis, data frames were created for each prefix length. For instance, a separate data frame for prefix 1 contains all trace prefixes that consist of just

the first step in the sequence. Each subsequent prefix length (prefix 2, prefix 3, etc.) is treated similarly. This segmentation aids in understanding the progression of events at different stages of the process and how early in the sequence accurate predictions can be made.

3.2.9. Encoding

The group applied a frequency-based encoding in the data frames used to train the models, this method basically replaces each category, in this case each activity, with the count of how often it appears in the dataset. This method not only simplifies the representation of our data but also helps in emphasizing the most common patterns that occur in our processes. 'Target' and other categorical variables, such as 'Sex' and 'IsOutSourcer', for instance, were label encoded.

3.2.10. Scaling

In order to scale the data, we resorted to the Standard Scaler, which was applied in each data frame. This method transforms the data, more precisely, each feature gets a the distribution of mean 0 and standard deviation of 1.

3.3. MODELING

The structured approach of using prefix-specific data frames plays a pivotal role in our predictive modeling. These data frames encapsulate the frequency-encoded values of process activities at various prefix lengths, which fundamentally influences the effectiveness of our predictive models. The data frames that contained only 1 or 2 activities (df_1 and df_2) were not considered, as well as the ones that contained less than 10 rows of data (from df_16 until df_25).

The group developed a function ('models_function') which would implement the models and other aspects for each one of the data frames, with efficiency and ease. This function includes the following steps.

- Firstly, data integrity assurance by removing columns entirely filled with NaN values (for instance, df_1 only had 1 column with 1 activity, but had NaN columns until the twenty-fifth activity);
- Isolation of the target variable isolated from predictors;
- Implementation of the standard scaler, explained before;
- Class imbalance management by implementing the Synthetic Minority Over-sampling Technique (SMOTE), improving the model's generalization capabilities across underrepresented classes (it was noted after encoding that target 0 was 4.21% of the time, target 1 was 73.42%, target 2 was 1.98% and target 3 was 20.39%);
- Model implementation, being the models selected for evaluation: Logistic Regression, Support Vector Machines (SVM), K-Nearest Neighbors (KNN), Naive Bayes and Random Forest;
- Each model was subjected to Stratified k-Fold cross-validation to maintain the proportion of classes across each fold, crucial for maintaining validation integrity given the class imbalance (train and test obtained);
- Parameter tuning was conducted via GridSearchCV in each model, which exhaustively searched through a predefined grid of parameters for each model. This approach not only helped in identifying the optimal parameters for each model but also ensured that the models were not biased towards the majority class;
- Application of model evaluation metrics, in order to compare and evaluate the models.

3.4. EVALUATION

The performance of each model was evaluated using metrics such as accuracy, precision, recall, and F1 Score. These metrics provided a comprehensive view of the models' abilities to predict accurately while also balancing the rate of false positives and false negatives (critical factors in the practical application of predictive models). It was decided to use the model with the best F1 Score as the model to create the predictions due to the imbalanced nature of the data to give a better explanation of the results. For this reason, the model Random Forest showed better results overall. The average F1 Score across all dataframes with Random Forest was 0.66 (image below).

	f1_score				
	LogReg	SVM	KNN	NB	RF
df_3	0,683816	0,755069	0,672075	0,714817	0,687117
df_4	0,729223	0,736495	0,655928	0,557508	0,733974
df_5	0,710884	0,63098	0,644264	0,737766	0,726496
df_6	0,71635	0,743881	0,726034	0,753316	0,726805
df_7	0,715023	0,706215	0,696441	0,638072	0,705847
df_8	0,691884	0,593318	0,528435	0,505189	0,721515
df_9	0,681214	0,628739	0,636076	0,668808	0,743874
df_10	0,55178	0,570417	0,581012	0,410635	0,58712
df_11	0,639896	0,632215	0,457115	0,616171	0,5654
df_12	0,556268	0,455086	0,55812	0,650893	0,577566
df_13	0,451852	0,400529	0,393292	0,531481	0,545899
df_14	0,573016	0,370723	0,389683	0,560847	0,607143
	2	2		3	5
				average RF	0,66073

Figure 2 - Model's Evaluation

This shows that the model is not biased and it is not overfitting. Which for the purposes of this project is a very good percentage.

4. RESULTS EVALUATION

The first question which the group was meant to answer was related to the prediction of the outcome, if the task would be rejected (at some point going to activity 101) or not (never went to activity 101). Using the predictions obtained from Random Forest and doing some engineering in the dataset, the conclusion taken was that 95.19% of the tasks did not go through 101 and, consequently, about 4.81% went through it. From this 5%, 40% of the requests were canceled and 20% were finished, being the rest of the tasks closed administratively.

Regarding the second and last inquiry, 'When executing the request (102 or 105) what is he predicted outcome?', it was seen that 94.93% of the tasks that went through 102 and 105, did not go through 101, and, consequently, 5.07% of the tasks that went through 102 and 105 went through 101. Having that in consideration, about 93% of all tasks had either 102 or 105.

To sum up, the grand majority of the tasks did not go through 101, so they were not cancelled or rejected. From the ones that went through 102 or 105, also the majority was not cancelled or rejected.

5. DEPLOYMENT

To ensure an effective deployment strategy for Millennium bcp's predictive modeling initiative, it's essential to integrate these models seamlessly into the bank's existing processes and decision-making workflows. Here's a revised deployment strategy that considers the best practices identified in the industry:

- **Integration with Business Processes:**

Ensure that predictive models are fully integrated with key business processes. This can be done by embedding analytics directly into workflow tools that are already in use, ensuring minimal disruption and enhancing user adoption.

- **User Adoption and Training:**

Focus on comprehensive training and support for the bank's staff to ensure they understand how to utilize the new predictive tools. This includes detailed training sessions, easy-to-access online resources, and ongoing support to address any questions or issues that arise.

- **Monitoring and Continuous Improvement:**

Continuously monitor the performance of deployed models to ensure they are providing the expected outcomes. Use real-time data to make adjustments as needed. Establish metrics for success and regularly review these metrics to evaluate the effectiveness of the predictive models.

- **Technology and Infrastructure:**

Evaluate and, if necessary, upgrade the technological infrastructure to support the deployment and smooth running of predictive models. This includes ensuring that data handling and processing capabilities are adequate and that data security measures are robust.

- **Stakeholder Engagement:**

Engage with stakeholders throughout the deployment process to ensure that the models meet the practical needs of the business and to foster a sense of ownership among those who will use these tools daily. Regular feedback sessions can help adjust the deployment strategy to better meet user needs.

- **Scalability and Expansion:**

Plan for the scalability of the models to accommodate future growth and potential expansions in their application. This may involve scaling up the models to handle larger datasets or refining them to enhance their precision and applicability to different types of customer interactions.

By following these guidelines, Millennium BCP can maximize the impact of its predictive analytics initiative, enhancing decision-making processes and ultimately improving customer satisfaction and business performance.

6. CONCLUSIONS

Having a good predictive analysis over Business Process Management allows for a strategic allocation of resources, that can save money to the company. With this project, the group managed to create an optimised model that can predict the outcome of a process depending on the activity. With 66% of confidence, the Random Forest model is the best model that can save resources and money to Millenium BCP. Therefore, the group predicts that about 95% of the tasks will end up not going through activity 101, out of all tasks only about 2% will be canceled, about 1% will be finished and the 3% of tasks will be closed administratively.

In other words, for a specific Request Identifier at Activity 102, the group predicts with 66% of confidence that the task will end up being “Closed administratively/Requester rejects accounting impact” about 73.7% of the time.

With these conclusions, Millenium BCP can allocate the necessary resources to make sure the attention is on the important requests.

6.1. CONSIDERATIONS FOR MODEL IMPROVEMENT

To improve the models, predictions and have better insights to enhance bank's operational capabilities and customer service, in terms of data, the group could have had more information about the clients involved in each task, for instance, demographic and behavioural data, and could as well have the possible monetary value involved in each process or task, from the side of the client.

7. REFERENCES

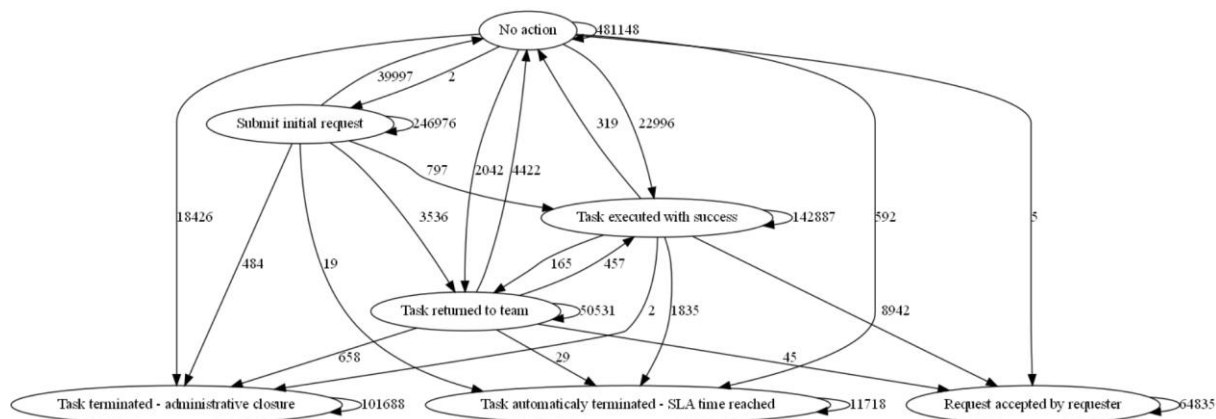
- Team, D. (2023, October 10). Top 5 outlier detection methods Every data enthusiast must know. DataHeroes. <https://dataheroes.ai/blog/outlier-detection-methods-every-data-enthusiast-must-know/>
- Ninja, N. (2024, April 1). Frequency encoding: counting categories for representation. Let's Data Science. <https://letsdatascience.com/frequency-encoding/>
- Quem somos - Conheça o Banco Millennium bcp - Millennium bcp. (n.d.). <https://ind.millenniumbcp.pt/pt/Institucional/quemsomos/Pages/quem.aspx>
- Millennium bcp. (2024, April 1). Millennium bcp is the best investment bank in Portugal. Millennium bcp. Available at https://ind.millenniumbcp.pt/pt/Institucional/imprensa/Documents/2024/Millenniumbcp_é_o_Melhor_Banco_Investimento_em_Portugal.pdf

- A Nossa História - Millennium bcp. (n.d.).

<https://ind.millenniumbcp.pt/pt/Institucional/quemsomos/Pages/historia.aspx>

8. ANNEXES

Annex 1 - Directly-Follows Graph for Actions



Annex 2 - Directly-Follows Graph for Activities

