

Credit Card Default Prediction Model

High Level Design (HLD)

*Deepraj Madhavrao
Vadhwane*

Revision Number 3.0

Revised Date 04/12/2022

Document Version Control

Date Issued	Version	Description
02/12/2022	1.0	Initial HLD
03/12/2022	2.0	Flowchart
04/12/2022	3.0	Final

Contents

Document Version

Control

Abstract

1. Introduction
 - 1.1. Why High-Level Document
 - 1.2. Scope
 - 1.3. Definition
2. General Description
 - 2.1. Problem Statement
 - 2.2. Proposed Solution
 - 2.3. Improvements.
 - 2.4. Requirements
3. Tools
4. Design
 - 4.1. Model Training and Evaluation
 - 4.2. Deployment Process
5. Model Management
 - 5.1. Event Log
 - 5.2. Error Handling
 - 5.3. Performance.
 - 5.4. Reusability
 - 5.5. Compatibility.
 - 5.6. Portability
 - 5.7. Deployment
6. Conclusion.

Abstract

In present day majority of people are using credit card for online shopping, EMI payments. Moreover banks are encouraging customers to use it by approving them freely without any initial charges. Because credit card interest income is the major share of profit for all commercial banks. Even if customers are not willing to buy, bank marketing team reaches out through phone call to sell it. At first the interest rate are normal to grope the customer into the practice of using it, later they will levy charges on everything to loot money from the customers. Due to this reason many customers are facing difficulty in repaying the credit amount on time. If they missed the deadline the interest rate will get multiplied leading to defaulting the account if left unpaid even for a short span of time. The model we are going to build will consider all inputs from user to predict whether the account will be defaulted or not based on various parameters.

1. Introduction

1.1 Why High Level Design?

The main purpose of HLD is to add the necessary details to the current project description to represent a suitable model for coding. This document is also intended to help detect contradictions prior to coding, and can be used as a reference manual for how the modules interact at a high level. Document includes the following aspects:

- ✓ Architecture Design and Definition
- ✓ User Interface Implementation
- ✓ Hardware/Software Interface
- ✓ Performance Requirements
- ✓ Non – functional attribute definitions
 - ✓ Security, Reliability, Maintainability, Portability, Reusability, Application Compatibility, Resource Utilization, Serviceability

1.2 Scope

Documentation provides the system structure, such as DB architecture, application architecture (layers), application flow (Navigation), and technology architecture. The HLD uses non– technical terms which should be understandable to the administrator.

1.3 Definitions

Term	Description
CCDD	Credit Card Default Detection Model
CIBIL	Credit Information Bureau Limited
IDE	Integrated Development Editor
MLA	Machine Learning Algorithm
EDA	Exploratory Data Analysis
ETL	Extract, Transform, Load

2. General Description

2.1 Problem Statement

To create machine learning model to detect anomalies in credit repayment transaction and predict whether the account will be defaulted or not.

2.2 Proposed Solution

The proposed model will be capable of studying customers based on various parameters like gender, education qualification, marital status, previous repayment history and current outstanding balance. Based on that it can predict the account default status in advance and notify the respective bank official for further action.

2.3 Further Improvements

Further we can enhance our model with additional information like customer expenditure pattern, CIBIL score, loan EMI or availability of any fixed deposit mapped to the account. These informations will improve the prediction accuracy and can suggest alternative methods to retrieve the pending dues in-case of future default.

2.4 Requirements

Model performance depends on various factors like dataset we use to train it, prediction algorithm, hyper tuning parameters

- ✓ Input Dataset – Data should be clean without any outliers, missing values, garbage values
- ✓ Outliers – can greatly impact the model performance
- ✓ Missing/Garbage values – Data collected often have missing and garbage values which has to be cleaned before training the model
- ✓ Choice of algorithm – Depending on the problem statement appropriate algorithm is used to train

3. Tools

Python programming language is our primary tool to build the model with few open-source libraries for various subtasks

3.1 Pandas

Python package providing fast, flexible and expressive data structure designed to make working with relational or labelled data both easy and intuitive

3.2 Numpy

Python library used for working with arrays with domain specific functionalities in Linear algebra, Fourier transform and Matrices.

3.3 Scikit-Learn

High level framework designed for supervised and unsupervised machines learning algorithms. It provides efficient tools for statistical and scientific models for illustrations.

3.4 Plot Visualization

Seaborn and Matplotlib libraries used to create 2D graphs and plots to visualize the data

3.5 HTML/CSS

HTML is used to create static web pages and web application whereas CSS is a style sheet responsible for presentation of documents written in HTML

3.6 Flask

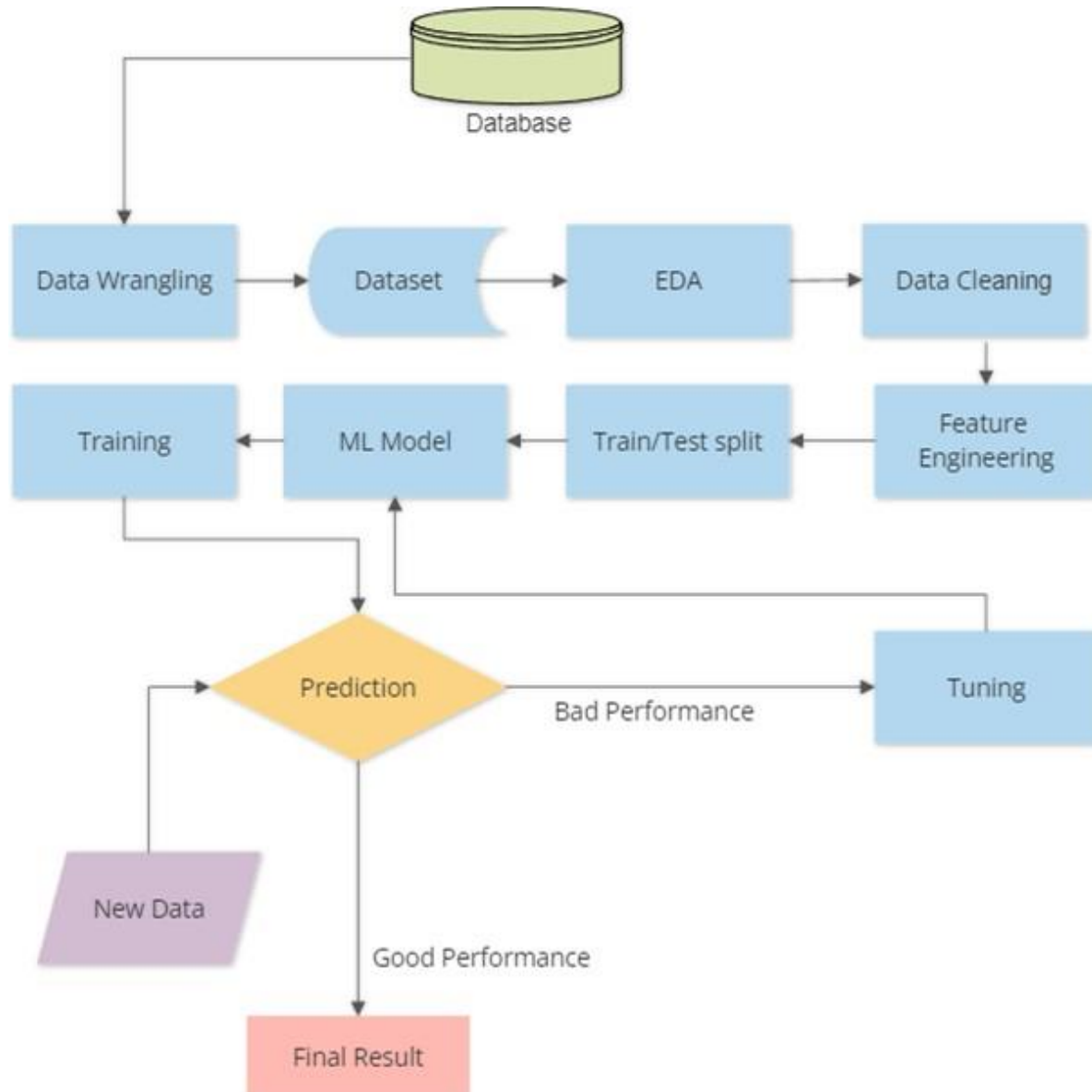
Small lightweight python web framework that provides useful tools and features to make API based web applications.

3.7 Github

Cloud based online software development platform used for storing, tracking, and collaborating on software projects. It enables developers to upload their own code files and to collaborate with fellow developers on open-source projects

4 Design Architecture

4.3 Model Training and Evaluation

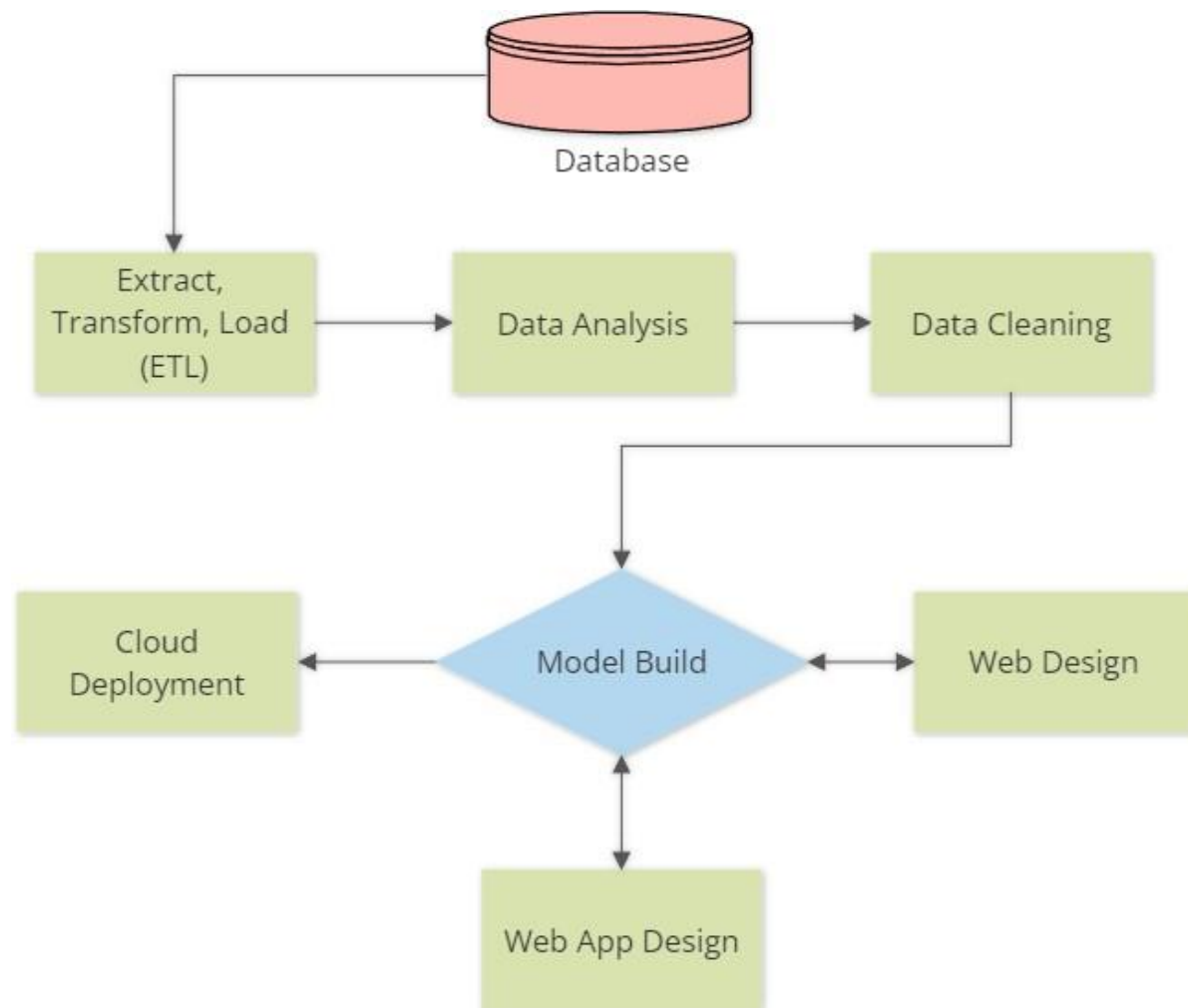


Data Wrangling – Process of transforming and mapping the raw data into a desirable format with the intent of making it more valuable dataset

EDA – Analysis of data to summarize their main characteristics often using statistical graphs and visualization methods

Feature Engineering – addition and construction of additional features to improve model performance and accuracy

4.4 Deployment Process



ETL – Three phase process where data is collected from multiple sources, transformed into desired format and loaded into multiple output containers

Web Design – Basic HTML/CSS design for the web application

Web App – Front end app interface using Flask python library

5 Model Management

5.3 Event Log

Model should log every event from start to end for debugging purpose. We can spot the error in logs and identify which part of the code need to be addressed and corrected. Logging is essential to see the internal execution flow and can be used to improve the architecture design for faster execution.

5.4 Error Handling

Error Handling makes a model robust, reliable ensuring the execution won't stop abruptly. Even in the worst case scenario the execution will terminate in between intimating the user the reason for termination. We should imagine all possible worst case scenarios and able to code the handlers to terminate safely or jump a step ahead in the execution after intimating the user.

5.5 Performance

As we discussed before model performance entirely depends on training dataset, quality of the data, authenticity of the data, algorithm used and hyper-parameter tuning. But can assure the performance will be good assuming the above criteria are all satisfied.

5.6 Reusability

All codes are written in modular fashion so that for future improvements codes can be reused or modified without affecting other modules. Improvements are necessary for better performance, faster execution and updated dataset training for the model.

5.7 Compatibility

No model is good if it not compatible. Model should be build in such a way it can perform in all scenarios provided with required input details. For this reason we need to use widely prevalent version of framework or library supporting maximum systems to run the model.

5.8 Portability

Models build in containers have high portability. Portable containers will make sure the model executes and performs well even if the system doesn't have necessary libraries pre-installed.

5.9 Deployment

Deploying the model in cloud platform makes it readily available for action. We can run and extract the results any time. Cloud platform will provide the necessary hardware and software support for the execution.

6 Conclusion

We have build multiple classification models, trained and tested with same datasets, and finally performed parameter tuning for decently performing model to achieve best prediction and recall accuracy.