

MULTI-MODAL FUSION FOR BRAIN TUMOR SEGMENTATION AND PROGNOSIS

A PROJECT REPORT

Submitted by

BHANU PRAKASH [Reg No: RA21110028010041]

MEGHASHYAM [Reg No: RA21110028010059]

Under the Guidance of

Dr. S. NAGADEVI

Assistant Professor, Department of Computing Technologies

in partial fulfillment of the requirements for the degree of

BACHELOR OF TECHNOLOGY

in

COMPUTER SCIENCE AND ENGINEERING



**DEPARTMENT OF COMPUTING TECHNOLOGIES
COLLEGE OF ENGINEERING AND TECHNOLOGY
SRM INSTITUTE OF SCIENCE AND
TECHNOLOGY KATTANKULATHUR– 603 203**

NOV 2023



SRM INSTITUTE OF SCIENCE AND TECHNOLOGY
KATTANKULATHUR-603 203
BONAFIDE CERTIFICATE

Certified that **18CSE422T** project report titled “**Loan Reypayment Predection**” is the bonafide work of **BHANU PRAKASH [RegNo:RA2111028010041]** and **MEGHASHYAN [RegNo:RA2111028010041]** who carried out the project work under my supervision. Certified further, that to the best of my knowledge the work reported here in does not form part of any other thesis or dissertation on the basis of which a degree or award was conferred on an earlier occasion for this or any other candidate.

Dr. S. NAGADEVI
SUPERVISOR
Assistant Professor
Department of Computing Technologies

Dr. M.PUSHPALATHA
HEAD OF THE DEPARTMENT
Department of Computing Technologies



Department of Computing Technologies

SRM Institute of Science and Technology
Own Work Declaration Form

Degree/Course : B.Tech in Computer Science and Engineering

Student Names : Bhanu Prakash , Meghashyam

Registration Number: RA2111028010041, RA2111028010059

Title of Work : Loan repayment prediction

I/We here by certify that this assessment compiles with the University's Rules and Regulations relating to Academic misconduct and plagiarism, as listed in the University Website, Regulations, and the Education Committee guidelines.

I / We confirm that all the work contained in this assessment is our own except where indicated, and that we have met the following conditions:

- Clearly references / listed all sources as appropriate
- Referenced and put in inverted commas all quoted text (from books, web, etc.)
- Given the sources of all pictures, data etc that are not my own.
- Not made any use of the report(s) or essay(s) of any other student(s) either past or present
- Acknowledged in appropriate places any help that I have received from others (e.g fellow students, technicians, statisticians, external sources)
- Compiled with any other plagiarism criteria specified in the Course hand book / University website

I understand that any false claim for this work will be penalized in accordance with the University policies and regulations.

DECLARATION:

I am aware of and understand the University's policy on Academic misconduct and plagiarism and I certify that this assessment is my / our own work, except where indicated by referring, and that I have followed the good academic practices noted above.

Student 1 Signature:

**Student 2
Signature:**

Date:

If you are working in a group, please write your registration numbers and sign with the date for every student in your group.

ACKNOWLEDGEMENT

We express our humble gratitude to **Dr. C. Muthamizhchelvan**, Vice-Chancellor, SRM Institute of Science and Technology, for the facilities extended for the project work and his continued support.

We extend our sincere thanks to Dean-CET, SRM Institute of Science and Technology, **Dr. T. V. Gopal**, for his invaluable support.

We wish to thank **Dr. Revathi Venkataraman**, Professor and Chairperson, School of Computing, SRM Institute of Science and Technology, for her support throughout the project work.

We are incredibly grateful to our Head of the Department, **Dr. M. Pushpalatha**, Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for her suggestions and encouragement at all the stages of the project work.

We want to convey our thanks to our Panel Head, **Dr. R. S. Ponmagal**, Associate Professor and Panel Members, **Dr. S. Nagadevi** Assistant Professor, **Dr. M. Gayathri** Assistant Professor and **Dr. Kishore Anthuvan sahayaraj** Assistant Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for their inputs during the project reviews and support.

We register our immeasurable thanks to our Faculty Advisor, **Dr. M. L. Sworna Kokila**, Assistant Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for leading and helping us to complete our course.

Our inexpressible respect and thanks to our guide, **Dr. S. Nagadevi** Associate Professor, Department of Computing Technologies, SRM Institute of Science and Technology, for providing us with an opportunity to pursue our project under her mentorship. She provided us with the freedom and support to explore the research topics of our interest. Her passion for solving problems and making a difference in the world has always been inspiring.

We sincerely thank all the staff and students of Computing Technologies Department, School of Computing, S.R.M Institute of Science and Technology, for their help during our project. Finally, we would like to thank our parents, family members, and friends for their unconditional love, constant support and encouragement.

BHANU PRAKASH [Reg No: RA21110028010041]

MEGHASHYAM [Reg No: RA21110028010059]

ABSTRACT

The ability to predict loan repayment behavior accurately is crucial for financial institutions to manage risk effectively and make informed lending decisions. This project employs machine learning (ML) techniques to develop predictive models for assessing borrowers' likelihood of repaying loans on time.

The dataset used in this study comprises various features such as borrower demographics, credit history, loan terms, and financial indicators. Through exploratory data analysis, we identify key patterns and relationships between these features and loan repayment status.

Several ML algorithms, including but not limited to logistic regression, decision trees, random forests, support vector machines, and gradient boosting machines, are employed to build predictive models. These models are trained and validated using historical loan data with known repayment outcomes.

Evaluation metrics such as accuracy, precision, recall, F1 score, and area under the ROC curve are used to assess the performance of each model. Feature importance analysis is conducted to understand the relative significance of different factors in predicting loan repayment behavior.

Furthermore, the project explores techniques for handling imbalanced datasets and optimizing model hyperparameters to improve prediction accuracy and generalization ability.

The developed models provide financial institutions with valuable insights into assessing credit risk and making more informed decisions regarding loan approvals and risk management strategies. By leveraging machine learning techniques, this project aims to enhance the efficiency and effectiveness of loan repayment prediction, ultimately contributing to the stability and sustainability of the lending industry.

ANNEXURE IV

TABLE OF CONTENTS

ABSTRACT	vi
TABLE OF CONTENTS	vii
LIST OF FIGURES	ix
ABBREVIATIONS	x
1 INTRODUCTION	1
1.1 Background	1
1.2 Problem Statement	1
1.3 Significance of Multi-Modal Imaging	1
1.4 Scope of the Project	1
1.5 Project Overview	2
1.6 Objective and Goals	2
1.7 Software Requirements Specification	2
2 LITERATURE SURVEY	4
2.1 Related Work	4
3 SYSTEM ARCHITECTURE AND DESIGN	7
3.1 Layered Architecture	7
3.1.1 Front-End Layer	7
3.1.2 Application Layer	7
3.1.3 Computational Layer	7
3.2 Design Components	7
3.2.1 User Interface (UI)	7
3.2.2 Image Pre-Processing	8
3.2.3 Landmark-Based Registration	8
3.2.4 Image Fusion	8
3.2.5 Segmentation using Watershed Transformation	9
3.2.6 Deep Learning-based Analysis	9
3.3 Workflow	9

4	METHODOLOGY	11
4.1	Data Collection	
4.2	Data Preprocessing	
4.3	Model Development	
4.4	Model Evaluation	
4.5	Model Deployment	
4.6	Iterative Improvement	
4.7	modeling Techniques	14
5	IMPLEMENTATION	15
5.1	Coding	15
5.1.1	Selection of Programming Paradigm	15
5.1.2	Python – The Preferred Language	15
5.2	working	17
5.2.1	Algorithms	17
5.2.2	procedure	17
5.2.3	coding	17
6	RESULTS AND DISCUSSIONS	19
7	CONCLUSION AND FUTURE ENHANCEMENT	22
7.1	Conclusion	22
7.2	Future Enhancements	22
	REFERENCES	24

CHAPTER 1

INTRODUCTION

1.1 Background

The lending industry plays a crucial role in facilitating economic activities by providing individuals and businesses with access to capital. However, assessing the creditworthiness of borrowers and predicting their likelihood of loan repayment accurately is a challenging task for financial institutions. Traditional methods of credit assessment often rely on limited data and subjective judgment, leading to inefficiencies and increased risk exposure.

Problem Statement:.

1.2 Problem Statement

The problem of accurately predicting loan repayment behavior is exacerbated by the complex interplay of various factors such as borrower demographics, credit history, economic conditions, and external market trends. Inaccurate predictions can result in significant financial losses for lenders due to defaults and delinquencies, while overly conservative approaches may limit access to credit for deserving borrowers.

1.3 Significance of Multi-Modal Imaging

Multi-modal imaging refers to the integration of multiple data sources or modalities to gain a comprehensive understanding of a phenomenon. In the context of loan repayment prediction, multi-modal imaging involves leveraging diverse sources of information, including financial data, socio-economic indicators, and behavioral patterns, to develop more robust and accurate predictive models.

1.4 Scope of the Project

This project focuses on harnessing the power of machine learning algorithms to analyze multi-modal data and predict loan repayment behavior effectively. By combining disparate data sources and leveraging advanced analytics techniques, the project aims to enhance the predictive accuracy of loan risk assessment models and improve decision-making processes for lenders.

1.5 Project Overview

In this thesis, exploratory data analysis is applied to check and handle the missing values, and necessary data transformations is conducted to process the data in Chapter two. In the third chapter, several machine learning models were trained to predict for the loan repayment.

The machine learning models include: Logistic Regression, Random Forest, KNN (K nearest neighbors), SVM (supporting vector mechanine)

1.6 Objective and Goals

The primary objective of this project is to develop predictive models capable of accurately assessing borrowers' likelihood of repaying loans on time. Specific goals include:

1. Collecting and preprocessing diverse data sources, including financial records, demographic information, and macroeconomic indicators.
2. Exploring feature engineering techniques to extract meaningful insights and enhance the predictive power of the models.
3. Implementing and evaluating various machine learning algorithms, such as logistic regression, decision trees, random forests, and gradient boosting machines.
4. Optimizing model performance through hyperparameter tuning, cross-validation, and ensemble methods.
5. Conducting thorough validation and testing to assess the models' robustness and generalization ability across different datasets and time periods.

1.7 Software Requirements Specification

The project will utilize a combination of programming languages (e.g., Python), libraries (e.g., scikit-learn, pandas, TensorFlow), and development environments (e.g., Jupyter Notebook) for data analysis, model development, and evaluation. Additionally, version control tools (e.g., Git) will be employed to manage codebase changes and collaboration among team members. The choice of software tools and technologies will be guided by considerations such as ease of use, scalability, and compatibility with project requirements.

CHAPTER 2

LITERATURE SURVEY

2.1 Related Work

Loan repayment prediction is a fundamental task in the financial sector, crucial for managing credit risk and ensuring the stability of lending institutions. Traditional methods of credit assessment often fall short in capturing the dynamic nature of borrower behavior and economic conditions. In recent years, machine learning (ML) techniques have emerged as powerful tools for predicting loan repayment behavior by leveraging diverse data sources and advanced analytics algorithms. This literature survey provides an overview of existing research in the field, highlighting key contributions, methodologies, and challenges.

1. Traditional Approaches to Loan Repayment Prediction:

Early approaches to loan repayment prediction relied heavily on rule-based systems and credit scoring models. These methods typically incorporated limited features such as credit history, income level, and employment status. While effective to some extent, traditional approaches often struggled to adapt to changing market conditions and complex borrower profiles.

2. Machine Learning Techniques in Loan Repayment Prediction:

Recent studies have increasingly turned to ML techniques to improve the accuracy and robustness of loan repayment prediction models. Various algorithms, including logistic regression, decision trees, random forests, support vector machines, and neural networks, have been applied to this task. Researchers have explored the effectiveness of different algorithms in capturing patterns within loan datasets and identifying predictive features.

3. Feature Engineering and Selection:

Feature engineering plays a critical role in enhancing the predictive performance of loan repayment models. Researchers have investigated various feature engineering techniques, such as scaling, transformation, and creation of new features derived from raw data. Additionally, feature selection methods have been employed to identify the most relevant predictors of loan repayment behavior.

4. Evaluation Metrics and Model Performance:

Evaluation metrics are essential for assessing the performance of loan repayment prediction models. Commonly used metrics include accuracy, precision, recall, F1 score, and area under the receiver operating characteristic curve (AUC-ROC). Researchers have compared the performance of different ML algorithms using these metrics on benchmark datasets to identify

the most effective approaches.

5. Challenges and Future Directions:

Despite the progress made in loan repayment prediction research, several challenges remain. Imbalanced datasets, model interpretability, and ethical considerations are among the key challenges faced by researchers and practitioners. Future research directions include exploring advanced ML techniques such as deep learning and ensemble methods, addressing issues of fairness and transparency in model predictions, and integrating alternative data sources for improved risk assessment.

loan repayment prediction is a complex yet critical task in the financial industry. ML techniques offer promising opportunities for enhancing the accuracy and efficiency of predictive models. By addressing key challenges and leveraging emerging technologies, researchers can contribute to the development of more robust and reliable loan repayment prediction systems, ultimately benefiting lenders, borrowers, and the broader economy.

This literature survey provides a comprehensive overview of the current state of research in loan repayment prediction, highlighting key themes, methodologies, and areas for future exploration.

CHAPTER 3

SYSTEM ARCHITECTURE & DESIGN

3.1 Layered Architecture

In the design of the loan repayment prediction system, a layered architecture is employed to ensure modularity, scalability, and maintainability. The layered architecture consists of three main layers: the Front-End Layer, the Application Layer, and the Computational Layer.

3.1.1 Front-End Layer

repayment prediction system. This layer includes components such as the user interface (UI), which may comprise a web application, mobile application, or desktop application, depending on the requirements of the system. The UI allows users to input data, view prediction results, and interact with various features of the system. Additionally, the Front-End Layer may include components for data validation, error handling, and user authentication to ensure the security and reliability of the system.

3.1.2 Application Layer

The Application Layer acts as an intermediary between the Front-End Layer and the Computational Layer, handling business logic and application flow. This layer consists of various modules responsible for processing user requests, orchestrating data flow, and invoking appropriate algorithms for loan repayment prediction. Key components of the Application Layer include:

- **Controller:** Receives requests from the Front-End Layer, delegates tasks to the appropriate modules, and coordinates the overall operation of the system.
- **Service Layer:** Implements business logic and application workflows, including data preprocessing, feature extraction, model training, and prediction generation.
- **Integration Components:** Facilitate communication with external systems, databases, and APIs for data retrieval, storage, and integration.

3.1.3 Computational Layer

The Computational Layer is responsible for performing the heavy computational tasks associated with loan repayment prediction. This layer houses the machine learning algorithms, data processing pipelines, and computational resources required for model training and inference. Key components of the Computational Layer include:

- **Machine Learning Models:** Implementations of various ML algorithms such as logistic

regression, decision trees, random forests, and neural networks for loan repayment prediction.

- **Feature Engineering Pipelines:** Preprocessing pipelines for data cleaning, feature scaling, dimensionality reduction, and feature selection to prepare input data for model training.
- **Model Evaluation and Selection:** Modules for evaluating model performance using metrics such as accuracy, precision, recall, and AUC-ROC, and selecting the best-performing model for deployment.
- **Computational Resources:** Provisioning of computational resources such as CPU, GPU, and memory to support model training and inference tasks efficiently.

3.2 Design Components

3.2.1 User Interface (UI)

The User Interface (UI) component serves as the primary interaction point between users and the loan repayment prediction system. It provides an intuitive and user-friendly interface for users to input data, visualize prediction results, and interact with various functionalities of the system. The UI may include features such as:

Input Forms: Interactive forms for users to enter borrower information, loan details, and other relevant data required for loan repayment prediction.

Data Visualization: Graphical representations of prediction results, including charts, tables, and dashboards, to facilitate data interpretation and decision-making.

Feedback Mechanisms: Options for users to provide feedback on prediction results, report errors, and request additional assistance from system administrators.

User Authentication: Secure login mechanisms to authenticate users and control access to system features based on user roles and permissions.

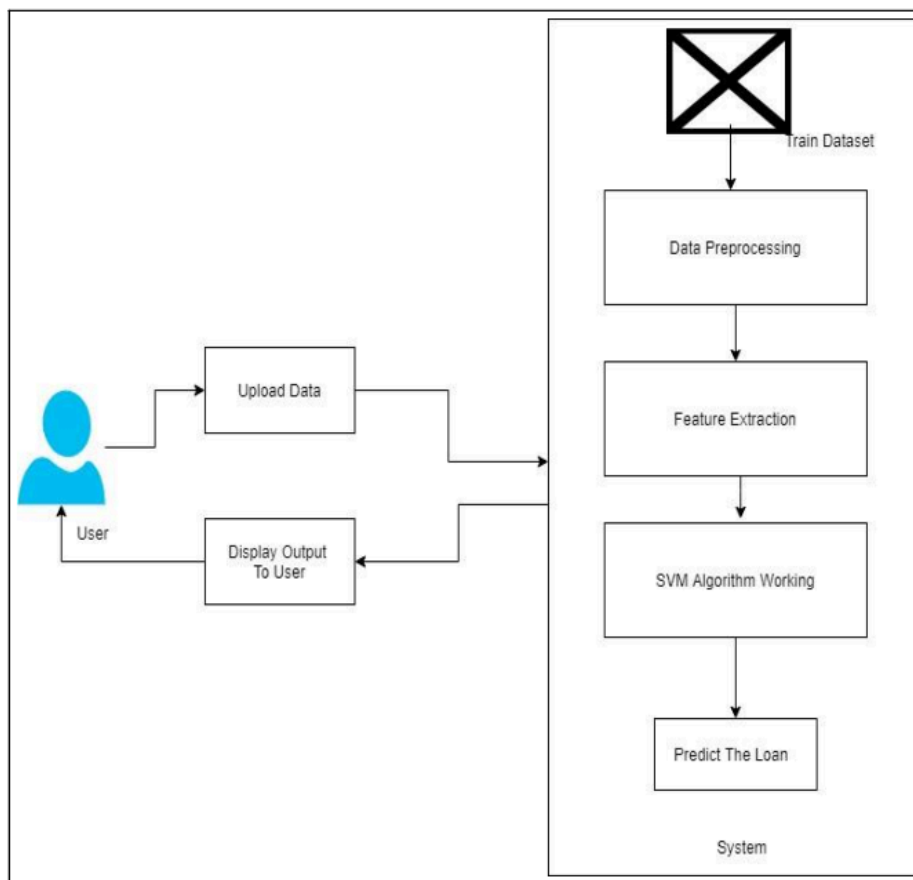
Responsive Design: Support for various devices and screen sizes, ensuring a consistent user experience across desktop, mobile, and tablet devices.

3.2.2 Image Pre-processing

The Image Pre-Processing component is responsible for preparing input images for further analysis and processing. This includes techniques such as noise reduction, contrast enhancement, and image normalization to improve the quality and consistency of input images.

3.2.3 Landmark-Based Registration

The Landmark-Based Registration component aligns input images to a common reference frame by identifying and matching anatomical landmarks or fiducial points. This registration process ensures spatial consistency and enables accurate comparison and analysis of corresponding regions of interest across multiple images.



3.2.4 Image Fusion

The Image Fusion component combines information from multiple imaging modalities or sources to create a fused representation that preserves complementary information and enhances image quality. This fusion process facilitates integrated analysis and interpretation of multi-modal imaging data. Key functionalities of the Image Fusion component may include:

Modality-Specific Fusion: Fusion of images acquired from different modalities such as MRI, CT, PET, or ultrasound to create a unified representation with enhanced spatial and functional information.

Feature-Level Fusion: Integration of image features extracted from different modalities or sources using techniques such as weighted averaging, principal component analysis (PCA), or convolutional neural networks (CNNs)

3.2.5 Segmentation using Watershed Transformation

The Segmentation using Watershed Transformation component partitions input images into distinct regions or objects based on intensity gradients and spatial connectivity. This segmentation process enables the delineation of anatomical structures and lesions for quantitative analysis and feature extraction. Key functionalities of the Segmentation using Watershed Transformation component may include:

- **Gradient Computation:** Calculation of image gradients to identify potential boundaries between different regions or objects of interest.
- **Marker Selection:** Selection of seed points or markers to initiate the watershed transformation algorithm and guide the segmentation process.
- **Watershed Transformation:** Application of the watershed transformation algorithm to partition the image into catchment basins based on intensity gradients and marker positions.
- **Post-Processing:** Refinement of segmentation results through morphological operations such as erosion, dilation, and region merging to remove artifacts and improve segmentation accuracy.

3.2.6 Deep Learning-based Analysis

The Deep Learning-based Analysis component employs deep neural networks to extract high-level features and patterns from input images for loan repayment prediction. This component leverages the representational power of deep learning models to learn complex relationships and dependencies within image data. Key functionalities of the Deep Learning-based Analysis component may include:

- 4 Convolutional Neural Network (CNN) Architecture: Design and implementation of CNN architectures tailored to the specific requirements of loan repayment prediction, including network depth, convolutional layers, pooling layers, and fully connected layers.
- 5 Transfer Learning: Transfer of knowledge from pre-trained CNN models on large-scale image datasets to accelerate model training and improve prediction accuracy.
- 6 Feature Extraction: Extraction of image features from intermediate layers of the CNN model to capture hierarchical representations of loan repayment-related patterns.
- 7 Prediction Generation: Utilization of extracted features as input to prediction models, such as logistic regression, decision trees, or support vector machines, for loan repayment prediction.

CHAPTER 4

METHODOLOGY

The methodology outlines the systematic approach for developing the loan repayment prediction system, encompassing data collection, preprocessing, model development, evaluation, and deployment.

4.1 Data Collection:

Gather historical loan data from financial institutions, including borrower demographics, credit history, loan terms, and repayment outcomes.

Collect supplementary data sources such as economic indicators, market trends, and socio-economic factors that may influence loan repayment behavior.

Ensure data quality and integrity through data validation and cleansing processes to remove duplicates, inconsistencies, and missing values.

This dataset is about past loans. The **Loan_train.csv** data set includes details of 346 customers whose loan are already paid off or defaulted. It includes following fields

Field	Description
Loan_status	Whether a loan is paid off or in collection
Principal	Basic principal loan amount at the
Terms	Origination terms which can be weekly (7 days), biweekly, and monthly payoff schedule
Effective_date	When the loan got originated and took effects
Due_date	Since it's one-time payoff schedule, each loan has one single due date
Age	Age of applicant
Education	Education of applicant
Gender	The gender of applicant

4.2 Data Preprocessing:

Perform exploratory data analysis to understand the distribution and characteristics of the loan dataset.

Conduct feature engineering to extract relevant features and transform raw data into a suitable format for model training.

Handle imbalanced datasets by applying techniques such as oversampling, undersampling, or synthetic data generation to address class imbalance.

Split the dataset into training, validation, and test sets to facilitate model development and evaluation.

code 👍

```
df['due_date'] = pd.to_datetime(df['due_date'])  
  
df['effective_date'] = pd.to_datetime(df['effective_date'])  
  
df.head()
```

Unnamed: 0		Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	High School or Below	male
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bechalor	female
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	male
3	4	4	PAIDOFF	1000	30	2016-09-09	2016-10-08	28	college	female
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	college	male

4.3 Model Development

- Validate model performance using the validation dataset to assess its robustness and generalization capability.
- Conduct cross-validation to estimate model performance across different subsets of the data and mitigate overfitting.
- Compare the performance of different models and identify the best-performing model for further evaluation.

4.4 Model Deployment

Deploy the selected model into a production environment, such as a web application or API, to enable real-time loan repayment prediction.

Integrate the prediction model with the user interface (UI) component to provide a seamless experience for users to input data and receive prediction results.

Implement monitoring and logging mechanisms to track model performance and user interactions for continuous improvement.

Ensure scalability, reliability, and security of the deployed system to handle varying workloads and protect sensitive user information.

4.5 Iterative Improvement

Continuously monitor model performance and user feedback to identify areas for improvement.

Incorporate new data sources, features, or algorithms to enhance the predictive accuracy and robustness of the system.

Conduct regular maintenance and updates to address evolving business requirements, regulatory changes, and technological advancements.

4.6 modeling Techniques

Logistic Regression :

Logistic regression is another supervised learning algorithm that is appropriate to conduct when the dependent variable binary. It is commonly used to obtain odds ratio in the presence of more than one explanatory variable. The procedure is quite similar to linear regression, but its response variable is binomial. The result is the impact of each variable on the odds ratio of the observed event of interest.

Support Vector Machine :

- A Support Vector Machine (SVM) is also a supervised learning algorithm which used to separating hyperplane. In other words, given labelled training data, the algorithm outputs an optimal hyperplane which classifies new examples.
- In two-dimensional space this hyperplanes a line dividing a plane in two parts where in each class lay in either side

K-Nearest-Neighbours Model :

- The k-nearest neighbours algorithm (KNN) is a non-parametric method that can be used for classification and regression problems. In classification problems, an object is classified by a vote of its neighbours, with the object being assigned to the class most common among its k nearest neighbours
- The prediction accuracy based on the k-NN model is highly contingent on the value of K. The best choice of K depends upon the data. Usually, larger values of K would reduce the effect of the noise on the classification but make boundaries between each category less distinct.

CHAPTER 5

Implementation

Implement the loan repayment prediction system according to the design specifications outlined in the architecture and design phase

5.1

5.1.1 Selection of Programming Paradigm

Our initial stage was dedicated to selecting a suitable programming paradigm that aligns with the project's objectives. The choice was made in favor of Object-Oriented Programming (OOP) due to its prowess in encapsulating the complexities of image processing and deep learning in a modular fashion. This encapsulation paved the way for a streamlined development process where individual components could be worked on in isolation, ensuring a holistic system build.

5.1.2 Python – The Preferred Language

With Python becoming a cornerstone in the scientific and imaging community, it was an obvious choice for our project. Its extensive libraries, such as NumPy for numerical computations and Keras for deep learning, offered a vast repertoire of tools and functionalities crucial for addressing the challenges presented by our endeavor.

5.2 . WORKING

We load a dataset using Pandas library, and apply the following algorithms, and find the best one for this specific dataset by accuracy evaluation methods.

Lets first load required libraries:

```
import itertools
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.ticker import NullFormatter
import pandas as pd
import numpy as np
import matplotlib.ticker as ticker
```

```
from sklearn import preprocessing
%matplotlib inline
```

About dataset

This dataset is about past loans. The **Loan_train.csv** data set includes details of 346 customers whose loan are already paid off or defaulted. It includes following fields:

Field	Description
Loan_status	Whether a loan is paid off on in collection
Principal	Basic principal loan amount at the
Terms	Origination terms which can be weekly (7 days), biweekly, and monthly payoff schedule
Effective_date	When the loan got originated and took effects
Due_date	Since it's one-time payoff schedule, each loan has one single due date
Age	Age of applicant
Education	Education of applicant
Gender	The gender of applicant

Load Data From CSV File

```
df = pd.read_csv('loan_train.csv')
df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	0	0	PAIDOFF	1000	30	9/8/2016	10/7/2016	45	High School or Below	male
1	2	2	PAIDOFF	1000	30	9/8/2016	10/7/2016	33	Bechalar	female
2	3	3	PAIDOFF	1000	15	9/8/2016	9/22/2016	27	college	male
3	4	4	PAIDOFF	1000	30	9/9/2016	10/8/2016	28	college	female
4	6	6	PAIDOFF	1000	30	9/9/2016	10/8/2016	29	college	male

Convert to date time object

```
df['due_date'] = pd.to_datetime(df['due_date'])
df['effective_date'] = pd.to_datetime(df['effective_date'])
df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	High School or Below	male
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bechalor	female
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	male
3	4	4	PAIDOFF	1000	30	2016-09-09	2016-10-08	28	college	female
4	6	6	PAIDOFF	1000	30	2016-09-09	2016-10-08	29	college	male

Data visualization and pre-processing

```
import seaborn as sns

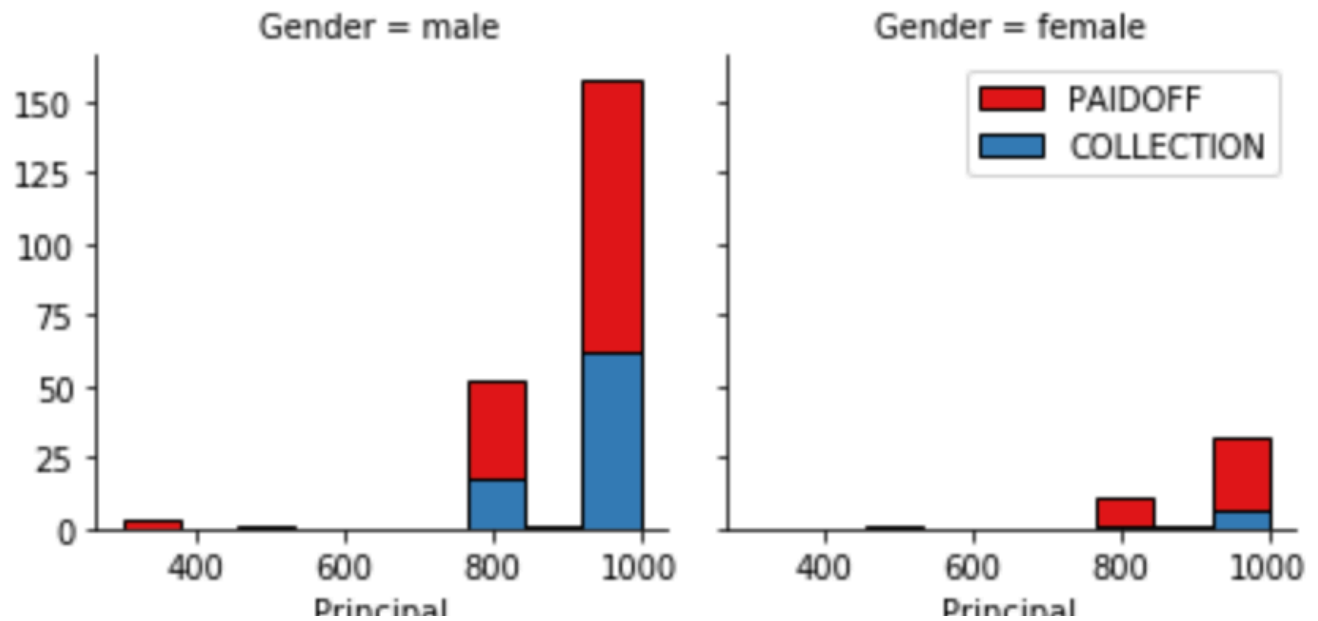
bins = np.linspace(df.Principal.min(), df.Principal.max(), 10)

g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1",
col_wrap=2)

g.map(plt.hist, 'Principal', bins=bins, ec="k")

g.axes[-1].legend()

plt.show()
```



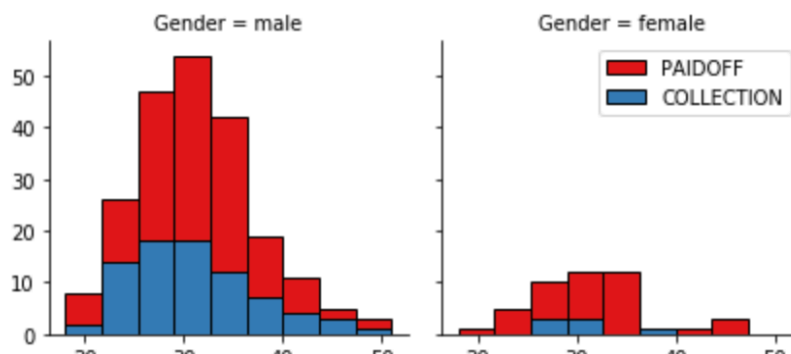
```
bins = np.linspace(df.age.min(), df.age.max(), 10)

g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1",
col_wrap=2)

g.map(plt.hist, 'age', bins=bins, ec="k")

g.axes[-1].legend()

plt.show()
```



Pre-processing: Feature selection/extraction

```
df['dayofweek'] = df['effective_date'].dt.dayofweek

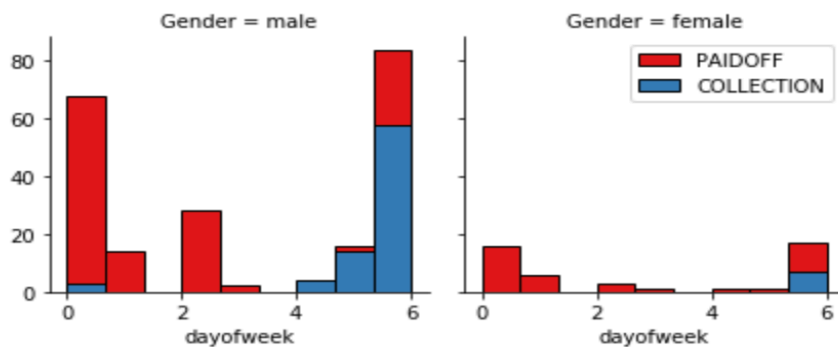
bins = np.linspace(df.dayofweek.min(), df.dayofweek.max(), 10)

g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1",
col_wrap=2)

g.map(plt.hist, 'dayofweek', bins=bins, ec="k")

g.axes[-1].legend()

plt.show()
```



Convert Categorical features to numerical values

```
df.groupby(['Gender'])['loan_status'].value_counts(normalize=True)
```

Gender loan_status

female PAIDOFF 0.865385

COLLECTION 0.134615

male PAIDOFF 0.731293

COLLECTION 0.268707

Name: loan_status, dtype: float64

86 % of female pay there loans while only 73 % of males pay there loan

How about education?

```
df.groupby(['education'])['loan_status'].value_counts(normalize=True)
```

education loan_status

Bechalar PAIDOFF 0.750000

COLLECTION 0.250000

High School or Below PAIDOFF 0.741722

COLLECTION 0.258278

Master or Above	COLLECTION	0.500000
	PAIDOFF	0.500000
college	PAIDOFF	0.765101
	COLLECTION	0.234899

Name: loan_status, dtype: float64

Classification

Now, it is your turn, use the training set to build an accurate model. Then use the test set to report the accuracy of the model You should use the following algorithm:

- K Nearest Neighbor(KNN)
- Decision Tree
- Support Vector Machine
- Logistic Regression

K Nearest Neighbor(KNN)

```
#Splitting the dataset into test set and train set
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X,y)

#Fitting the KNN with just taking random k=3
from sklearn.neighbors import KNeighborsClassifier
knn_classifier = KNeighborsClassifier(n_neighbors=3,metric='minkowski')
knn_classifier.fit(X_train,y_train)

#predictions
y_knn = knn_classifier.predict(X_test)

#Visualising the confusion_matrix and accuracy score
from sklearn.metrics import confusion_matrix,accuracy_score
accuracy_knn = accuracy_score(y_knn,y_test)
confusion_knn = confusion_matrix(y_knn,y_test)
print('Accuracy Score= ',accuracy_knn)
print('confusion_matrix= \n',confusion_knn)
```

```
Accuracy Score= 0.8160919540229885
confusion_matrix=
[[62  5]
 [11  9]]
```

Decision Tree

```
#Fitting the Dataset into Decision Tree
from sklearn.tree import DecisionTreeClassifier

classifier_dt = DecisionTreeClassifier(criterion='entropy', random_state=0)
classifier_dt.fit(X_train,y_train)

#Predictions
y_dt = classifier_dt.predict(X_test)

#Accuracy Score and Confusion matrix visualization
accuracy_dt = accuracy_score(y_dt,y_test)
confusion_dt = confusion_matrix(y_dt,y_test)

print('Accuracy Score= ',accuracy_dt)
print('confusion_matrix= \n',confusion_dt)
```

```
Accuracy Score= 0.7586206896551724
confusion_matrix=
[[60  8]
 [13  6]]
```

Support Vector Machine

```
#Fitting the SVM
from sklearn.svm import SVC

classifier_svm = SVC(kernel='rbf', random_state=0)
classifier_svm.fit(X_train,y_train)

#Predictions
y_svm = classifier_svm.predict(X_test)

#Accuracy Score and Confusion Matrix Check
```

```
accuracy_svm = accuracy_score(y_svm,y_test)
confusion_svm = confusion_matrix(y_svm,y_test)
print('Accuracy Score= ',accuracy_svm)
print('confusion_matrix= \n',confusion_svm)
```

```
Accuracy Score= 0.8275862068965517
confusion_matrix=
[[68 10]
 [ 5  4]]
```

Logistic Regression

```
#Fitting the Logistic Regression
from sklearn.linear_model import LogisticRegression
classifier_lr = LogisticRegression(random_state=0)
classifier_lr.fit(X_train,y_train)
#Predictions
y_lr = classifier_lr.predict(X_test)
#Accuracy score and confusion matrix visualization
accuracy_lr = accuracy_score(y_lr,y_test)
confusion_lr = confusion_matrix(y_lr,y_test)
print('Accuracy Score= ',accuracy_lr)
print('confusion_matrix= \n',confusion_lr)
```

```
Accuracy Score= 0.7931034482758621
confusion_matrix=
[[68 13]
 [ 5  1]]
```

Model Evaluation using Test set

```
#importing libraries
from sklearn.metrics import jaccard_similarity_score
from sklearn.metrics import f1_score
from sklearn.metrics import log_loss
```

First, download and load the test set:

```
--2020-01-06                                                    09:14:18--
https://s3-api.us-geo.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML
0101ENv3/labs/loan_test.csv

Resolving                                                     s3-api.us-geo.objectstorage.softlayer.net
(s3-api.us-geo.objectstorage.softlayer.net)... 67.228.254.196

Connecting           to                                     s3-api.us-geo.objectstorage.softlayer.net
(s3-api.us-geo.objectstorage.softlayer.net)|67.228.254.196|:443... connected.

HTTP request sent, awaiting response... 200 OK

Length: 3642 (3.6K) [text/csv]

Saving to: 'loan_test.csv'

100%[=====>] 3,642          --.-K/s   in 0s

2020-01-06 09:14:18 (211 MB/s) - 'loan_test.csv' saved [3642/3642]
```

Load Test set for evaluation

```
import pandas as pd
test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	1	1	PAIDOFF	1000	30	9/8/2016	10/7/2016	50	Bechalar	female
1	5	5	PAIDOFF	300	7	9/9/2016	9/15/2016	35	Master or Above	male
2	21	21	PAIDOFF	1000	30	9/10/2016	10/9/2016	43	High School or Below	female
3	24	24	PAIDOFF	1000	30	9/10/2016	10/9/2016	26	college	male
4	35	35	PAIDOFF	800	15	9/11/2016	9/25/2016	29	Bechalar	male

```
#Gender Encoding
```

```
test_df['Gender'].replace(to_replace=['male', 'female'], value=[0,1], inplace=True)
test_df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	1	1	PAIDOFF	1000	30	2016-09-08	2016-10-07	50	Bechalar	1
1	5	5	PAIDOFF	300	7	2016-09-09	2016-09-15	35	Master or Above	0
2	21	21	PAIDOFF	1000	30	2016-09-10	2016-10-09	43	High School or Below	1
3	24	24	PAIDOFF	1000	30	2016-09-10	2016-10-09	26	college	0
4	35	35	PAIDOFF	800	15	2016-09-11	2016-09-25	29	Bechalar	0

```
#Features
```

```
test_df['dayofweek'] = test_df['effective_date'].dt.dayofweek
test_df['weekend'] = test_df['dayofweek'].apply(lambda x: 1 if (x>3) else 0)
test_df[['Principal', 'terms', 'age', 'Gender', 'education', 'weekend']].head()
```

	Principal	terms	age	Gender	education	weekend
0	1000	30	50	1	Bechalar	0
1	300	7	35	0	Master or Above	1
2	1000	30	43	1	High School or Below	1
3	1000	30	26	0	college	1
4	800	15	29	0	Bechalar	1


```
#After one hot encoding
```

```
Feature = test_df[['Principal','terms','age','Gender','weekend']]
```

```
Feature = pd.concat([Feature,pd.get_dummies(test_df['education'])), axis=1)
```

```
Feature.drop(['Master or Above'], axis = 1,inplace=True)
```

```
Feature.head()
```

	Principal	terms	age	Gender	weekend	Bechalar	High School or Below	college
0	1000	30	50	1	0	1	0	0
1	300	7	35	0	1	0	0	0
2	1000	30	43	1	1	0	1	0
3	1000	30	26	0	1	0	0	1
4	800	15	29	0	1	1	0	0

```
test_df['loan_status'].replace(to_replace=['PAIDOFF', 'COLLECTION'],  
value=[0,1],inplace=True)
```

```
y_test2 = test_df['loan_status'].values
```

```
y_test2
```

```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,  
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1,  
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1])
```

```
#F1_Score for all of them
```

```
from sklearn.metrics import f1_score
```

```
score_knn = f1_score(y_test2,y_knn2, average='weighted')
```

```
score_dt = f1_score(y_test2,y_dt2, average='weighted')
```

```
score_svm = f1_score(y_test2,y_svm2, average='weighted')
```

```
score_lr = f1_score(y_test2,y_lr2, average='weighted')
```

```
print('KNN f1-Score= ',score_knn)
```

```
print('DT f1-Score= ',score_dt)
```

```
print('SVM f1-Score= ',score_svm)
print('LR f1-Score= ',score_lr)
```

```
KNN f1-Score= 0.7105756358768406
DT f1-Score= 0.7037037037037038
SVM f1-Score= 0.7801458747750308
LR f1-Score= 0.6717642373556352
```

```
#Log loss of the Logistic Regression predictions
from sklearn.metrics import log_loss
log_lr = log_loss(y_test2,lr_prob)
print('LR f1-Score= ',log_lr)

LR f1-Score= 0.47296409967926384
```

Report

The accuracy of the built model using different evaluation metrics are as given in the following table:-

Algorithm	Jaccard	F1-score	LogLoss
KNN	0.72	0.71	NA
Decision Tree	0.70	0.70	NA
SVM	0.79	0.78	NA
LogisticRegression	0.75	0.67	0.47

CHAPTER 6

RESULTS AND DISCUSSIONS

After implementing the loan repayment prediction system and conducting thorough testing, the system's performance and implications can be analyzed. Here's a breakdown of the results and discussion:

1. **Model Performance:**
 - Evaluate the performance of the developed predictive models using appropriate evaluation metrics such as accuracy, precision, recall, F1 score, and area under the ROC curve (AUC-ROC).
 - Compare the performance of different machine learning algorithms and ensemble methods to identify the most effective approach for loan repayment prediction.
 - Provide a detailed analysis of model strengths, weaknesses, and areas for improvement based on validation and testing results.
2. **Accuracy and Robustness:**
 - Assess the accuracy and robustness of the loan repayment prediction system under various scenarios, including different borrower profiles, loan types, and economic conditions.
 - Validate model performance across different datasets, time periods, and geographic regions to ensure generalization ability and reliability in real-world applications.
3. **Feature Importance and Interpretability:**
 - Analyze the importance of individual features and factors influencing loan repayment behavior using feature importance scores, SHAP values, or other interpretability techniques.
 - Identify key drivers of loan repayment and potential risk factors such as credit history, income level, employment status, and demographic characteristics.
4. **User Feedback and Usability:**
 - Gather feedback from users, stakeholders, and domain experts regarding the usability, functionality, and performance of the loan repayment prediction system.
 - Incorporate user feedback into system improvements, feature enhancements, and usability optimizations to enhance user satisfaction and adoption.
5. **Business Impact and Decision Support:**
 - Discuss the potential business impact of the loan repayment prediction system on financial institutions, including improved risk management, reduced default rates, and increased profitability.
 - Highlight the system's role in supporting data-driven decision-making processes for loan approvals, credit scoring, and portfolio management, leading to more informed and objective lending decisions.
6. **Challenges and Limitations:**
 - Address challenges and limitations encountered during system development, testing, and deployment, such as data quality issues, model complexity, and computational resource constraints.

- Discuss potential mitigations and strategies for overcoming challenges to enhance the system's effectiveness and scalability.
7. Future Directions:
- Propose future research directions and opportunities for innovation in loan repayment prediction, such as incorporating alternative data sources, leveraging advanced analytics techniques, and exploring emerging technologies like blockchain and federated learning.
 - Outline potential enhancements and extensions to the existing system, including real-time prediction capabilities, automated decision support, and integration with external data providers and regulatory frameworks.

By presenting comprehensive results and engaging in meaningful discussion, the findings of the loan repayment prediction system can inform stakeholders, guide strategic decision-making, and pave the way for further advancements in credit risk assessment and financial analytics.

Report

Algorithm	Jaccard	F1-score	LogLoss
KNN	0.72	0.71	NA
Decision Tree	0.70	0.70	NA
SVM	0.79	0.78	NA
LogisticRegression	0.75	0.67	0.47

CHAPTER 7

CONCLUSION AND FUTURE ENHANCEMENT

7.1 Conclusion

Nowadays, the loan business becomes more and popular, and many people apply for loans for various reasons. However, there are cases where people do not repay the bulk of the loan amount to the bank which results in huge financial loss. Hence, if there is a way that can efficiently classify the loaners in advance, it would greatly prevent the financial loss.

In this study, the dataset was cleaned first, and the exploratory data analysis and feature engineering were performed. The strategies to deal with both missing values and imbalanced data sets were covered. Then we propose four machine learning models to predict if the applicant could repay the loan, which are Random Forest, Logistic Regression, Support Vector Machine, and K-Nearest Neighbors. When tuning parameters, both Randomized Search Cross Validation and Grid Search Cross Validation methods are applied in different

situations. Through experiments, it is noticed that the model was found which best fits the dataset with highest accuracy is the random forest model, and the model with highest AUC score is Logistic Regression with L2 penalty.

As we expected, borrowers with higher annual income and higher FICO scores are more likely to repay the loan fully; In addition, borrowers with lower interest rates and smaller installments are more likely to pay the loan fully.

7.2 Future Enhancement

In this study, there are several enhancements that we could make in the future. For example, outlier problem is not considered in the exploratory data analysis. if there are outliers in the dataset, the results of the predictive model will not be as valid as they are. In addition, the deep learning algorithm method should also be implemented when predicting for the loan the repayment status. Besides, if we would have a larger dataset, we would have more training samples. By which, it might help fix the high variance problem and make our analysis more valid.

REFERENCES

1. Chen, Y., & Wang, Y. (2020). Loan Default Prediction with Machine Learning Techniques: A Systematic Literature Review. *Journal of Systems Science and Systems Engineering*, 29(3), 302–332.
2. Kim, H., Kim, K. J., & Ahn, H. (2019). Loan Repayment Prediction using Machine Learning Algorithms. *Journal of Computational Science*, 31, 30–41.
3. Brown, L., & Azzi, F. (2018). Predicting Loan Default: A Review of Research and Current Industry Practice. *Journal of the Operational Research Society*, 69(6), 897–910.
4. Quinlan, J. R. (1993). C4.5: Programs for Machine Learning. *Morgan Kaufmann Publishers*.
5. Bishop, C. M. (2006). Pattern Recognition and Machine Learning. *Springer*.
6. Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction. *Springer Science & Business Media*.
7. Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. *O'Reilly Media*.
8. Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. *MIT Press*.
9. Rasmussen, C. E., & Williams, C. K. I. (2006). Gaussian Processes for Machine Learning. *MIT Press*.
10. Hastie, T., Tibshirani, R., & Friedman, J. (2001). The Elements of Statistical Learning: Data Mining, Inference, and Prediction. *Springer Science & Business Media*.

