

Machine learning approach to support ticket forecasting from software logs

Matti Haukilintu

School of Electrical Engineering

Thesis submitted for examination for the degree of Master of Science in Technology.

Espoo 31.07.2022

Supervisor

Prof. Arto Visala

Advisor

MSc Petri Pyöriä

Copyright © 2022 Matti Haukilintu

Author Matti Haukilintu

Title Machine learning approach to support ticket forecasting from software logs

Degree programme Automation and electrical engineering

Major Control, Robotics and Autonomous Systems **Code of major** ELEC3025

Supervisor Prof. Arto Visala

Advisor MSc Petri Pyöriä

Date 31.07.2022 **Number of pages** 32+1 **Language** English

Abstract

Your abstract in English. Keep the abstract short. The abstract explains your research topic, the methods you have used, and the results you obtained. In the PDF/A format of this thesis, in addition to the abstract page, the abstract text is written into the pdf file's metadata. Write here the text that goes into the metadata. The metadata cannot contain special characters, linebreak or paragraph break characters, so these must not be used here. If your abstract does not contain special characters and it does not require paragraphs, you may take advantage of the abstracttext macro (see the comment below). Otherwise, the metadata abstract text must be identical to the text on the abstract page.

Keywords Machine learning algorithms, Robotic process automation, log analyzing, random delay, Microsoft Azure ML Studio



Tekijä Matti Haukilintu

Työn nimi Sovelluslokien ja vikatikettien yhteyden löytäminen koneoppimista hyödyntäen

Koulutusohjelma Automaatio- ja sähkötekniikka

Pääaine Ohjaus, robotiikka ja autonomiset järjestelmät

Pääaineen koodi ELEC3025

Työn valvoja Prof. Arto Visala

Työn ohjaaja FM Petri Pyöriä

Päivämäärä 31.07.2022

Sivumäärä 32+1

Kieli Englanti

Tiivistelmä

Tiivistelmässä on lyhyt selvitys kirjoituksen tärkeimmästä sisällöstä: mitä ja miten on tutkittu, sekä mitä tuloksia on saatu.

Avainsanat Koneoppiminen, koneoppimisalgoritmit, ohjelmistorobotiikka, loki, lokin analysointi, satunnainen viive, satunnaisviive, Microsoft Azure ML Studio

Preface

Pitkäksi venynyt työ
Kiitokset

Otaniemi, June 21, 2022

Matti Haukilintu

Contents

Abstract	3
Abstract (in Finnish)	4
Preface	5
Contents	6
Symbols and abbreviations	8
1 Introduction	9
1.1 Background and motivation	10
1.2 Research objectives	11
1.3 Scope	12
1.4 Structure	13
2 Background	14
2.1 Machine learning field	14
2.2 Cloud ML platforms	15
2.3 Azure ML Studio	16
2.4 Regression algorithm	16
2.5 PCA-based anomaly detection	16
2.6 N-gram features and feature hashing	17
2.7 Data sensitivity	17
2.8 Log data analyzing with ML	17
2.9 Random delay in log event analyzing	18
2.10 Hybrid machine learning approach in anomaly detection	18
3 Research material and methods	20
3.1 Support ticket data	20
3.2 RPA log data	20
3.3 Data anonymization	21
3.4 Data formatting	22
3.5 Azure environment	23
3.6 Machine learning pipeline	23
4 Results	26
4.1 Azure and Azure ML Studio	26
4.2 ML training and validation	27
5 Summary	30
5.1 Discussion	30
References	31

A Esimerkki liitteestä**33**

Symbols and abbreviations

Symbols

P placeholder-symbol
A another placeholder-symbol

Operators

$\nabla \times \mathbf{A}$ curl of vectoring **A**
+ yep, a plus

Abbreviations

AI	Artificial Intelligence
ML	Machine Learning
HML	Hybrid Machine Learning
RPA	Robotic Process Automation
SQL	Structured Query Language
JSON	JavaScript Object Notation
CSV	Comma-Separated Values
GDPR	General Data Protection Regulation
ADA	Anomaly Detection Algorithm

1 Introduction

Artificial intelligence (AI) and machine learning (ML) has found their way to more and more fields of business. In banking business they are already used in fraud detection, risk management and service recommendations.[1] Even though these modern big data utilizing technologies are widely used abroad, in Finnish banking field AI and ML are not popularly utilized. Instead, many self-acting solutions are being used to streamline manual labor which could be called intelligent, but are merely highly automated processes and thus cannot be included in the AI category. One of these technologies used in Finnish banking systems is Robotic Process Automation (RPA).

RPA operates “on the user interface of other computer systems in the way a human would do”,[2] but is strictly bounded by predefined operations thus being prone to unforeseen situations such as faulty input. RPA, like generally all other software, produces log to “register the automatically produced and time-stamped documentation of events, behaviors and conditions relevant to a particular system”[3]. Logs don’t have any standards or form guidelines to follow which tends to make log analysis and log based problem-solving troublesome. This is also the case with RPA’s developed by Oy Samlink Ab.

Oy Samlink Ab (Samlink from now on) was founded in 1994 and is now owned by Kyndryl. From the early years while going by the name of Samcom the company was owned by several Finnish banks developing all sorts IT solutions for them. Nowadays, Samlink offers a wide variety of banking solutions from basic banking system to Codeapp-mobile software.

Besides banking, Samlink develops multiple other IT solutions to extensive range of customers, for example, entertainment platform solutions for DNA. Even though Samlink can be considered a modern technology company, the most modern AI technologies has not yet been adopted in the variety of tools used in development. However, RPA has been actively used in some banking solutions to reduce the amount of manual labor required from banking clerks.

In addition to continuous development as well as product maintenance services, Samlink also offers a technical help desk regarding the software solutions produced. As no IT solutions comes without bugs or misbehaviour, Samlink service desk has to use considerable amount of labor to resolve if technical support request tickets received are due to the problems in programming and require fixing in production. In many cases, the problem-solving starts by reading the log and analyzing the data written by processes in question.

In this study, we aim to find if it is possible to utilize ML methods in analyzing logs created by Samlink RPA’s. Ultimately, we intend to train ML which is able to predict the arrival of a technical support ticket thus giving a warning for developers about possible issues in the production.

1.1 Background and motivation

In the field of information technology logging is one of the most important methods in problem-solving, be it software or operating system related.[3] Typically, at least in Samlink processes, logging is a bit more verbose than it needs to be. This is usually because when the problem occurs it is easier to already have the verbose logs available than trying to replicate the issue after setting logging to more verbose mode. Too verbose logging, however, leads into two problematic issues for Samlink.

First of, the size of log is huge and finding the critical lead for solving the problem in hand takes more time. Of course, with more strict logging pinpointing the issue from within the logs would be faster, but then again, solving the problem with only critical error messages could be more time-consuming if crucial context is missing.

Secondly, a well-designed software is able to retry the process after first failure, but logging is done in real time, not after final results of the process has been determined. This means, that each process failure is logged even though said process succeeds eventually. Thus, logs may include dozens of rows of information about a problem encountered, which are not critical information after all. These issues make log analyzing considerably laborious.

Production logs are usually not viewed if everything is presumably working as intended. Technical support tickets are both last and most visible indicator that something is wrong. When a technical support ticket is received from banking clerks it means that something is wrong in a very visible way. Roughly speaking, technical tickets that are due to clear misbehavior of the RPA systems and not, for example, user errors, can be divided in two categories. First are the tickets that uncover an unknown bug in the system which can be either fixed or instructed to user how to avoid. Second type of tickets are somewhat pre-known issues that occur from time to time and are either fixed with updating parts of the system or by rebooting the process.

Typically, in software systems, if issue is known and can be fixed by rebooting something, developers can create log monitors that search for certain keywords and raise an alert if encountered. However, when it comes to RPA's and technical tickets considering them, it hard to say which type of issue is in question by reading the RPA logs word by word without context. New kind of issues can be more frequent than already confronted ones, and clear keyword linked to a certain problem may not exist without considerable amount of false positive matches.

ML algorithms are widely used to find patterns from massive amount of data making it an ideal tool for log analyzing. Patterns, however, need a connection to a visible issue. If RPA system has encountered an error but is able to retry successfully, then no issue has practically happened that needs immediate concern. Hence, RPA log analyzing with ML can find meaningful patterns only if they relate to actual technical tickets.

If Samlink support has received a help request the issue behind the request is not fresh anymore. In the event of RPA job failing, it takes some time for the clerk to notice the issue, write a help request to first support level, which then redirects the ticket to corresponding team. Furthermore, if the issue is noticed during friday,

it takes few more days to be handled by RPA developers. This leads to noticeable delay in processes that were supposed to be dealt by RPA but which now have to be manually taken care of by said clerks.

TODO: Add something about why ticket forecasting would be desirable

TODO: something about log format?

Comment: Otherwise, this chapter is good to go!

1.2 Research objectives

This research aims to kickstart the ML application usage in Samlink operations. Multiple obstacles need to be tackled as most of the phases in this study has not yet been encountered inside the company.

First and foremost, it is crucial to construct some basic form for log data to make it usable by ML algorithms. Log data forming is one of the key elements in automatic log analyzing applications as it is not for just machines but also for people to read.

As today more and more concern is set on anonymization not only due to GDPR, the data used for machine learning must be sanitized. Because of this, one major objective is to create a clean dataset that is safe to use in cloud environment without raising concern around security and privacy issues. In addition to this, data must also be clean enough so that ML algorithms are able to process it.

Finally, the main objective this research is aiming to answer is whether it is possible to use ML algorithms in such ways that combine support ticket timestamps and software run log so that an algorithm can predict if new support ticket might be coming.

- * Counting anomaly "probability" for individual log rows

- * Random delay! Solving issue by grouping

- * Hybrid ML with anomaly detection and regression

=> Can this give any results for forecasting?

?? Assumption to lean on:

- !! Ticket causing log rows are ANOMALIES in log data.

TODO: explain why this assumption is necessary!!!

- Log amount is not necessarily proportional to ticket creating error amount, e

- Considerable amount of rows may be related to one, or several different, tick

What really happens in this research:

<Anonymization scope>

<Azure ML studio setup and prerequisites>

<Log data and timestamp combining>

<Connection for data, ML estimates>

1.3 Scope

In order to limit the study to feasible length and content, it is necessary to define the scope for the thesis. As major part of the time consumed for the study was used for data acquiring and anonymizing yet the most value comes from content considering machine learning, the length of the content per section does not reflect the amount of time used for each section.

Before diving in to the scope of research objectives, we must first make one assumption regarding the data that ML is going to find some meaning from. In order to find a connection between log data and support tickets, we make a hypothesis that errors leading to tickets are visible to or parseable by ML algorithm. In addition, as we are going to utilize anomaly detection algorithm in log analyzing, we must also assume that these errors in log are, in fact, anomalies. We will, however, compare the results against these assumptions to test this hypothesis.

Data anonymization

Anonymization in the context of this thesis refers to data sanitization process purposed to edit the data into more secure form in the privacy point of view. In this study we aim to create a dataset usable in ML training. In this respect anonymization is not in the main focus of the study but only treated as a sub-phase of the data preprocessing in whole. Nevertheless, anonymization is from the privacy point of view the most important phase of data preprocessing.

Keeping this in mind, anonymization is covered rather superficially, only enough to explain the reasons behind actions taken during anonymization process.

Azure setup

The ML training and result scoring is done in Azure ML environment. Azure is used because... As part of the study aims to create an initial guideline for ML process commissioning the Azure setup phase is documented in such detail reflecting the importance of this information for future developers starting Azure ML projects in Samlink.

TODO: Maybe something more

Data requirements

Data purity in a sense of how easy it is to be used by ML algorithms creates challenges at the beginning of ML training. If data is not consistent, has lots of missing values or is formed in unanticipated way, it requires considerable amount of preprocessing slowing the training process and causing errors in pipeline runs.

In order to create a baseline for Samlink ML projects this study aims to give basic criterion what is required from the data, so it is easily analyzable by ML algorithm.

Machine learning methods

Several different machine learning algorithms exist that are aimed for different applications in mind. For example, to make an algorithm that can predict the price of an apartment listed[4] we could be using linear regression, and in order to detect possible cyber threats from network traffic[5] a two-class support vector machine could be utilized. These two methods are very different in usage and has their pros and cons in different applications.

As different methods can be used in creative ways in very different applications depending on how the data is presented and how the ML problem is formed, this study focuses on just a few easily approachable training methods.

When it comes to anomaly detection, only principal component analysis (PCA) is considered because PCA-based Anomaly Detection component is the only one usable from existing anomaly detection algorithms in Azure ML Studio. As purely Azure ML Studio is used during this study, no other anomaly detection algorithms are debated.

1.4 Structure

TODO: Open up what we are going to discuss here and in what order

2 Background

Machine learning, or ML, is a subcategory of the AI field and data science. Typically, ML refers to a set of technologies used to “build computers that improve automatically through experience”.[6] This is generally considered a machine way to simulate human learning process. ML usage has become more common and is nowadays widely used in many fields, not just in general information technology and computer science. This is because data can be gathered everywhere, and where there is data to be processed, ML can be there to process it. Computer algorithms are able to find statistical correlation and patterns from places overlooked by human mind, or where amount of data is just too much for people to process. This is why ML has proved its power in various empirical science fields, such as biology, cosmology or social science.[6]

In this section, key concepts of ML are explained briefly and several ML features are explored that are most relevant to this study. We also discuss shortly about data sensitivity and how it had to be addressed during this study.

2.1 Machine learning field

Machine learning is built with algorithms that operate

ML algorithms

Algorithm means a finite sequence of (typically) mathematical operations that are used to solve a specific problem, typically by repetition of some steps until the problem resolves.[7] Algorithms are the main “magic” inside ML where repeating patterns are searched from within the data by iterating through all the data points.

In ML, algorithms can be generally divided into <four> categories

TODO: Check categories, create list or table

Regression algorithms predict values and are typically used with supervised learning.

TODO: Examples: House market price?

Classification algorithms predict categories. Depending on the algorithm, they can predict between two or several categories.

TODO: Examples?

Clustering algorithms use unsupervised learning to find structures inside data.

TODO: Examples?

Anomaly detection algorithms work also unsupervised and try to find unusual or rare data points from data.

TODO: Examples?

TODO: Something more generally about algorithms... Intro to next section (training).

ML training

ML algorithms can be trained in two basic ways: supervised and unsupervised learning.

In supervised learning, algorithm is given data with ready answers on how the data needs to be interpreted. Algorithm then tries to figure out how given data and the correct answers are related.

In unsupervised learning, on the other hand, algorithm does not get model data from which to train itself, but instead it tries to find clusters or groups inside the data that are linked together more closely than to other parts of data.

TODO: References!

Typically, when training an algorithm, some predefined portion of the data is used as training data.

TODO: Sources? Amounts?

The rest is used to validate the results so that validation data and training data do not overlap. Instead, trained algorithm is given data it has not seen before and the result it produces with it is then validated. For example, in supervised learning the key values the algorithm is trained to find out are hidden from the validation data. The resulting values produced by the algorithm are compared to those hidden values and the difference between the estimate and the real value can be used to determine how well the current algorithm compares to others.

TODO: Next one might change a bit depending on the actual results,
which are still to be tested.

However, in this study, we are going to break that rule about non-overlapping training and validation data. The reason for this is explained further in section [3.6](#))

2.2 Cloud ML platforms

TODO: Briefly about Azure, Google and AWS

Machine learning algorithms are not light to operate. Depending on the amount of data, it can be a serious

TODO: something something seriously big

Especially with online applications where real time analysis of new input data is required, cloud computing resources can make a huge difference in terms of processing speed.

TODO: referencing

Online market offers several solutions for ML computing in cloud.

TODO: open up some facts about these. BRIEFLY

Google

Amazon AWS(?)

Microsoft Azure[\[8\]](#)

2.3 Azure ML Studio

TODO: More info about Azure ML studio

Comment: some unorganized text:

Microsoft Azure offers a Machine Learning Studio environment for easy ML pipeline designing.

Explain pipeline briefly.

ML Studio gives ML designer a possibility to train algorithms and publish cloud endpoints utilizing all Azure resources connecting the power of ML to all other Azure possibilities like data storages, IoT-services and cloud computing.

With drag-and-drop pipeline designer it is easy to get started with ML programming, and visualizing the process helps understand all pipeline components and their relations to each other.

TODO: picture of azure pipeline

TODO: Intro to next subsections

Each component in pipeline can be tuned to certain extent. ML Studio has a predefined set of ready algorithms to use. In this study we focus on R-script execution component, regression algorithm components, PCA-based anomaly detection component, N-Gram feature extraction component and feature hashing component.

2.4 Regression algorithm

TODO: science and math behind regression algorithm

Regression analysis is typical approach in statistical science. It is used to find relationships with a set of variables.

2.5 PCA-based anomaly detection

TODO: Explain PCA and mention other ADA algorithms

Comment: unorganized text below:

Principal Component Analysis, or PCA, is an ML technique used to analyze data and explain the variance inside it.

Other anomaly detection methods exist, but they are not supported by ML Studio in a ready component level.

TODO: Something about Anomaly and Novelty detection differences?

2.6 N-gram features and feature hashing

TODO: cover basic n-gram features and ML connection

Comment: unorganized text below:

As stated before, features are the key elements in ML algorithm training. As textual input does not have any meaning to machines as itself, it is necessary to create a connection between words and features for algorithm. In ML training, one typical approach is to convert textual input to numerical features. For example, by creating a dictionary of words used in the input and assigning each word an identification number, we can express words as a count of certain words used. In addition, as words include meanings not only individually but also with relation to each other and in their order, we can add more information for the algorithm by creating word pairs and groups in the dictionary. These groups are referred as word grams, where **n** in n-gram refers to the maximum number of words in a group of consecutive words in the input sentence.

TODO: explain feature hashing component functionality briefly

As the number of word grams in a dictionary can increase significantly in complex input cases, it is necessary to limit the resource usage by decreasing the features analyzed. One way to do this is use feature hashing. This means that instead of pure n-gram count we use hashed value of several n-grams thus reducing the amount of features. As a drawback, the amount of information might also get reduced as the data is “compressed” but this way we can include more features for algorithm training without significant resource demands.

2.7 Data sensitivity

TODO: some basic stuff about anonymization, data sensitivity and data protection

TODO: Just practicalities, not much about theory or background

During this study, it was necessary to make sure no sensitive data was moved out of the production environment. This was mostly due to regulations described in GDPR.

Data anonymization was executed in production environment with Powershell script. Several predefined identification features were searched with regular expression (or regex), patterns and replaced with default keys.

2.8 Log data analyzing with ML

TODO: Some research should exist

TODO: In this section, we look what studies and cases of ML log data analyzing exists in the IT field. Only briefly, nothing too deep

2.9 Random delay in log event analyzing

TODO: Anything about the topic?

Random delay in input data features is not unusual aspect in time-series forecasting. Time-series in ML context refers to data features that varies over time and is usually affected by past values. As an example, ML algorithm could try to predict future weather based on measured temperature and air pressure. Both these features change over time and also affect to their own future values.

Random delay in such an example could be due to some other local or global features like

This study, however, is not time-series because the majority of the log rows are not affected by previously logged events. Random delay in this case is caused by banking clerks finding the issue and writing a technical support request. This delay can span from hours to several days depending on the weekday and time the issue occurred.

As random delay of such does not seem to be trivial to take into account with ML algorithms, a simple method to solve this was used we call “time frame compression method”. This means that in order to eliminate the effects of random delay we compress some features in certain time frame which is at least as long as the longest estimated delay. Simply put, as we count possible anomalies during one hour of log, we cannot compare this number to actual tickets received at the same hour or the next. What we can do, with time frame compression, is that we count some statistical values of anomaly estimates, for example, the mean and median values of a week, and then compare these numbers with the tickets received during the same week.

2.10 Hybrid machine learning approach in anomaly detection

Hybrid machine learning (HML) refers to an ML technique where two or more ML methods are combined to overcome the limitations of or to boost the estimation capabilities of a single method alone.[9] In this study, we combine PCA-based anomaly detection algorithm with regression algorithm in order to amplify the prediction powers of our ML algorithm when trying to determine the possible ticket count based on log events.

Hybrid machine learning is not rare technique in ML field.[10, 11, 12, 13, 14, 15, 16, 17]

TODO: Something about the HML in ADA

In order to clarify whether hybrid approach is suitable for the current study problem we will compare the results of hybrid ML technique with a single ML algorithm usage.

With ML algorithm utilizing n-gram features combined with time frame compression it is possible to get estimates about the support tickets based on the log events. It is not feasible to use anomaly detection on its own to do this as plain sum of anomalies detected is not correlating with tickets received.

We can, however, amplify our ticket estimating algorithm with anomaly value features. As we first count the anomaly numbers with anomaly detection algorithm and their calculated statistical features with another algorithm, like regression algorithm, we get more relative information to use when creating the final ticket number estimations.

Comment:

Hybrid ML as a term is used here to explain that we use two different ML algorithms in two separate phases. In first phase we try to give an anomaly certainty value for each log row using PCA-based anomaly detection component. In second phase we use this value as a feature to estimate ticket amount in time range by utilizing regression algorithm.

Dual algorithm approach should not be unusual in ML field, but how existing studies or case examples relate to our way is uncertain.

If nothing exists about the topic (at least nothing easily to be found) it should be worth to mention. But if there is a lot of case examples about this, it feels unnecessary to discuss about it in more detail.

3 Research material and methods

In the next section we explain in more detail what the data used in the study consists of and what methods were used in attempt to answer the research goals.

The data in the research is mainly made up of two parts. The most important part is, obviously, the log data produced by the numerous RPA processes. The second part complementing the study is the support ticket data written by clerks of customer banks. In order to use the data safely in the cloud environment it was necessary to sanitize the data from any sensitive information. This was done by anonymizing the log data and using only timestamps from the support tickets.

After confirming the results of anonymization, the data was preprocessed into a better form to make it more usable by algorithms. More processing was done inside the pipeline as ML Studio offered several usable components for this but main cleaning was easier and to execute in local environment. This was also done with Powershell scripting.

ML pipeline was created in Azure ML Studio and several ML algorithms were compared in order to find the most feasible for our goal in mind. ML training was organized in two different phases in order to find the relation between log anomalies and technical tickets.

Finally, the results of the trained algorithms were validated against newly acquired production data in order to estimate how well the initial goals of the study were fulfilled. These results are presented in the next section [4](#).

3.1 Support ticket data

Like all other software, RPA components fail from time to time. As described before, RPA logs are verbose making possible error identification from among it hard. Due to that, it is not feasible to create log parsers that would be able to identify critical errors from within thousands of lines of log. When critical error happens causing the RPA process to fail, the banking clerks need to finish manually the job left by the RPA robot. Every time this happens, these clerks then send a support request ticket to Samlink technical help desk and ask to fix the issue.

When clerks send the ticket to technical support a verbose description of the situation is written to help developers to identify the problem. This description often contains sensitive end customer information like bank account details and social security numbers. To avoid privacy issues when processing this data, it was decided to use only timestamps of the tickets. The resulting data was practically a list of date and time values. More about the issue from privacy point of view is described in section [3.3](#).

3.2 RPA log data

Robotic process algorithms used in Samlink are designed to ease the workload of bank clerks. RPA robots work mostly with loan applications among other routine tasks that require mostly manual labor.

TODO: Check!

Like other software, RPA also produces log data during runtime. As several RPA robots are running the amount of log produced daily is also significant.

TODO: How many RPA:s? How many customers? How much log data?

This log data is not in consistent structure being formed out of typical CSV data injected with even more inconsistent JSON data that varies in contents vastly.

TODO: refer to appendices, include examples of data

RPA log data is stored in SQL database. The database is split in live production log that is gathered for few months and then moved to archive that has several years worth of log.

TODO: check live timescale

In this study we used archived data as it was easier to acquire in one run without the need to merge different parts together. Archive also had data that was considered as sufficient amount for machine learning algorithm training.

TODO: how much really?

Comment:

Next up we'll explain some key features about the log data such as timestamps, robot names, messages etc.

The important part to keep in mind here is that message is in two forms: message and rawmessage. Without rawmessage, the data was in pure CSV-format. However, it was not certain that rawmessage would not hold usable data for ML. In the end, the usage of rawmessage was mostly skipped as it demanded too much resources in ML Studio to process. In the data preprocessing phase the rawmessage was kept along and it lead to some interesting and most likely reusable preprocessing scripts.

3.3 Data anonymization

Support ticket data privacy

Samlink handles highly sensitive banking customer data in its processes, such as personal identification numbers, home addresses, email addresses and bank account numbers. All possibly sensitive data had to be removed before data could be transferred out from production environment to cloud. Due to bureaucratic reasons, technical support tickets were under more strict policies. Because of this, they were allowed to be used in the research on condition that no business critical nor customer sensitive information was processed in the first place. Only way to assure this was to select solely timestamp fields from ticket data. Thus, no sanitation for ticket data was needed as ticket data consisted of only list of datetime values.

RPA log data sanitization

Information privacy is one of the key values in Samlink business promise as company develops high security banking applications and processes sensitive customer data. Thus, several aspects were needed to take into consideration before log data could be authorized for thesis study usage. To improve privacy, it was decided to assume that personal customer details are not critical information for ML algorithm training if goal is to find possible problems in RPA runtime and not detect individual customer related problems. This way it was not necessary to achieve adequate security by less secure and more effort consuming ways such as pseudonymization or k-anonymization, which would have also required strict inspections before data could have been approved for cloud processing.

TODO: References?

Comment:

Here we could explain more about pseudonymization and k-anonymization, but is it necessary as they were not used or considered? It would bring something more to the study, of course, but is it worth the time?

As production environment is built on Microsoft Server based solution, and because it was highly unrecommended to install additional software to the production server, data acquiring and anonymization tools were chosen based on what was already usable in the RPA production environment. Microsoft Powershell offers sufficient tools for database SQL querying and stream editing. The amount of data was significant which made straight file editing impossible due to the memory limitations. Thus, stream editing was necessary for finding and replacing sensitive information from the data.

TODO: maybe references?

Anonymization took good proportion of the time in workdays as processes were slow, the amount of data was huge and multiple re-runs were needed before the results were seemed adequate.

TODO: Appendix of the script used. Does this need more explaining?

3.4 Data formatting

At the beginning of the research, the log data from RPA was in SQL database. However, the database used was not “pure” in a way that typical relational databases are, but some columns included JSON-formatted data in them. For ML algorithms to be able to read the given data with ease this kind of impurities needed to be cleared from the data.

When feeding the log data to anomaly detection algorithm it was necessary that all the rows were as minimally unique as possible in order to use the pattern finding abilities of the algorithm. Too unique data points would have made all of them anomalies compared to each other. Thus, all unique features were stripped from the data, such as the fingerprint value that was unique for each data point.

Comment:

Also, job ID information and timestamps of the rows were removed momentarily from within the rawmessage.

Comment:

Some interesting formatting solutions were used in this phase. Rawmessage was practically created again so that it fitted in the CSV format. This should be explained in some level of detail.

TODO: check above section so it makes sense

3.5 Azure environment

TODO: This section is under construction! Here are some things we are discussing here:

Comment:

Azure resources, such as virtual networks, storage spaces, connections to ML studio etc.

More info about ML studio, computing clusters, jobs, datasets and such. What kind of resources were usable based on the issue at hand and limitations by the company (cost, basically). We discuss some things about Azure ML Studio but only in extent of what parts needed configuring. More general info about ML Studio should be in chapter 2.

3.6 Machine learning pipeline

TODO: this section is under construction

As stated in section 2.1, the approach in this thesis is, if expression is allowed, unorthodox. Typically, it is not a good idea to use same data points in ML algorithm training and validating. Acting otherwise leads to algorithm processing with same data it was trained with thus creating a situation where algorithm already knows what to do with the current datapoint. If the results were validated after this the algorithm would get unreliable score as it had the validation data already in the training phase. This could be compared to giving some right answers to students during test and scoring test results as if no help was given.

TODO: explain why we did it anyway!

Comment:

Depending on the results (still working on it...) it is possible that this unorthodox training method didn't bring much results compared to proper training pipeline.

Initial plan when starting the ML pipeline testing was to feed the log data to anomaly detection algorithm and try to get some sort of estimate of possible anomaly count. This plan had several problems.

First, as stated, logging is very abundant and several thousands of rows is logged during a single day. Some errors encountered are not critical and RPA robot is able to recover from them finalizing the initial task. This means, that errors that could be deemed anomalous may not result to a ticket in the end.

In addition, one single error case may be linked to several problems in runtime, meaning that one ticket received is, in fact, linked to multiple, dozens or even hundreds of log rows.

Two different algorithms were needed. In phase 1, algorithm defines how likely one datapoint is to be considered an anomaly. In phase 2, another algorithm aims to predict how many tickets are to be expected to receive within a time frame.

TODO: Under construction below:

Comment:

Next in this subsection, ML studio pipeline used is explained in more detail. This means pictures about the pipeline, explanations about the components used and their parameters, etc. etc.

Some of the content might be wise to move to the Results section. This needs more somewhat more thinking, also based on the results...

Memory issues and limitations

Memory is crucial in ML training as multiple steps happen and data is formatted etc. . While building ML pipeline in Azure ML studio, a memory issue emerged that affected several components and caused serious limitations in terms of usable components and data size. Due to the time limits of this study this issue was not resolved and the problem behind it was not found. As several conditions considering the environment costs were already issued by the company, the issue was declared to be linked with compute instance property limitations. However, this was not certain.

Feature format for PCA-ADA

In Azure ML Studio there is only one module selectable for anomaly detecting, the PCA-based anomaly detection module, which is explained in section 2.5. However, with textual input like logs it can be used at least in two ways. First, input data can be fed to the algorithm trainer as is, letting the PCA-based ADA component do the work without further modifying the log rows. This way, the component tries to recognize the anomalies based on all the information included in the row. Practically this means that the component processes data in textual format making each row in the input a feature as a whole to consider.

TODO: PCA-ADA should be explained in background-section

Second option is to convert the textual features into numerical n-gram features. Each word or n-gram is now a number of said instances found on the row being processed, and each row can be presented as a sequence of numbers indicating the number of those features.

N-grams can in addition have a weight based on the frequency they appear in the entire data. Different weights usable in Azure ML component are listed below.

1. Binary Weight
2. TF Weight
3. IDF Weight
4. TF-IDF Weight

TODO: Explain different weights and open up more Azure ML Studio
PCA-component

Anomaly probability

Comment:

Here we explain what PCA-component outputs and how the result is used in the pipeline.

Statistical features

Comment:

This part explains some statistical features used when starting regression training in hybrid ML phase 2. These features consist of both log row amounts and PCA-component output values. All these values are grouped together in timescales.

Regression based estimating

Comment:

Here is more information about different regression algorithms used in ML pipeline. Some basic information about all of them is given so the results are understandable by reader in the Result section.

1. Linear regression
2. Decision forest regression
3. etc.
4. etc.

4 Results

We start this section by looking at the Azure resources needed for ML training both on common Azure environment and ML Studio. Next we analyze the results of different ML algorithms and pipelines.

4.1 Azure and Azure ML Studio

Azure resources

Comment:

what resources was needed inside Azure?

virtual machines etc.

Unless these were discussed in previous sections.

Azure ML Studio components

Comment:

clusters and data

Memory problems

During the initial pipeline runs the execution came to an abrupt stop and Azure notified about memory issues. These problems were linked to the data amount which had to be reduced to 600 megabytes before any pipeline could be finished using the data. This reduction was against the initial goal where preferably all the data could have been used.

Considerable amount of time was used to fix or avoid this issue but nothing clear was found that would explain the error received. While working with the issue it was also noted that data needed more cleaning in order to ease the preprocessing phase as described with more detail in section 3.3 Thus, the data had to be imported from log archive and anonymized once more.

Two choices was possible to take:

1. Continue working with full data and attempting to fix the memory issue by consulting Azure experts
2. Trim the data to reduce the data size by declaring info-type log messages as unnecessary and working with vastly diminished data until the memory issue would be solved one way or another

To advance the study more efficiently it was decided to trim info-type log messages from data hence reducing the data amount considerably. Final data included 8.6 million log rows which was about 10% of the original data size.

Even with this data size, some Azure ML components faced this memory issue and forced us to choose such components that were able to handle these data amounts.

In addition to the data we used to train the algorithm, we needed to set up computing instance in Azure ML studio. Some predefined resource limitations affected the computing instance choosing and the memory issue encouraged us to pick memory prioritized instances. Single computing instance did not work, but we needed to choose a computing cluster instead to be able to run ML training pipeline.

4.2 ML training and validation

Comment:

anomaly detection

N-Gram Feature extracting

Some regression algorithm to predict event count. Poisson only for poisson distributed data, probably not usable here

two-class classification support vector machine etc

random forest regression is probably best based on the initial results

preformatting data!

1. remove fingerprint etc unique values from raw message
2. calculate anomaly probability per line

! CAN WE COMBINE THIS PER JOB-ID?

3. create new table consisting:
 - (a) amount of rows per timeframe
 - (b) amount of unique job ID's per said timeframe
 - (c) anomaly probability value (median, mean etc)
 - (d) efecte tickets received in said timeframe

<Integrating with timestamps>

To avoid the problem with random delays between log rows and technical ticket timestamps, log rows were grouped by time stamp into certain time frame groups.

Comment: Next is some initial values from pipeline runs. Subjected to change!

N-Gram feature extraction

Using Decision forest regression algorithm.

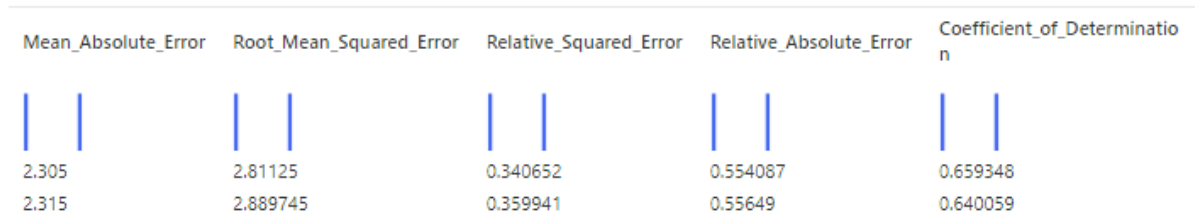


Figure 1: Message with N-Gram feature extraction using unconventional training method in phase 1, Decision forest regression in phase 2, compared to method without anomaly values.

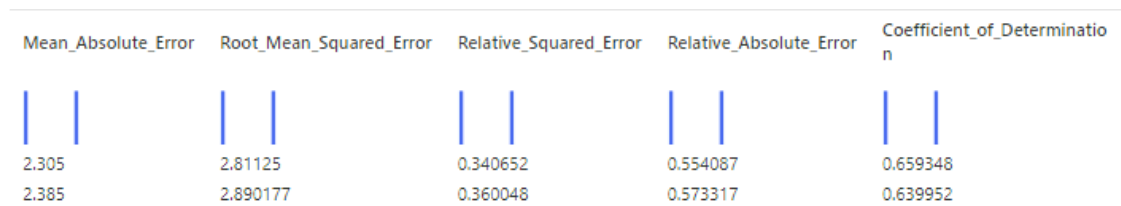


Figure 2: Message with N-Gram feature extraction using proper training method in phase 1, Decision forest regression in phase 2, compared to method without anomaly values.

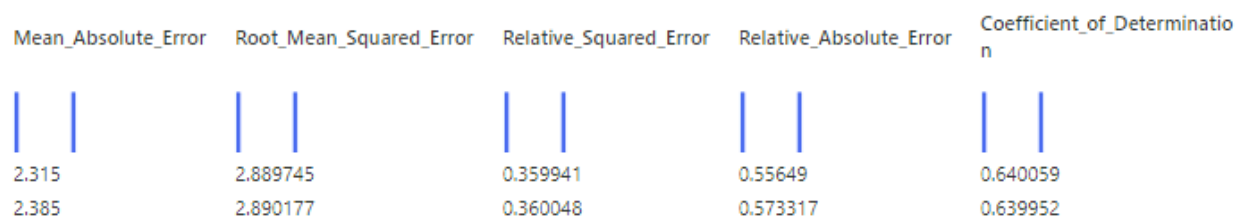


Figure 3: Message with N-Gram feature extraction using Decision forest regression in phase 2, comparing unconventional training to proper training.

Anomaly detection with pure textual data

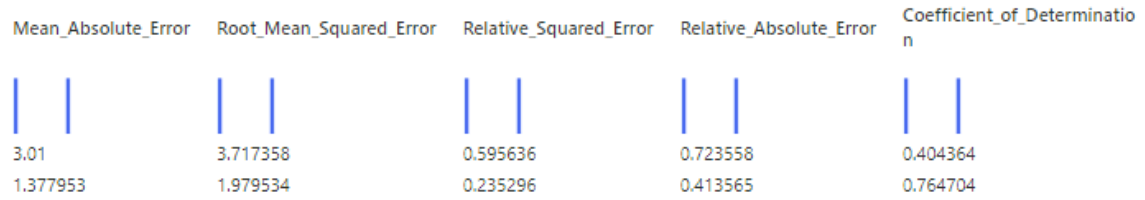


Figure 4: Message fed pure to ADA-component, using unconventional training method in phase 1, Decision forest regression in phase 2, compared to method without anomaly values.

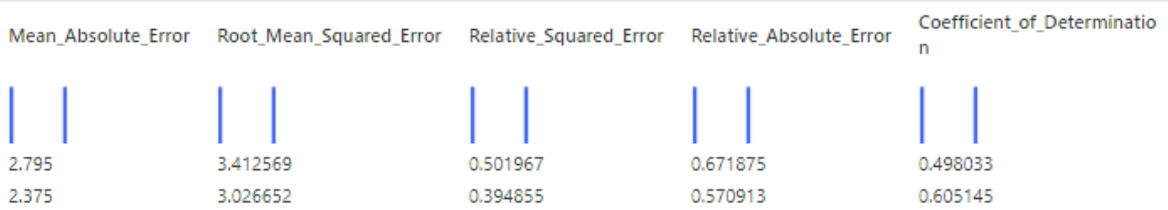


Figure 5: Message fed pure to ADA-component, using proper training method in phase 1, Decision forest regression in phase 2, compared to method without anomaly values.

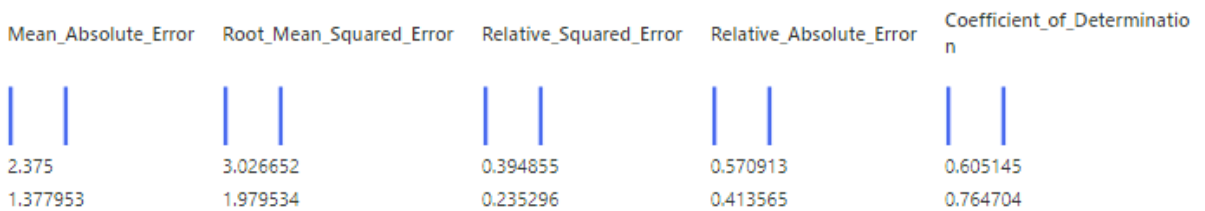


Figure 6: Message fed pure to ADA-component, using Decision forest regression in phase 2, comparing unconventional training to proper training.

5 Summary

TODO: Very under construction...

<Sum up here what we did and why>

5.1 Discussion

<Here some thinking what should have been improved>

Integrating with real time logging?

Data formatting

The most time-consuming tasks in the study was without a doubt the anonymization and preformatting of the data. Although sensitive information may sometimes be crucial in error fixing as problems may consider just one client, it is necessary that the data sanitation is possible to do in order to use the data in less secure environment. By preformatting the data in such way that all different personal information types do not differ between use cases.

Possible ML methods

Some sort of time delay forecasting?[18] Could we estimate the number of tickets based on some log metrics in time frame?

Memory error fixing would have given more options.

TODO: Something else specifically needed?

References

- [1] P. K. Donepudi, "Machine learning and artificial intelligence in banking," *Engineering International*, vol. 5, no. 2, pp. 83–86, 2017.
- [2] W. M. Van der Aalst, M. Bichler, and A. Heinzl, "Robotic process automation," pp. 269–272, 2018.
- [3] A. DeLaRosa, "Log monitoring: not the ugly sister," *Pandora FMS*, 2018. [Online]. Available: <https://web.archive.org/web/20210901031146/https://pandorafms.com/blog/log-monitoring/>
- [4] W. K. Ho, B.-S. Tang, and S. W. Wong, "Predicting property prices with machine learning algorithms," *Journal of Property Research*, vol. 38, no. 1, pp. 48–70, 2021. [Online]. Available: <https://doi.org/10.1080/09599916.2020.1832558>
- [5] K. Ghanem, F. J. Aparicio-Navarro, K. G. Kyriakopoulos, S. Lambotharan, and J. A. Chambers, "Support vector machine for network intrusion and cyber-attack detection," in *2017 Sensor Signal Processing for Defence Conference (SSPD)*, 2017, pp. 1–5.
- [6] M. I. Jordan and T. M. Mitchell, "Machine learning: Trends, perspectives, and prospects," *Science*, vol. 349, no. 6245, pp. 255–260, 2015.
- [7] "Definition of algorithm," <https://web.archive.org/web/20220510183749/https://www.merriam-webster.com/dictionary/algorithm>, accessed: 2022-05-19.
- [8] "Comparing machine learning as a service: Amazon, microsoft azure, google cloud ai, ibm watson," <https://www.altexsoft.com/blog/datascience/comparing-machine-learning-as-a-service-amazon-microsoft-azure-google-cloud-ai-ibm-watson/>, accessed: 2022-04-04.
- [9] F. Anifowose, "Hybrid machine learning explained in nontechnical terms," *JPT*, 2020. [Online]. Available: <https://web.archive.org/save/https://jpt.spe.org/hybrid-machine-learning-explained-nontechnical-terms>
- [10] T. Shon and J. Moon, "A hybrid machine learning approach to network anomaly detection," *Information Sciences*, vol. 177, no. 18, pp. 3799–3821, 2007.
- [11] C.-F. Tsai and M.-L. Chen, "Credit rating by hybrid machine learning techniques," *Applied soft computing*, vol. 10, no. 2, pp. 374–380, 2010.
- [12] S. Mohan, C. Thirumalai, and G. Srivastava, "Effective heart disease prediction using hybrid machine learning techniques," *IEEE access*, vol. 7, pp. 81 542–81 554, 2019.
- [13] N.-C. Hsieh, "Hybrid mining approach in the design of credit scoring models," *Expert Systems with Applications*, vol. 28, no. 4, pp. 655–665, 2005.

- [14] A. Jain and A. M. Kumar, “Hybrid neural network models for hydrologic time series forecasting,” *Applied Soft Computing*, vol. 7, no. 2, pp. 585–592, 2007.
- [15] H.-j. Kim and K.-s. Shin, “A hybrid approach based on neural networks and genetic algorithms for detecting temporal patterns in stock markets,” *Applied Soft Computing*, vol. 7, no. 2, pp. 569–576, 2007.
- [16] T.-S. Lee, C.-C. Chiu, C.-J. Lu, and I.-F. Chen, “Credit scoring using the hybrid neural discriminant technique,” *Expert Systems with applications*, vol. 23, no. 3, pp. 245–254, 2002.
- [17] R. Malhotra and D. Malhotra, “Differentiating between good credits and bad credits using neuro-fuzzy systems,” *European journal of operational research*, vol. 136, no. 1, pp. 190–211, 2002.
- [18] G. H. Erharter and T. Marcher, “On the pointlessness of machine learning based time delayed prediction of tbm operational data,” *Automation in Construction*, vol. 121, p. 103443, 2021.

A Esimerkki liitteestä