# FIT5145: ASSIGNMENT-3

Predictive Modelling for Loan Approval

Bhavna Balakrishnan
33954437

# Project Description

The aim of this project is to develop a predictive model that will help determine the eligibility of loan applicants based on various criteria such as income, employment status etc. The project will commence by gathering insights on applicant demographics and conducting exploratory data analysis that will assist in identifying trends in the data. Following this, regression and other tests will be conducted to establish a predictive model.

## Business Model

1. Application Areas:

- Financial Services Sector: The project is specifically designed for financial institutions such as banks, credit unions, and other money lending organizations. These entities require efficient and accurate assessment tools to evaluate the creditworthiness of loan applicants.
- Regulatory Compliance and Risk Management: Besides direct financial services, the model can also serve regulatory bodies and risk management departments by providing tools that ensure compliance with financial regulations and manage lending risks effectively.

2. Benefits or Values Created:
- Efficiency and Accuracy: By automating the loan eligibility process, the predictive model reduces manual processing time and errors, leading to faster and more accurate loan decisions.
- Cost Reduction: Automating processes typically results in significant cost savings in terms of human labour and the associated expenses of manual errors and delays.
- Risk Mitigation: The model helps in identifying potential defaulters early in the loan application process, thereby reducing the likelihood of bad debt and improving the overall credit risk profile of the portfolio.
- Regulatory Compliance: The model ensures that all loan approval processes are in compliance with existing financial regulations, helping institutions avoid legal and reputational risks.

- Improved Customer Experience: Customers benefit from quicker loan processing times and a more transparent assessment process, which can improve their overall satisfaction and trust in the financial institution.
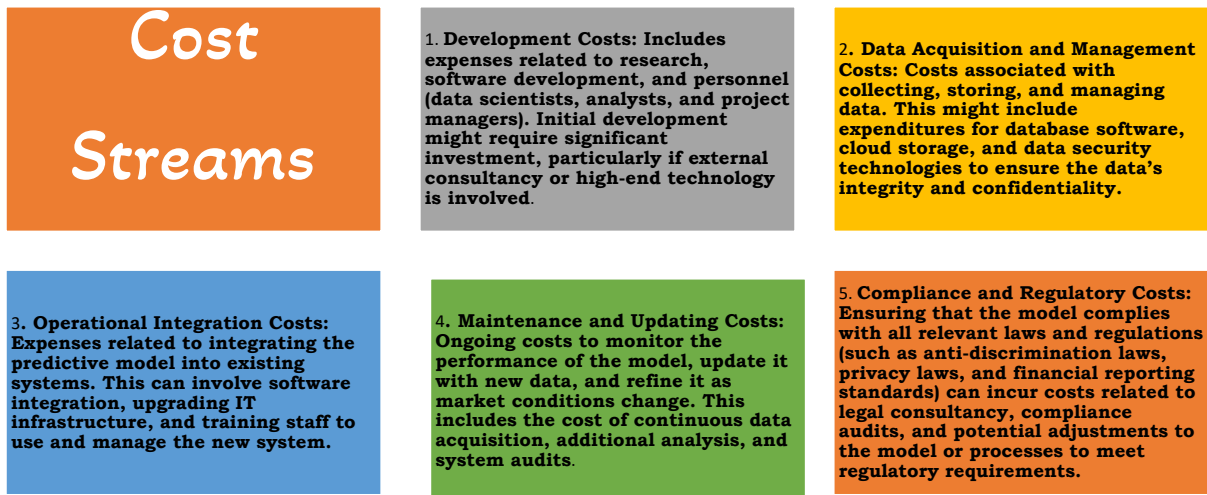
3. Stakeholders Who Can Benefit:

- Financial Institutions: Banks, credit unions, and other lenders can streamline their loan processing workflows, enhance customer satisfaction, and ensure compliance with regulatory standards.
- Customers: Loan applicants receive faster responses and fair treatment due to the unbiased nature of the automated decision-making process.
- Regulators: Regulatory bodies can rely on the use of such models to ensure that financial institutions are complying with guidelines and managing risks appropriately.

4. Challenges:

- Data Quality and Availability: High-quality, accurate, and comprehensive data is crucial for building effective predictive models. The availability and integrity of data can be a significant hurdle.
- Model Bias and Fairness: Ensuring the model does not perpetuate existing biases and discriminates against certain groups of people is a critical challenge. The model must be regularly audited for fairness.
- Scalability and Reliability: The model must be scalable to handle large volumes of applications and reliable enough to maintain performance over time without frequent downtimes or errors.
- Technical Expertise: Developing, deploying, and maintaining predictive models require a high level of technical expertise in data science and machine learning, which can be a resource constraint.
- Regulatory Compliance: Navigating the complex landscape of financial regulations and ensuring the model complies with all relevant laws and guidelines is an ongoing challenge.

By addressing these points, the business model section of the proposal would clearly articulate where the project fits within the broader industry, the specific benefits it offers, who it serves, and the obstacles it faces. This level of detail will likely enhance stakeholder understanding and support for the project.

## High-level Overview

| | | |
|---|---|---|
| **Cost Streams** | 1. **Development Costs: Includes expenses related to research, software development, and personnel (data scientists, analysts, and project managers). Initial development might require significant investment, particularly if external consultancy or high-end technology is involved**. | 2. **Data Acquisition and Management Costs: Costs associated with collecting, storing, and managing data. This might include expenditures for database software, cloud storage, and data security technologies to ensure the data's integrity and confidentiality.** |
| 3. **Operational Integration Costs: Expenses related to integrating the predictive model into existing systems. This can involve software integration, upgrading IT infrastructure, and training staff to use and manage the new system.** | 4. **Maintenance and Updating Costs: Ongoing costs to monitor the performance of the model, update it with new data, and refine it as market conditions change. This includes the cost of continuous data acquisition, additional analysis, and system audits**. | 5. **Compliance and Regulatory Costs: Ensuring that the model complies with all relevant laws and regulations (such as anti-discrimination laws, privacy laws, and financial reporting standards) can incur costs related to legal consultancy, compliance audits, and potential adjustments to the model or processes to meet regulatory requirements.** |

| **Revenue Generation Streams** | | |
|---|---|---|
| | **Increased Loan Disbursements** | **Efficiency**: By automating and accelerating the loan approval process, the institution can handle a higher volume of loan applications without a corresponding increase in staff or resources. |
| | | **Scale**: As the process becomes more streamlined, the institution may expand its reach, potentially entering new markets or customer segments previously considered too risky without robust predictive capabilities. |
| | **Interest Revenue** | **Higher Approval Rates**: Improved predictive accuracy means more loans can be approved with confidence, increasing the overall loan portfolio. |
| | | **Risk-Adjusted Pricing**: The model can help in better assessing the risk associated with each loan, allowing for more accurately priced loans in terms of interest rates, which can increase profitability on a per-loan basis. |
| | **Customer Acquisition and Retention** | **Attracting New Customers**: An efficient, reliable loan approval process can be a strong selling point in attracting new customers who value quick and fair service. |
| | | **Customer Loyalty**: Customers who have a positive borrowing experience are more likely to return for future financial needs and refer others, growing the customer base organically. |
| | **Cross-Selling and Upselling** | **Data Insights**: The analytics used in the predictive model can reveal insights about customers' financial behaviors and needs, providing opportunities to offer tailored financial products such as overdrafts, credit cards, or investment products. |
| | | **Targeted Offers**: With a deeper understanding of the customer base, the institution can design and target offers more effectively, increasing the uptake of additional services. |
| | **Risk Reduction and Default Minimization** | **Reduced Defaults**: By better predicting which applicants are likely to default, the institution can avoid costly bad debts, maintaining a healthier balance sheet and reducing the need for debt recovery actions. |
| | | **Capital Efficiency**: Lower default rates mean that less capital needs to be held in reserve against possible loan losses, freeing up capital for other revenue-generating activities. |
| | **Regulatory Compliance and Reputation** | **Compliance**: Efficient compliance through automated systems can reduce the risk of fines or sanctions from failing to meet regulatory standards, indirectly protecting revenue. |
| | | **Brand Value**: Being known for using advanced, fair, and transparent lending practices enhances the institution's reputation, which can be a significant competitive advantage in attracting quality borrowers and investors. |

## Roles

A Data Analyst handles data collection, cleaning, and exploration. A Data Scientist designs predictive models like logistic regression for tasks like loan eligibility. A Data Engineer supports pre-processing and optimization.

## Potential Data Sources

- Financial Records: Data from bank statements, credit card histories, and previous loan records which can provide insights into the financial behaviour of applicants.
- Employment Information: Details about employment status, duration, and income level from databases maintained by employers or collected directly from applicants.
- Credit Bureaus: Comprehensive reports from credit bureaus that include credit scores, existing debts, and credit history.
- Public Records: Information from public records like bankruptcy filings, legal judgments, or tax liens that might impact creditworthiness.
- Online Behaviour Data: Data from online activities that may indirectly reflect financial behaviour, such as spending habits or stability in financial planning, if privacy laws allow.

Realistically speaking, if the project to be released in real time, then the data source would be a live data base of real time applicants as mentioned previously in Financial Records.

For this assignment I will be using a data source from Kaggle, licensed by MIT from the below source-

MIT (n.d.). *Loan-Approval-Prediction-Dataset*. Kaggle.

https://www.kaggle.com/datasets/architsharma01/loan-approval-prediction-dataset?resource=download

This dataset contains mock data that has data about applicant graduation, income and assets as well as approval status.
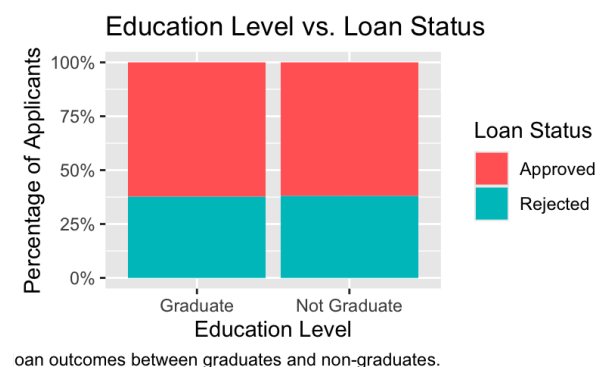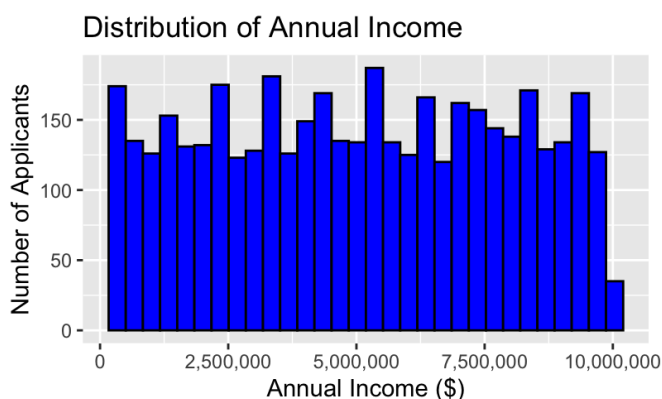
## Characteristics of the Data (The 4 V's)

- Volume: For such a project, the volume would be tremendous in realistic scenarios. For the purpose of this assignment, my dataset has about 4300 rows.
- Velocity: Data flow will be continuous as new applications come in and additional data (like credit scores) updates regularly.
- Variety: Data will include a mix of structured data (e.g., income levels, credit scores) and unstructured data (e.g., employment history narratives).
- Veracity: The accuracy and truthfulness of the data can vary, particularly with self-reported data, requiring verification processes.

## Required Platforms, Software, and Tools

- Data Storage and Management: Use of relational databases like MySQL for structured data and NoSQL databases like MongoDB for unstructured data.
- Data Processing: Apache Hadoop for handling large volumes of data, Apache Spark for fast processing tasks.
- Data Analysis Software: Python and R for statistical analysis, with libraries like Scikit-learn for machine learning and Pandas for data manipulation.
- Data Visualization: Tools like Tableau or PowerBI to visualize data trends and insights for stakeholders.
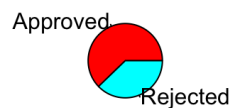
**DEMONSTRATION**

**Demographics**



oan outcomes between graduates and non-graduates.
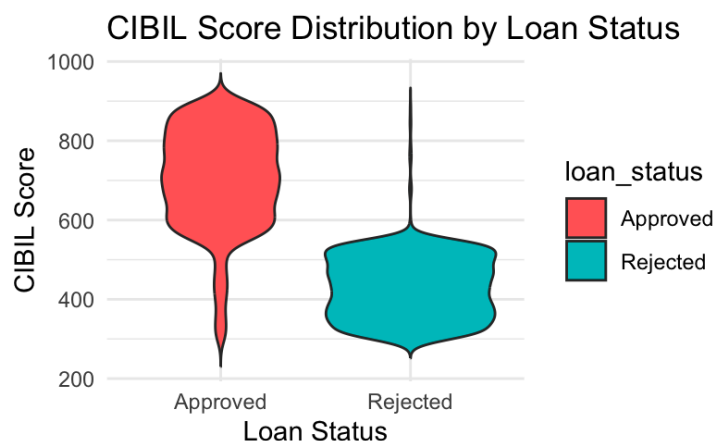
Distribution of Loan Amount by Number

From above, we can conclude that applicants have varied annual income ranges but there appears to be no great difference between their education levels or any significant relationship between loan amount requested and number of dependents. Furthermore, we can explore what the success rate of getting a loan is as seen below-
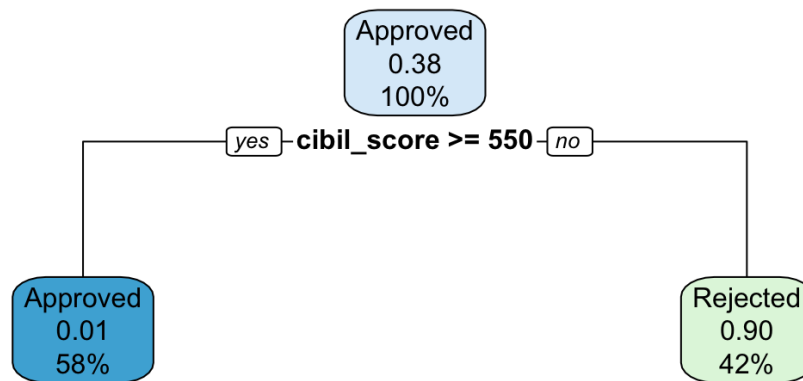


**Loan Status Distribution**

We can see that a lot of the loan applicants are having their loans approved from the above pie chart.
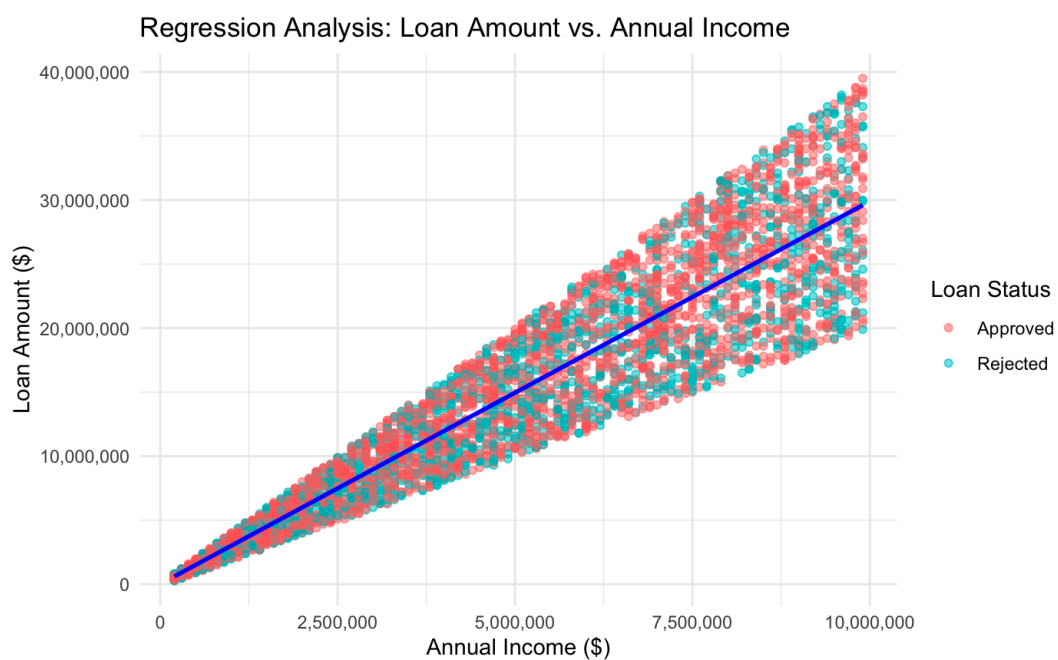


CIBIL Score Distribution by Loan Status

From above, we can also conclude that the higher an applicant's CIBIL score is the more chances they have of having their loan approved. This is logical since a credit score is a good indicator of the goodwill of an applicant to make timely and correct payments for their credit cards which makes them less likely to default and become credit risks to a lender. This would

therefore make them more eligible for sanctioning loans. Below we can see a decision tree depicting the same-



Therefore, we can say that applicants with scores equal to or higher than 550 have 58% chance of getting loans approved.

## Regression Testing

```
Signif. codes:  0 '***' 0.001 '**'
0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3382000 on
4267 degrees of freedom
Multiple R-squared:  0.8602,
Adjusted R-squared:  0.8602
F-statistic: 2.626e+04 on 1 and 4267
DF,  p-value: < 2.2e-16
```

The regression model analysing the relationship between annual income and loan amount is highly effective, explaining a significant proportion (86.02%) of the variability in loan amounts. The model's predictions vary from actual values by an average of about $3,382,000, which may be notable, suggesting further refinement could be beneficial depending on the application's precision needs. The strong statistical significance of the model indicates that annual income is a crucial predictor of loan amount. Financial institutions can rely on this model for a robust estimation of loan amounts based on applicants' incomes, though they should also consider other factors potentially influencing loan decisions due to the observed residual error.

## Standard for Data Science Process

CRISP-DM (Cross-Industry Standard Process for Data Mining) is commonly used in data science projects. This standard provides a structured approach to planning, implementing, and reviewing the process of data mining, including predictive modelling. The stages of CRISP-DM are:

1. Business Understanding: Define project objectives and requirements from a business perspective. For this project:
    a. Objective: Develop a predictive model to enhance loan approval accuracy, streamline processes, and ensure compliance with financial regulations.
    b. Key Goals: Improve decision accuracy, automate processes, and reduce bias.

2. Data Understanding: Start collecting data and proceed with activities to get familiar with the data, identify data quality issues, discover first insights into the data, or detect interesting subsets to form hypotheses for hidden information.
    a. Activities: Collect historical loan data, inspect for quality and completeness, and perform preliminary statistical analyses to identify trends and data quality

issues.

3. Data Preparation: Data is cleaned and formatted; it involves dealing with missing values, encoding categorical variables, and possibly scaling and normalizing data.
   a. Tasks: Clean data by addressing missing values and outliers, encode categorical variables, and normalize numerical values to prepare the dataset for modelling.

4. Modelling: Various modelling techniques are selected and applied, and their parameters are calibrated to optimal values.
   a. Approach: Evaluate and select the best algorithms (e.g., logistic regression, decision trees), tune parameters, and conduct feature selection to build an optimal predictive model.

5. Evaluation: Evaluate the model, reviewing the steps executed to construct the model and be certain it properly achieves the business goals.
   a. Validation: Assess the model using metrics such as accuracy, ROC-AUC, and confusion matrix. Utilize techniques like cross-validation to ensure robustness.
   b. Adjustment: Refine the model based on evaluation results and business feedback to align closely with operational goals.

6. Deployment: Depending on the requirements, the deployment phase can be as simple as generating a report or as complex as implementing a repeatable data scoring or data mining process.

## Data Governance and Management

Data Governance encompasses the overall management of the availability, usability, integrity, and security of the data employed in an enterprise, with the following practices being critical:

1. Accessibility:
- Role-based Access Control (RBAC): Ensures that only authorized personnel have access to specific levels of data based on their roles within the organization.

- Audit Trails: Keep logs of who accessed the data and when, which aids in maintaining a clear record of data usage.

2. Security:

- Encryption: Data should be encrypted both at rest and in transit to protect sensitive information from unauthorized access.

- Secure Backup Systems: Regularly updated backups that are stored securely to prevent data loss and allow recovery in case of data corruption or loss.

3. Confidentiality:

- Data Anonymization: Removing personally identifiable information from the data sets to protect individual privacy before processing the data.

- Data Masking: Implementing masking techniques to hide data, which ensures that confidentiality is maintained while data is in use.

4. Ethical Concerns:

- Bias and Fairness: Monitor models to ensure they do not propagate or amplify biases, which could lead to unfair treatment of individuals based on prohibited characteristics.

- Transparency: Maintain transparency about how data is collected, stored, used, and shared, especially when using data for predictive modelling in sensitive areas like loan approval.

- Compliance with Regulations: Adhere to applicable laws and regulations such as GDPR, HIPAA, or others that apply to data privacy and protection.

## Implementing Data Governance:

- Implement robust data handling and security measures to protect sensitive financial information.
- Ensure that the data collection and modelling processes are compliant with financial regulations and ethical standards.
- Conduct regular audits and reviews of data use and model performance to identify potential issues in data handling or model bias.

## Conclusion

The predictive modelling project for loan approval aims to refine the loan decision-making process, enhancing both efficiency and accuracy. By employing sophisticated data analytics, this project significantly reduces manual processing time, minimizes the risk of defaults, and boosts customer satisfaction through quicker and fairer loan approvals. The implementation of this model presents an initial financial outlay and requires ongoing investment in data management and system updates. However, the potential for increased loan volume, improved customer retention, and cross-selling opportunities offers a promising avenue for substantial long-term revenue growth. Reflecting on this Endeavor, it is clear that while the challenges are non-trivial, the strategic integration of this predictive model could be transformative, offering competitive advantages and aligning with modern, data-driven business practices in the financial sector.

# References

(n.d.). *CRISP-DM Help Overview*. IBM. https://www.ibm.com/docs/en/spss-modeler/saas?topic=dm-crisp-help-overview

(n.d.). *Imbalanced Data : How to handle Imbalanced Classification Problems*. https://www.analyticsvidhya.com/blog/2017/03/imbalanced-data-classification/

Deloitte (n.d.). *Model Risk Management*. https://www2.deloitte.com/content/dam/Deloitte/fr/Documents/risk/deloitte_model-risk-management_plaquette.pdf

MIT (n.d.). *Loan-Approval-Prediction-Dataset*. Kaggle. https://www.kaggle.com/datasets/architsharma01/loan-approval-prediction-dataset?resource=download

bankwest (n.d.). *Understanding the home loan approval process*. https://www.bankwest.com.au/personal/home-buying/guides/home-loan-approval#:~:text=Look%20into%20your%20credit%20score,credit%20cards%2C%20bills%20or%20loans.

Link to dataset:

https://drive.google.com/file/d/1ROUdDVUS3oKYMGocSZN1Hgekk5KyO1vb/view?usp=sharing