

DATA DRIVEN BUSINESS DEVELOPMENT - LOAN DATA

EIRIK DUESUND HELLAND



Michelle Wiggins | <https://www.kwikcashionline.com/quick-online-loans/> | <https://creativecommons.org/licenses/by-sa/4.0>

EXECUTIVE SUMMARY

A 35-year-old man arrives at a bank downtown. He lines up in the queue stretching through the reception and sighs at the fading opportunity of a swift visit when he sees the slow banker. She is counting the coins of an old lady at the counter. He thinks of his wife and children waiting at home and wishes that this process consumed less time and was not so nerve-wracking. He is there to apply for a loan.

Almost everyone must apply for a loan at some point. Billions of applications have been processed manually until now, consuming lots of time and resources. Is there a way to free up these resources, streamline the application process and create a better experience for the applicants?

This project has created a model for predicting the loan status of loan applications, using loan data. This model has been integrated into a user interface using streamlit to create a user-friendly experience, which makes the loan application process much faster.

The project demonstrates the automation of the loan application process and the added business value it provides. To implement this project in real-life-scenarios, the model is dependent on high quantities of quality data. If this is provided, then the time of the manual processing of loan applications might come to an end.

LIST OF FIGURES

Figure 1 - Molka, T., Gilani, W. & Zeng, X.-J. (2013). Dotted Chart and Control-Flow Analysis for a Loan Application Process. Business Process Management Workshops, Tallinn, Estonia, pp. 219-220. Berlin, Heidelberg: Springer.....	8
Figure 2 – System Architecture Flow Chart.....	10
Figure 3 – System Architecture part 1.....	11
Figure 4 – System Architecture part 2.....	11
Figure 5 – Features and Missing Values.	12
Figure 6 – Histogram showing the Data’s Distribution.	12
Figure 7 – Boxplot of Data.....	13
Figure 8 – Numerical Variables plotted vs Loan Status.....	13
Figure 9 – Categorical Variables plotted vs Loan Status.	13
Figure 10 – Correlation Matrix.	13
Figure 11 – Pre-Processing script.	14
Figure 12 – Model training script.	15
Figure 13 – Confusion Matrix.	16
Figure 14 – Output from console after modelling.	16
Figure 15 – Imports for the application script.....	16
Figure 16 – The rest of the application script.	17
Figure 18 – Declined Loan in Streamlit.	18
Figure 17 – Approved Loan in Streamlit.....	18

CONTENT

Introduction	5
Background	5
Problem statement	5
Goals and Objectives.....	5
Limitations.....	6
Theory and key concepts	7
Data driven business development.....	7
Robotic process automation	7
Machine learning	8
Loan Application process	8
State-of-the-art analysis.....	9
Develeopment.....	10
Design description.....	10
Concept Overview	10
System Architecture.....	10
Exploring the data	12
system development.....	14
PreProcessing	14
Modelling	15
Application	16
Evaluation and results.....	19
Data	19
Model	19
Implementation	19
Business value.....	20
Conclusion and Recommendations.....	20
Further work	20
Attachments.....	20
Source code.....	20
References	21

INTRODUCTION

BACKGROUND

This paper is written as a study for the course TIN200 at NMBU. The idea and data material were provided as course material.

Data driven business development is currently a hot topic. The advances in programming and artificial intelligence have pushed us onto another step of the automation process that the industrial revolution started. We are now automating the previously labour-intensive task in banking, known as the loan application process.

This paper will address the problem of machines making intelligent, fair decisions that can be compared to decisions made by human loan officer. It will also address the useability of this system by non-coders, and the factors that will make this process better, such as good data and models.

This paper will streamline the process of loan applications and create an effective and robust solution, without the need of human interference.

PROBLEM STATEMENT

Applying for a loan is a tedious process that demands that the applicants are interviewed by a loan officer. This takes up unnecessary time and resources of both the customers and the bank. Is applied machine learning a way for the applicants to get instant response on their loan applications, without having to physically interact with the bank? Will this yield a positive result for the bank, freeing resources and doing as good of a job as the loan officers?

GOALS AND OBJECTIVES

The overall goal of this report is to explore the possibility of automating the approval process of loan applications using machine learning.

The following objectives have been set to achieve this:

- Research
 - Evaluate key concepts to be included in the theoretical framework of this paper.
 - Chart the standard application process of loans.
 - Research the existing field of Robot Process Automation.
 - Identify state-of-art solutions in Robot Process Automation. Analyse their evaluations to find elements of interest for this project.
 - Explore similar data sets to the one used in this project to consider further implementation of the model.
 - Investigate the business value of the development of similar systems.
- Development
 - Collect data.
 - Explore the data set.
 - Structure the RPA system.
 - Pre-process data.

- Create model.
 - Create app for automatic processing of new loan data.
 - Deploy app.
- Evaluation
 - Evaluate data used for training and testing. Evaluate the similarity to other loan data.
 - Evaluate the model for predicting approval.
 - Evaluate the app and its implementations.
 - Evaluate the business value of this system.
 - Compare machine learning vs programming.
- Recommendations and further work

LIMITATIONS

This paper is limited by the data used and the short project scope. The principal of garbage in, garbage out takes full effect in every machine learning project. A model can only be as good as the quality of the data it uses to train. The project spans over three weeks and this will also impact decisions regarding improving the model, and the paper in general.

THEORY AND KEY CONCEPTS

This chapter presents data as a resource, automation in business processes, machine learning and the loan application process. The goal is to show the value that must be extracted from data to streamline business processes and increase efficiency and free up resources.

DATA DRIVEN BUSINESS DEVELOPMENT

The Economist published an article named “The world’s most valuable resource is no longer oil, but data”, informing us that the five most valuable listed companies in the world are data companies. Companies such as Google and Facebook have a tremendous value, but their most popular services are free (The world’s most valuable resource is no longer oil, but data, 2017). What is their value?

Their value lies in the data they accumulate about you, the consumer. Google and Facebook, revenue-wise, are mainly advertisement companies, and profits from selling ads to other companies (Sweney & Canon, 2021). Their effectiveness as advertisement companies comes from their ability to tailor ads to different users on their platform, therefore increasing the value of each advertisement. This shows that there is value in data.

Other businesses can take part of this discovered value, by utilizing their data to its fullest potential. Data gathered in a business can be used internally to help the business to separate themselves in a competitive ecosystem, providing better products and services than their competitors. Data can also be used commercially by selling it to external businesses (Werger et al., 2020).

Data and physical assets have a lot in common. It can be bought and sold. It makes us more efficient and able to win, and keep, customers. Having it, without using it is a waste. It degrades over time and must be maintained. The big differences are that it does not become less when used and it is easily made accessible to everyone (Treder, 2019).

Data should be treated as every other asset. Give it a price, inventory it, maintain it, refine it and increase its value (Treder, 2019).

ROBOTIC PROCESS AUTOMATION

“RPA is the technological imitation of a human worker with the goal of automating structured tasks in a fast and cost efficient manner”(Aguirre & Rodriguez, 2017). RPA is not a physical robot, but a software created to do repetitive operational processes that usually are performed by humans. RPA can be used to automate processes based on structured data and deterministic outcomes. One of the huge advantages of RPA is that it can be implemented on top of already existing systems, across platforms. They are also made to not require programming skills to use them (Aguirre & Rodriguez, 2017).

Some benefits of RPA are as following: Reduced costs by automating processes, increasing productivity, and needing less employees. Better customer experience by freeing up resources, giving more time to focus on customers. Lowering risk by removing the highest cause of errors, the humans. RPA can leverage your existing systems, so you don’t need to replace your infrastructure (Mitra, 2018).

MACHINE LEARNING

Machine learning is the application of algorithms that makes sense of data. By using self-learning algorithms, we can turn data into something much more valuable; knowledge, by spotting patterns and making predictions. Instead of relying on humans to make complicated models by making rules, machine learning offers a faster way of capturing knowledge and provides you with the ability to make data-driven choices (Raschka & Mijalili, 2019).

There are three different kind of machine learning. Supervised learning uses labelled data to train up a model, giving direct feedback to make predictions. Unsupervised learning has no labels or feedback, but it works by finding hidden structure in the data. Reinforcement learning is a decision process that works by rewarding wanted behaviour. Reinforcement learning knows nothing to start with, so decisions are random, but every time it does something right, its rewarded, and will therefore do it similarly the next time (Raschka & Mijalili, 2019).

LOAN APPLICATION PROCESS

The loan application summarized by a control-flow analysis:

- Application is submitted, but not necessarily finished.
- All applications are either internally pre-accepted, declined, cancelled by applicant, or contacted by call.
- Now that all applications are completed by call, the same decision process is initiated, and the application is finalized.
- Now remains approval, registration, and activation.

After this application process, the loan is offered, and a new process starts around the applicants answer (Molka et al., 2013).

The Norwegian newspaper “Finansavisen” published an article last year about the long processing times of loan applications. You could expect to wait up to three weeks for processing in some Norwegian banks, at the time (Parr, 2020).

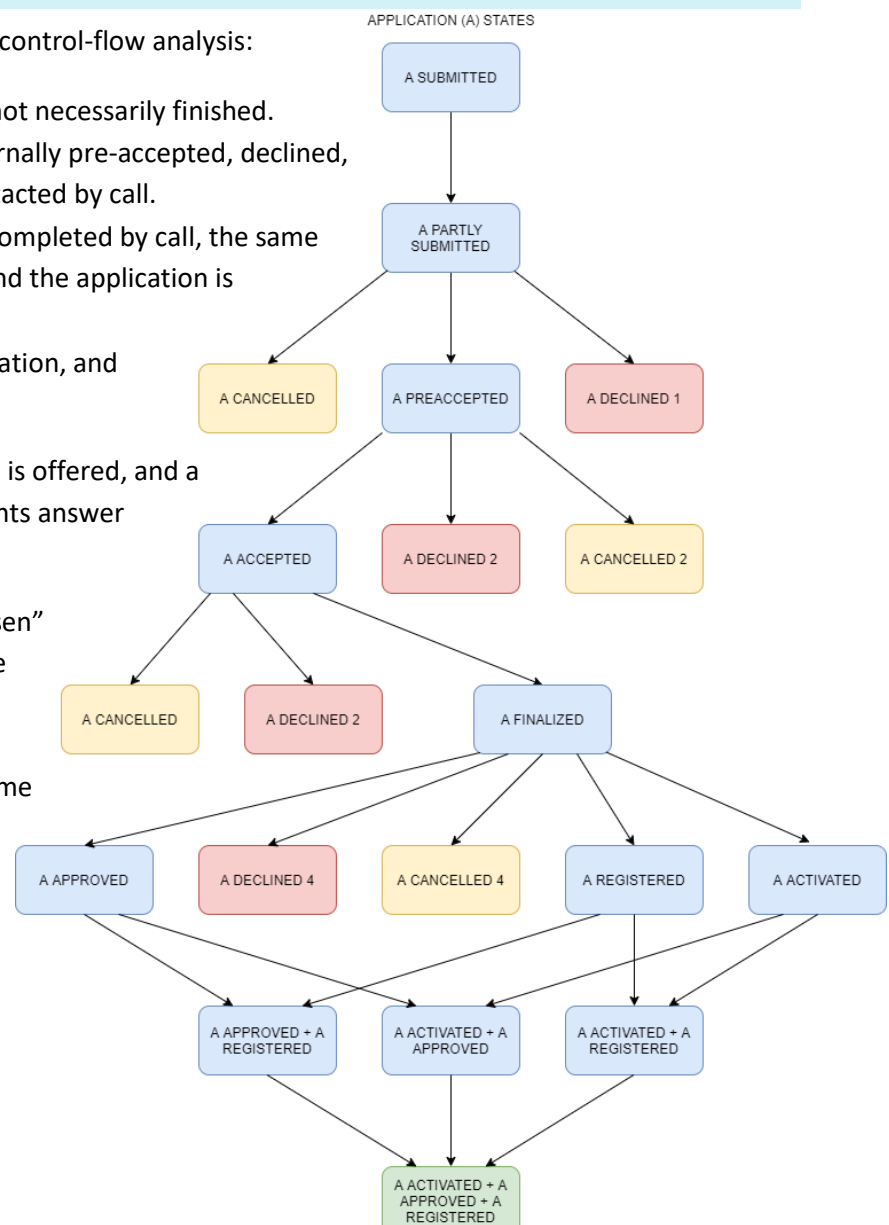


Figure 1 - Molka, T., Gilani, W. & Zeng, X.-J. (2013). Dotted Chart and Control-Flow Analysis for a Loan Application Process. *Business Process Management Workshops*, Tallinn, Estonia, pp. 219-220. Berlin, Heidelberg: Springer.

STATE-OF-THE-ART ANALYSIS

Some of the recent papers from the Business Process Management conference in 2020 gives an insight to the relevant state of art of the field of study RPA. The proceeding includes papers such as “A Conversational Digital Assistant for Intelligent Process Automation”, “How to Trust a Bot: An RPA User Perspective” and “From Robotic Process Automation to Intelligent Process Automation” (*Business Process Management: Blockchain and Robotic Process Automation Forum*, 2020).

The papers explore the use of machine learning combined with RPA interfaces to streamline the experiences of business processes. An example is how “A Conversational Digital Assistant for Intelligent Process Automation” solves the inexperience with the use an RPA with digital assistant, much like the ones we have at home, or in our phones. An example used is the simplification of the loan process. The assistant can help a loan officer without experience with machine learning, to automate the process of approving a loan, just by telling the digital assistant some key information (Rizk et al., 2020).

DEVELOPMENT

This chapter contains the development of the automation process of loan applications.

DESIGN DESCRIPTION

CONCEPT OVERVIEW

CONCEPT

An application for instant response on loan applications.

DETAILS

This app will use a model created using machine learning to decide the outcome of loan applications. This will streamline the whole process and be beneficial for both the bank and the customers.

BENEFITS

- Freeing resources for the bank by streamlining a time-consuming process.
- Giving the customers instant feedback on their loan applications.

ISSUES

- Hard to understand how the model decides the outcome.
- Lack of customer service.

SYSTEM ARCHITECTURE

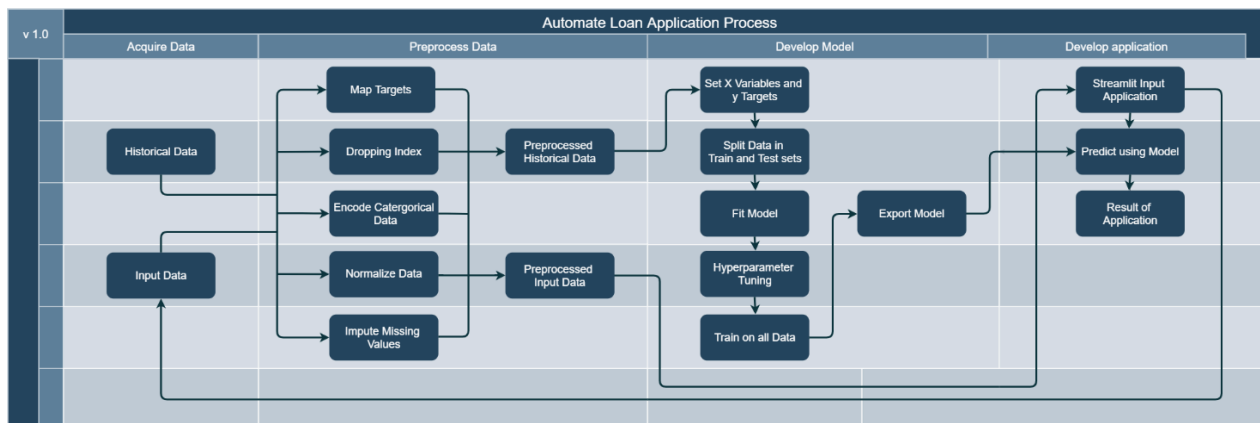


Figure 2 – System Architecture Flow Chart.

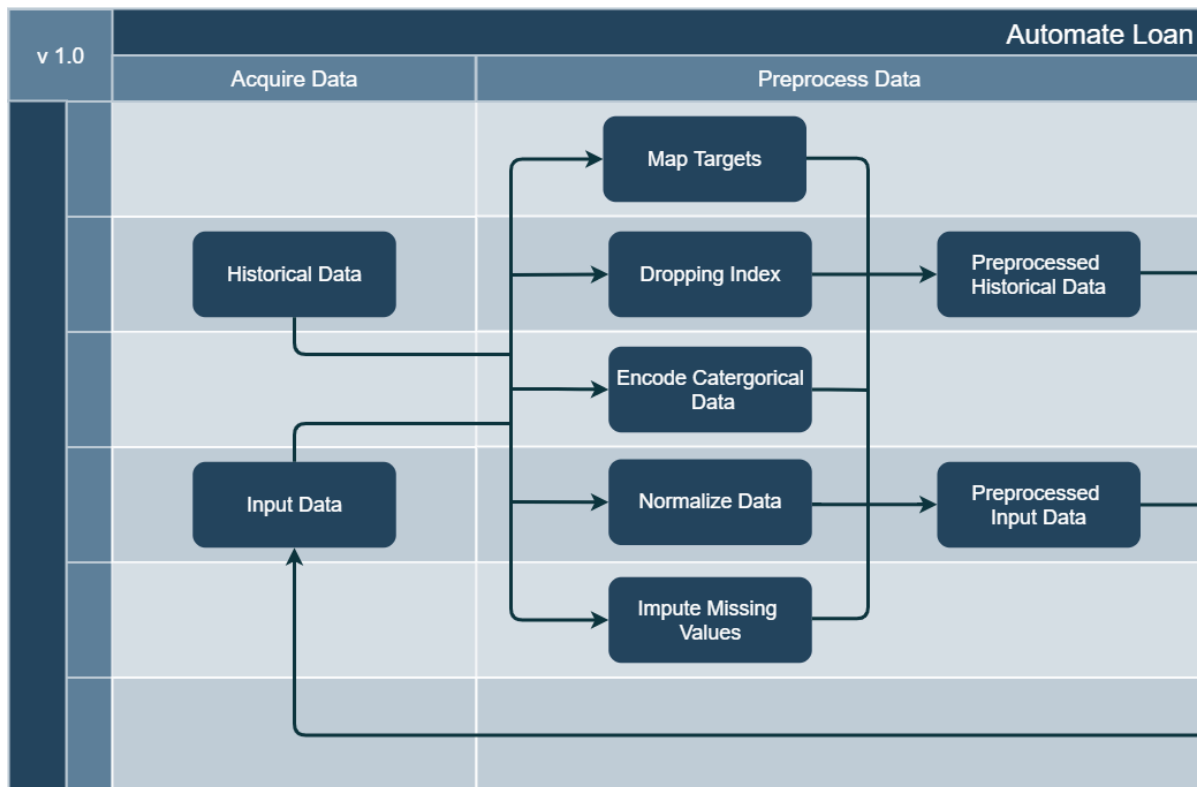


Figure 3 – System Architecture part 1.

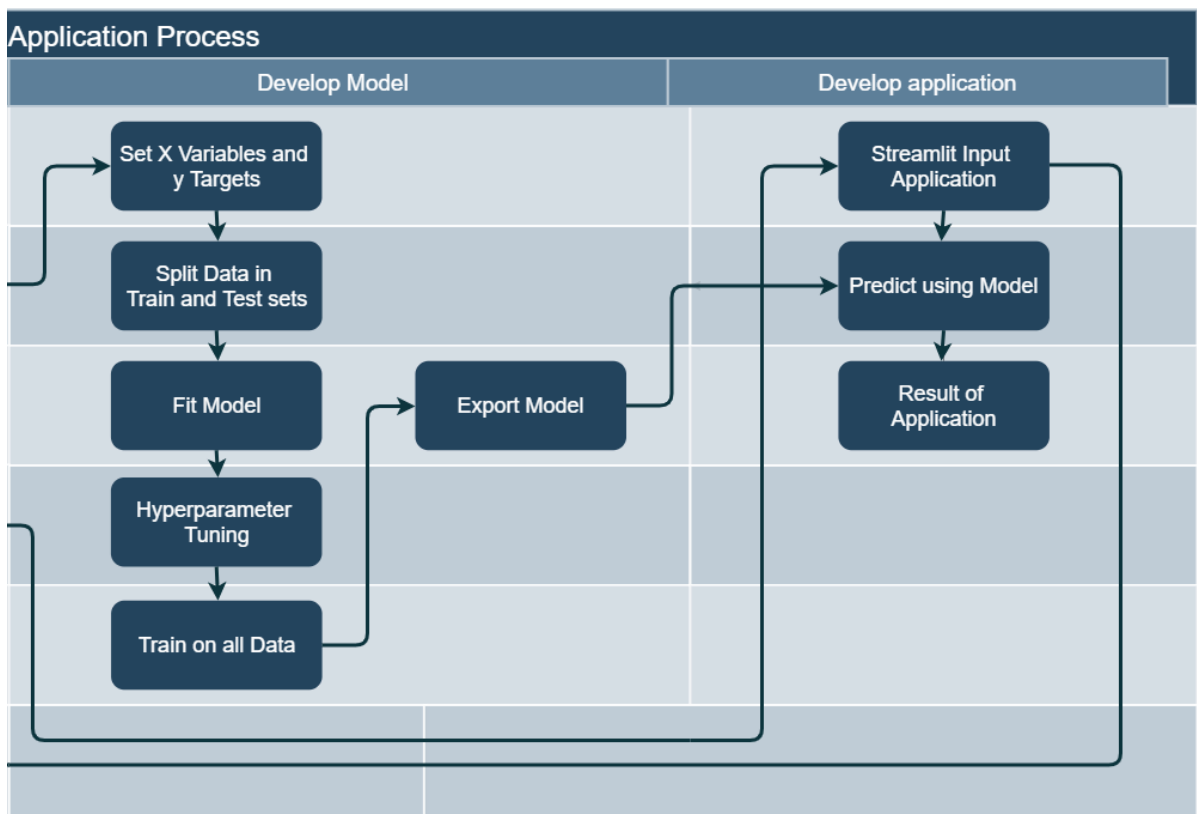


Figure 4 – System Architecture part 2.

EXPLORING THE DATA

The data used in this project was handed out as part of the TIN200 project material. It contains the information of 614 loan applicants and their loans, and their loans status.

“data_exploration.py” is a python script, found in the source code, that explores the attributes of the data. The script uses pandas, matplotlib and seaborn to display different characteristics such as missing values, distributions, and correlation.

Figure 5 shows the features the loan dataset contains. It contains information about the applicant’s gender, marriage status, number of dependents, education, employee status, income, co-applicant’s income, amount of loan applied for, terms of payback in months, if they have credit history and their property area.

Figure 5 also shows the number of missing values in the dataset. There are 22 missing values in the LoanAmount column, and 14 missing values in the Loan_Amount_Term column, in the training dataset.

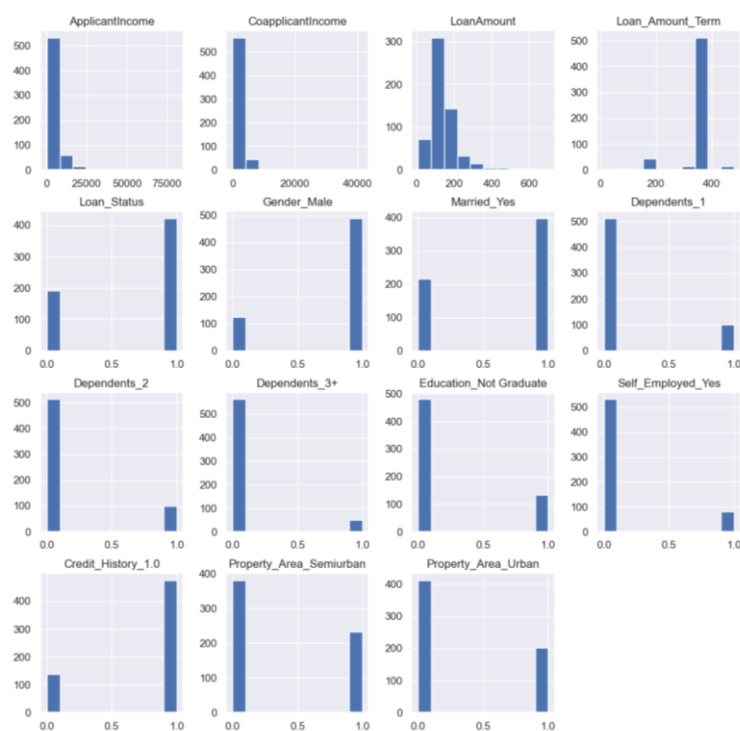


Figure 6 – Histogram showing the Data’s Distribution.

Figure 6 shows the distribution of all the data. We can see that most of the columns are categorical and have only two different values, 0 or 1.

```
Shape of data: (614, 13)

Features:
Loan_ID
Gender
Married
Dependents
Education
Self_Employed
ApplicantIncome
CoapplicantIncome
LoanAmount
Loan_Amount_Term
Credit_History
Property_Area

Target:
Loan_status: ['Y' 'N']
NaN in train:

Loan_ID          0
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount       22
Loan_Amount_Term 14
Loan_Status      0
Gender_Male      0
Married_Yes      0
Dependents_1     0
Dependents_2     0
Dependents_3+    0
Education_Not_Graduate 0
Self_Employed_Yes 0
Credit_History_1.0 0
Property_Area_Semiurban 0
Property_Area_Urban 0
dtype: int64

NaN in test:

Loan_ID          0
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount       5
Loan_Amount_Term 6
Gender_Male      0
Married_Yes      0
Dependents_1     0
Dependents_2     0
Dependents_3+    0
Education_Not_Graduate 0
Self_Employed_Yes 0
Credit_History_1.0 0
Property_Area_Semiurban 0
Property_Area_Urban 0
```

Figure 5 – Features and Missing Values.

Error! Reference source not found. also shows the distribution of the variables in the dataset. The boxplot also gives us an idea of outliers in the dataset.

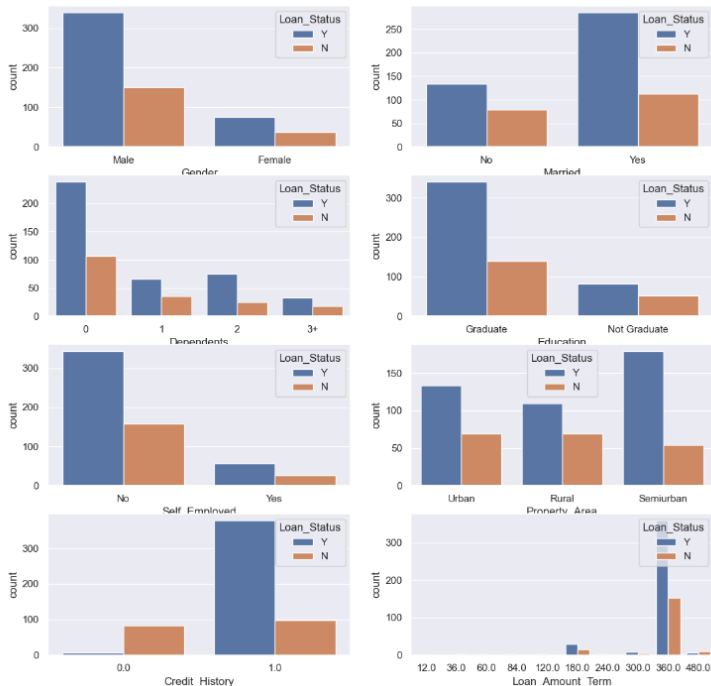


Figure 8 – Numerical Variables plotted vs Loan Status.

Figure 8 shows the distribution between the loan

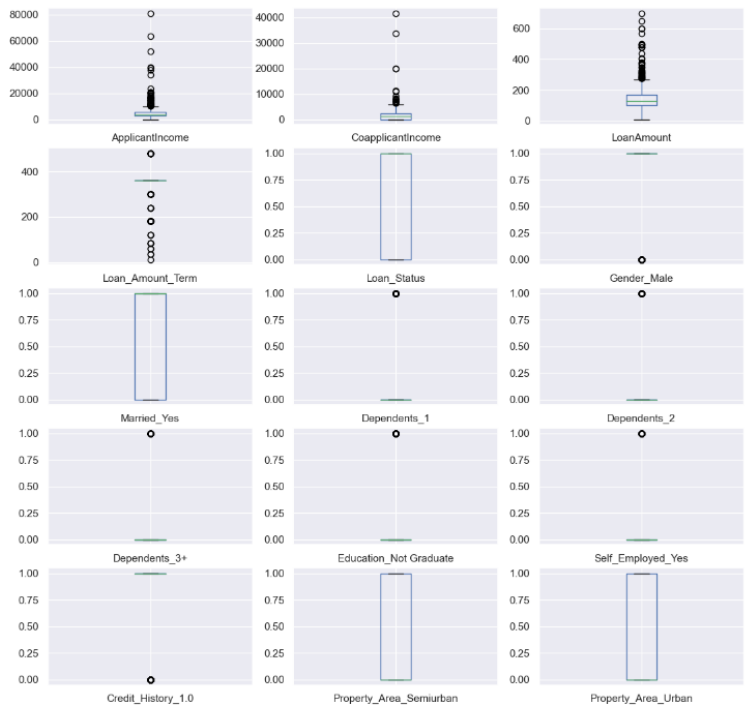
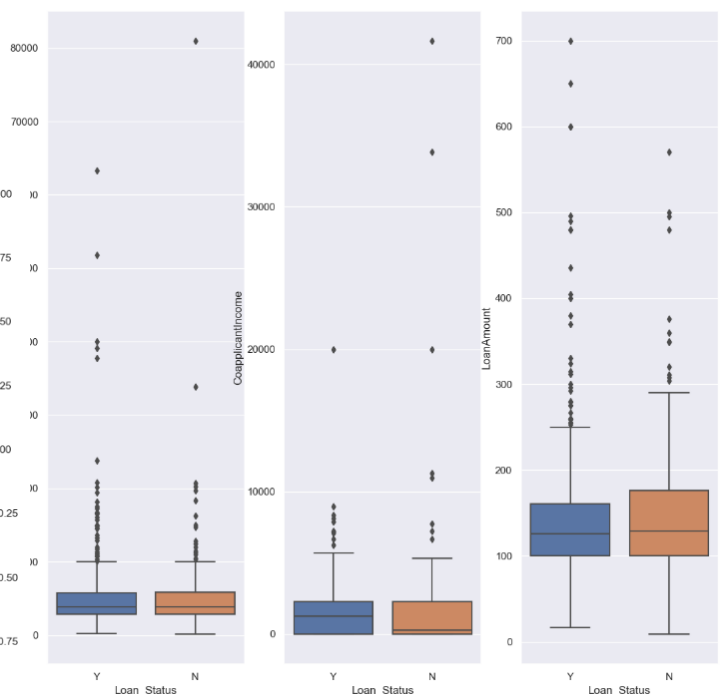


Figure 7 – Boxplot of Data.

Figure 9 shows the distribution between the loan status of each of the categorical variables.



status of each of the numerical variables.

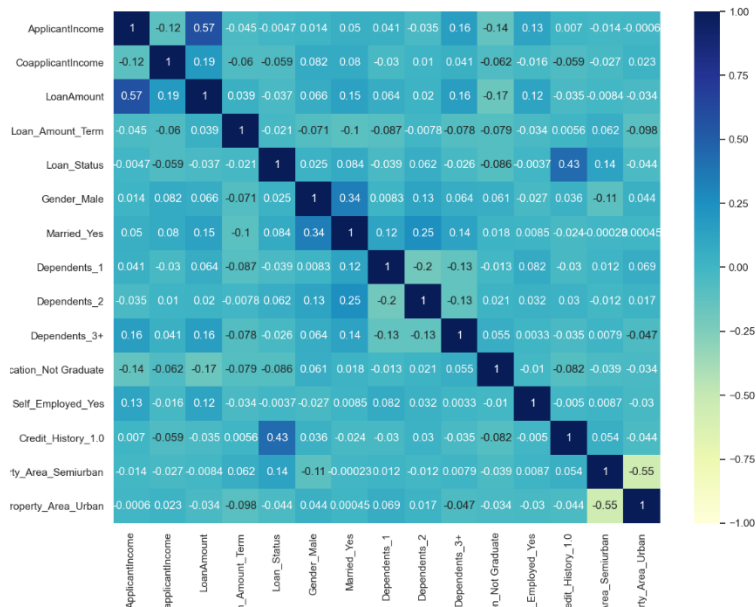


Figure 10 shows the correlation between all the variables.

SYSTEM DEVELOPMENT

PREPROCESSING

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Mon May 31 21:41:36 2021
4
5  @author: Eirik
6  """
7
8  import pandas as pd
9  import numpy as np
10
11  from sklearn.impute import SimpleImputer
12  from sklearn.preprocessing import MinMaxScaler
13
14
15  # Reading data and making dataframe
16  data = pd.read_csv('train.csv')
17  test_data = pd.read_csv('test.csv')
18
19  # Making dataframes
20  df = pd.DataFrame(data)
21  test_df = pd.DataFrame(test_data)
22
23  #Target mapping - N: 0, Y = 1
24  target_mapping = {target: idx for idx, target in
25                    enumerate(np.unique(df['Loan_Status']))}
26  df['Loan_Status'] = df['Loan_Status'].map(target_mapping)
27
28  # Dropping Loan_ID
29  df = df.drop(['Loan_ID'], axis=1)
30  test_df = test_df.drop(['Loan_ID'], axis=1)
31
32  # Create dummies
33  df_encoded = pd.get_dummies(data=df, drop_first=True,
34                             columns=['Gender', 'Married', 'Dependents',
35                                     'Education', 'Self_Employed',
36                                     'Property_Area'])
37  test_df_encoded = pd.get_dummies(data=test_df, drop_first=True,
38                                  columns=['Gender', 'Married', 'Dependents',
39                                          'Education', 'Self_Employed',
40                                          'Property_Area'])
41
42  # Normalizing data using MinMaxScaler
43  scaler = MinMaxScaler()
44  df_encoded = pd.DataFrame(scaler.fit_transform(df_encoded),
45                           columns = df_encoded.columns)
46  test_df_encoded = pd.DataFrame(scaler.fit_transform(test_df_encoded),
47                                columns = test_df_encoded.columns)
48
49  # Filling NaN data/Missing data with Simple Imputer
50  imputer = SimpleImputer(strategy='median')
51  df_encoded = pd.DataFrame(imputer.fit_transform(df_encoded), columns=df_encoded.columns)
52  test_df_encoded = pd.DataFrame(imputer.fit_transform(test_df_encoded), columns=test_df_encoded.columns)
53
54  # Save dataframes to csv
55  df_encoded.to_csv('preprocessed_train.csv', encoding='utf-8', index=None)
56  test_df_encoded.to_csv('preprocessed_test.csv', encoding='utf-8', index=None)

```

Figure 11 – Pre-Processing script.

Figure 11 shows how the data is imported from csv using Pandas and pre-processed. The target, “Loan Status”, which is either yes or no (Y or N), is mapped to 0 for no, and 1 for yes. The first column, “Loan_ID”, is removed because it is obviously not important for developing an accurate model. The categorical variables are encoded to multiple binary variables instead because the model does not know the difference between

categorical variables. The data is normalized using MinMaxScaler from the scikit-learn library. The missing values are imputed over the median of the other values in each column. The newly pre-processed data is saved as a csv file to be accessed by the model.

MODELLING

- Pre-processed data is imported.
- Data is split into y, targets, and X, features.
- The data is split random into 2/3 training data and 1/3 testing data. The testing data is later used to validate the model.
- XGBoostClassifier, or Extreme Gradient Boosting Classifier, is used to make the model.
- The model is fitted with the training data from the earlier split.
- The hyperparameters of the model is tuned using RandomizedSearchCV.
- A confusion matrix is created to show the accuracy of the model in an understandable way.
- The model is trained on all the data, without doing a train-test-split.
- The model is exported.

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Wed Jun  2 12:41:44 2021
4
5  @author: Eirik
6  """
7
8  import pandas as pd
9  import matplotlib.pyplot as plt
10 import pickle
11 from scipy.stats import uniform, randint
12
13 from sklearn.model_selection import train_test_split
14 from sklearn.model_selection import RandomizedSearchCV
15 from sklearn.metrics import confusion_matrix
16 import xgboost as xgb
17
18 # Reading data and making dataframe
19 df = pd.read_csv('preprocessed_train.csv')
20
21 # Setting the y values (target)
22 y = df['Loan_Status'].values
23
24 # Setting up the x values (features)
25 X = df.drop(['Loan_Status'], axis=1).values
26
27 # Split data into training and test data
28 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
29                                                    random_state=1, stratify=y)
30 model = xgb.XGBClassifier(objective="binary:logistic", random_state=1,
31                           use_label_encoder=False)
32 model.fit(X_train, y_train)
33
34 # Hyperparameter tuning
35 param_grid = {
36     "colsample_bytree": uniform(0.7, 0.3),
37     "gamma": uniform(0, 0.5),
38     "learning_rate": uniform(0.03, 0.3),
39     "max_depth": randint(2, 6),
40     "n_estimators": randint(100, 150),
41     "subsample": uniform(0.6, 0.4)
42 }
43
44 rs_model = RandomizedSearchCV(estimator=model,
45                               param_distributions=param_grid,
46                               scoring='roc_auc',
47                               cv=10,
48                               n_jobs=-1,
49                               n_iter=200,
50                               verbose=1)
51
52 rs_model = rs_model.fit(X_train, y_train)
53 print(rs_model.best_params_)
54 print('Test accuracy: %.3f' % rs_model.score(X_test, y_test))
55
56 # Creating confusion matrix
57 y_pred = rs_model.predict(X_test)
58 confmat = confusion_matrix(y_true=y_test, y_pred=y_pred)
59
60 fig, ax = plt.subplots(figsize=(2.5, 2.5))
61 ax.matshow(confmat, cmap=plt.cm.Blues, alpha=0.3)
62 for i in range(confmat.shape[0]):
63     for j in range(confmat.shape[1]):
64         ax.text(x=j, y=i, s=confmat[i, j], va='center', ha='center')
65
66 plt.xlabel('Predicted Label')
67 plt.ylabel('True Label')
68 plt.tight_layout()
69 plt.show()
70
71 # Training on all train data
72 rs_model = rs_model.fit(X, y)
73
74 # Saving model
75 filename = 'loan_application_model.sav'
76 pickle.dump(rs_model, open(filename, 'wb'))
77

```

Figure 12 – Model training script.


```
[21:28:57] WARNING: ..\src\learner.cc:1061: Starting in XGBoost
1.3.0, the default evaluation metric used with the objective
'binary:logistic' was changed from 'error' to 'logloss'. Explicitly
set eval_metric if you'd like to restore the old behavior.
Fitting 10 folds for each of 200 candidates, totalling 2000 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent
workers.
[Parallel(n_jobs=-1)]: Done 48 tasks      | elapsed: 3.2s
[Parallel(n_jobs=-1)]: Done 348 tasks    | elapsed: 9.7s
[Parallel(n_jobs=-1)]: Done 848 tasks    | elapsed: 20.7s
[Parallel(n_jobs=-1)]: Done 1548 tasks   | elapsed: 36.2s
[21:29:44] WARNING: ..\src\learner.cc:1061: Starting in XGBoost
1.3.0, the default evaluation metric used with the objective
'binary:logistic' was changed from 'error' to 'logloss'. Explicitly
set eval_metric if you'd like to restore the old behavior.
{'colsample_bytree': 0.7618710631033927, 'gamma':
0.4924722485783868, 'learning_rate': 0.13887550893978923,
'max_depth': 3, 'n_estimators': 105, 'subsample':
0.7972050064712304}
Test accuracy: 0.758
```

Figure 14 – Output from console after modelling.

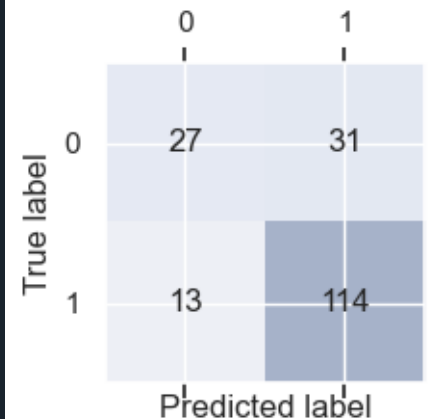


Figure 13 – Confusion Matrix.

Figure 14 and Figure 13 shows the accuracy of the model. The model has an accuracy of 79.7%. 114 of the approved loans have been successfully predicted. 13 of the approved loans have been misclassified as not approved. 27 of the not approved loans has been successfully predicted. 31 of the not approved loans has been misclassified as approved loans.

APPLICATION

Figure 15 shows the imports used in the application script.

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Tue Jun 1 10:56:48 2021
4
5 @author: Eirik
6 """
7
8 import streamlit as st
9 import pandas as pd
10 from sklearn.impute import SimpleImputer
11 from sklearn.preprocessing import MinMaxScaler
12 import pickle
13 import numpy as np
14
```

Figure 15 – Imports for the application script.

The application is created using Streamlit. Streamlit is used to make an easy-to-use interface to interact with the model and to predict the outcome of a loan application. This can be done without having to directly interact with the Python script.

The inputted data is collected and added to a test dataset which does not contain target labels. This is done so that the inputted data will fit into a DataFrame with the same number of columns as the model takes.

The data is scaled the same way as the data the model is trained on, and eventual missing values are imputed.

The outcome of the application is displayed in the Streamlit app.

```

15 # Reading data and making dataframe - Needed because of encoder
16 test_df = pd.read_csv('test.csv')
17
18 # Filling in information in streamlit
19 st.title('Loan application')
20 st.markdown('Please fill in information below: ')
21 gender = st.selectbox('Gender: ', ('Male', 'Female'))
22 married = st.selectbox('Married? ', ('No', 'Yes'))
23 dependents = st.selectbox('How many are dependent on you? ', ('0', '1', '2', '3+'))
24 education = st.selectbox('Education: ', ('Graduate', 'Not Graduate'))
25 self_employed = st.selectbox('Are you self employed? ', ('No', 'Yes'))
26 slider_a_income = st.slider('What is your income?', value=3000, min_value=0,
27                             max_value=10000, step=100)
28 applicant_income = st.number_input('', value=slider_a_income)
29 slider_co_income = st.slider('What is your coapplicant's income?', value=3500,
30                             min_value=0, max_value=10000, step=100)
31 coapplicant_income = st.number_input('', value=slider_co_income)
32 slider_loan_amount = st.slider('How much do u want to loan?', value=200,
33                                min_value=0, max_value=600, step=10)
34 loan_amount = st.number_input('', value=slider_loan_amount)
35 slider_loan_amount_term = st.slider('How Long is the term of the loan in months?',
36                                     value=360, min_value=0, max_value=700, step=10)
37 loan_amount_term = st.number_input('', value=slider_loan_amount_term)
38 credit_history = st.selectbox('Your credit history: ', ('0', '1'))
39 property_area = st.selectbox('Your property area? ', ('Urban', 'Rural',
40                                                         'Semiurban'))
41
42 # Saving information
43 new_data = {
44     'Gender': gender,
45     'Married': married,
46     'Dependents': dependents,
47     'Education': education,
48     'Self_Employed': self_employed,
49     'ApplicantIncome': applicant_income,
50     'CoapplicantIncome': coapplicant_income,
51     'LoanAmount': loan_amount,
52     'Loan_Amount_Term': loan_amount_term,
53     'Credit_History': credit_history,
54     'Property_Area': property_area
55 }
56
57 # Dropping Loan_ID
58 test_df = test_df.drop(['Loan_ID'], axis=1)
59
60 # Making dataframe with inputted data
61 new_df = pd.DataFrame.from_records(new_data, index=[0], columns=new_data.keys())
62
63 # Adding new data to old dataframe because of encoding
64 test_df = test_df.append(new_df, ignore_index=True)
65
66 # Encoding to match model
67 test_df_encoded = pd.get_dummies(data=test_df, drop_first=True,
68                                 columns=['Gender', 'Married', 'Dependents',
69                                         'Education', 'Self_Employed',
70                                         'Property_Area'])
71
72 # Normalizing data using MinMaxScaler
73 scaler = MinMaxScaler()
74 test_df_encoded = pd.DataFrame(scaler.fit_transform(test_df_encoded),
75                               columns=test_df_encoded.columns)
76
77 # Filling NaN data/Missing data with KNNImputer
78 imputer = SimpleImputer(strategy='median')
79 test_df_encoded = pd.DataFrame(imputer.fit_transform(test_df_encoded),
80                               columns=test_df_encoded.columns)
81
82 # Setting up X
83 X = test_df_encoded.values
84
85 # Load model
86 filename = 'loan_application_model.sav'
87 model = pickle.load(open(filename, 'rb'))
88 y_pred = model.predict(X)
89
90 # Saving predictions
91 pred_df = pd.DataFrame(y_pred, columns=['Loan_Status'])
92 pred_df['Loan_Status'] = np.where(pred_df['Loan_Status']==0, 'No', 'Yes')
93
94 result = st.button('Check for Loan approval')
95 if result:
96     if pred_df.Loan_Status.iat[-1] == 'No':
97         st.markdown('Your Loan has been declined')
98     else:
99         st.markdown('Your Loan has been approved')
100

```

Figure 16 – The rest of the application script.

Loan application

Please fill in information below:

Gender:
Male

Married?
No

How many are dependent on you?
0

Education:
Graduate

Are you self employed?
No

What is your income?
3000

What is your coapplicant's income?
3500

How much do u want to loan?
200

How long is the term of the loan in months?
360

Your credit history:
1

Your property area?
Semiurban

Check for loan approval

Your loan has been approved

Figure 18 – Approved Loan in Streamlit

Loan application

Please fill in information below:

Gender:
Female

Married?
No

How many are dependent on you?
3+

Education:
Not Graduate

Are you self employed?
Yes

What is your income?
900

What is your coapplicant's income?
0

How much do u want to loan?
420

How long is the term of the loan in months?
580

Your credit history:
1

Your property area?
Urban

Check for loan approval

Your loan has been declined

Figure 17 – Declined Loan in Streamlit.

EVALUATION AND RESULTS

This chapter presents the results of this study. It includes an evaluation of different choices made in this project.

DATA

The data used for this project was handed out as course material in TIN200 and was to be used despite unfortunate features that might have been discovered during the exploration of the data. The project was based on the concept of the data being good enough to create a working model. If this project were to be implemented in a real-life-scenario, the data would have been reconsidered. Only the train data handed out was used for this project, as the test data was created by filling out the form in the streamlit module.

The data exploration shows that the data contains 614 loan applications, including 13 features each. The data also includes a target label telling the status of the loan. Some of the applications are missing some features, either the "LoanAmount" or the "Loan_Amount_Term". These missing values was imputed using the median, but in hindsight, they might as well have been removed. It does not make sense to apply for a loan, but not give the amount you apply for, or the period of back payment. A solution to remove such data would be to dismiss all applications that is not filled out completely.

A quick overview in the chapter: Exploring the data, shows us by looking at the distributions and the confusion matrix, that the feature with the biggest impact on the loan status seems to be "Credit_History".

MODEL

The model is trained by the help of XGBoostClassifier with an accuracy of 79.7%. More classifiers were tested, some giving better results on the train-test-split, but to avoid overfitting, this classifier was chosen. This was also decided after testing predictions with inputted data which gave an unrealistic big value to the credit history when using other models. This error lies in the data used for training, but preventing overfitting makes the solution more robust, and more applicable to other data.

Figure 13 – Confusion Matrix. shows that we have some room for improvement since 31 loans which was not approved by a loan officer, was approved by the machine. One of the issues of this paper is that we cannot know for sure if the loan officers have made the right decision or not. A model is never better then the people who labels the targets. The better the data, the better the model. This model works very good for demonstrative purposes.

IMPLEMENTATION

The model is implemented in a scrip that is run through the Streamlit module made for python. This module makes it easy to visualise code, and can be used to make an easy-to-use interface, which does not demand any computer knowledge. This is excellent for this demonstration.

The Streamlit module creates an app with a fill-in form that is easy to navigate and understand. After the form is filled in, it gets processed in the background and tested by the model. This takes a few milliseconds, and from there you can see if your application has been approved or declined.

BUSINESS VALUE

We can conclude that this solution, correctly implemented, will give added business value to any bank. We have freed up a lot of capacity and resources by automating the loan process. We have given the customers an effective solution which gains everyone.

We can imagine that this solution can be implemented fast, but it needs good supervision, and an opportunity to solve problems ad hoc. Customer service must also be maintained.

CONCLUSION AND RECOMMENDATIONS

The goal of this project was to research the process of loan applications, create a model which can automate the loan application process, and make a user interface to easily use the model.

This paper explored the basics of the state of the loan application process as is. It discovered that it is a tedious process which in some cases takes up to three weeks, taking a hold of resources and time.

A model was created using course material data, to predict the outcome of loan applications. This model had an accuracy of 79.7%, which meant that in almost 4 of 5 cases, the machine would do the same as the loan officer. This was considered good enough for demonstrative purposes.

The model was implemented to a user interface using Streamlit, making a working demonstration on how an easy fill out form could give instant response to a loan application.

The project is considered a success in the way it demonstrates the possible application of machine learning in the automation of the loan application process. It is easy to use, and incredibly fast. The project demonstrates a clear business value potential.

FURTHER WORK

The following points should be considered done in the continuation of this project:

- Cooperate with a bank to get real life, quality data that would be authentic in a loan application process. The quantity should be increased, and other factors should be considered such as inflation, to future proof the model.
- The model should be trained, tuned, and tested to perfection. A hybrid between a normal programmed program and a machine learning model should be considered, to make the prediction more robust to special cases.
- Integration between the chosen interface, and the existing interfaces of the bank, to make it easy to use this model for non-coders. First step would be to implement something like pyinstaller, which can make python scripts into executables, so you do not need python on your machine to run the model.

ATTACHMENTS

SOURCE CODE

https://github.com/hellund/loan_application_processing.git

REFERENCES

- Aguirre, S. & Rodriguez, A. (2017). *Automation of a Business Process Using Robotic Process Automation (RPA): A Case Study*. Applied Computer Sciences in Engineering, Cartagena, Colombia, pp. 65-71: Springer International Publishing.
- Business Process Management: Blockchain and Robotic Process Automation Forum*. (2020). Business Process Management, Seville, Spain. Cham, Switzerland: Springer.
- Mitra, M. (2018). *Robotic Process Automation(RPA) and Benefits*. Web: Mantra Labs. Available at: <https://www.mantralabsglobal.com/blog/robotic-process-automationrpa-and-benefits/> (accessed: 10.06).
- Molka, T., Gilani, W. & Zeng, X.-J. (2013). *Dotted Chart and Control-Flow Analysis for a Loan Application Process*. Business Process Management Workshops, Tallinn, Estonia, pp. 219-220. Berlin, Heidelberg: Springer.
- Parr, O. S. (2020, 1. juli 2020). Travle boligkjøpere stresser bankene. *Finansavisen*. Available at: <https://finansavisen.no/nyheter/bolig/2020/07/01/7540549/rekordaktivt-boligmarked-gir-lanerush-dnb-beklager-lang-ventetid> (accessed: 10.06.2021).
- Raschka, S. & Mijalili, V. (2019). *Python Machine Learning*. Third ed. Birmingham, UK: Packt Publishing Ltd.
- Rizk, Y., Isahagian, V., Boag, S., Khazaeni, Y., Unuvar, M., Muthusamy, V. & Khalaf, R. (2020). *A Conversational Digital Assistant for*
- Intelligent Process Automation*. Business Process Management: Blockchain and Robotic Process Automation Forum, pp. 85–100: Springer Nature Switzerland AG 2020.
- Sweney, M. & Canon, G. (2021, 28. april 2021). Alphabet: revenue soars for Google owner as Covid lockdown boom continues. *The Guardian*. Available at: <https://www.theguardian.com/technology/2021/apr/27/alphabet-google-revenue-quarterly-earnings> (accessed: 10.06.2021).
- Treder, M. (2019). *Becoming a data-driven Organisation*
- Unlock the value of data*. 1 ed. Berlin, Heidelberg: Springer Vieweg.
- Werger, K., Kenedy, J., Peckham, D., Mather, S., Ginsberg, R., Jooste, A., Robinson, A. & Knappenberger, D. (2020). Data valuation: Understanding the value of your data assets. 12. Available at: <https://www2.deloitte.com/content/dam/Deloitte/global/Documents/Finance/Valuation-Data-Digital.pdf> (accessed: 10.06.2021).
- The world's most valuable resource is no longer oil, but data. (2017, May 6th 2017). *The Economist*. Available at: <https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data> (accessed: 10.06.2021).