Solution Step 1: Exploratory Data Analytics

1. **What is the shape of the dataset?**
   Shape of the dataset: (5960, 13)

2. **Are there any missing values in the dataset?**

   Yes, there is missing values in the dataset.
   Missing values in the dataset:
   BAD          0
   LOAN         0
   MORTDUE    518
   VALUE      112
   REASON     252
   JOB        279
   YOJ        515
   DEROG      708
   DELINQ     580
   CLAGE      308
   NINQ       510
   CLNO       222
   DEBTINC   1267
   dtype: int64

3. **How does the distribution of years at present job "YOJ" vary across the dataset?**

   - Loan applicants with a lower number of years at their present job are more common than those with a longer job tenure. For example, there are 415 applicants who have just started their current job (0 years), which is the largest group of applicants.

   - As the number of years at the present job increases, the frequency of loan applicants tends to decrease. This pattern suggests that there are fewer loan applicants with a longer job tenure in the dataset.

   - The distribution of YOJ values appears to be positively skewed, with the majority of the data concentrated at the lower end of the range. This indicates that most loan applicants in the dataset have relatively shorter job tenures.

4. **Is there a relationship between the REASON variable and the proportion of applicants who defaulted on their loan?**

The relationship between the REASON variable and the proportion of applicants who defaulted on their loan:

- For loan applicants who requested a loan for debt consolidation (DebtCon), the proportion of applicants who defaulted on their loan is approximately 18.97%.
- For loan applicants who requested a loan for home improvement (HomeImp), the proportion of applicants who defaulted on their loan is approximately 22.25%.
- These results suggest that there is a slightly higher proportion of applicants who defaulted on their loans when the reason for the loan is home improvement compared to those who took the loan for debt consolidation.

5. **Do applicants who default have a significantly different mortgage amount compared to those who repay their loan?**

- The mean mortgage amount for defaulted loans is approximately 69,460.45.
- The mean mortgage amount for repaid loans is approximately 74,829.25.
- A t-test was performed to compare the mean mortgage amounts for the two groups, and the results show a t-test statistic of -3.377 and a p-value of 0.00075. Since the p-value is less than the commonly used significance level (alpha) of 0.05, we can conclude that there is a statistically significant difference between the mean mortgage amounts for defaulted loans and repaid loans.
- In summary, applicants who default on their loans have, on average, a significantly lower mortgage amount compared to those who successfully repay their loans.

## Solution Step 2: Data Preprocessing:

6. **Are there any patterns in the missing values for particular variables?**

- The DEBTINC variable has the highest number of missing values (1,267). It might indicate that applicants are less likely to provide their debt-to-income ratio or this information is harder to collect by loan providers.

- The variables MORTDUE, YOJ, DEROG, DELINQ, CLAGE, and NINQ have a significant number of missing values. These variables are related to the applicant's credit history, and it might suggest that this information is either difficult to collect or some applicants have not provided the information.

- The categorical variables REASON and JOB have 252 and 279 missing values, respectively. This may imply that some applicants have not provided this information, or there were issues during data collection.

7. **How to impute missing values for numerical and categorical variables?**

- Numerical variables: For numerical variables, I created a SimpleImputer object with the strategy set to 'mean'. This object will be used to impute missing values with the mean value of each numerical variable. The list of numerical variables includes 'LOAN', 'MORTDUE', 'VALUE', 'YOJ', 'DEROG', 'DELINQ', 'CLAGE', 'NINQ', 'CLNO', and 'DEBTINC'. We use the fit_transform() function to impute the missing values in the DataFrame for these variables.

- Categorical variables: For categorical variables, I created a SimpleImputer object with the strategy set to 'most_frequent', which means it will use the mode (most frequent value) for imputing missing values. The list of categorical variables includes 'REASON' and 'JOB'. We use the fit_transform() function to impute the missing values in the DataFrame for these variables.

8. **Do we need to normalise the data before splitting into train and test set?**

- Yes we need it.Normalizing the data can help improve the performance of some machine learning models that are sensitive to the scale of input features. It's recommended to normalize the data before splitting it into train and test sets to ensure that both sets have the same scale.

## Solution Step 3: Model Building:

9. **Which algorithm can used to predict the target variables?**
   - I will use logistic regression as a linear algorithm and random forest as a non-linear algorithm.

10. **Is there any significant difference between the results linear and non linear algorithms**

- Yes, there is a significant difference between the results of the linear and non-linear algorithms on this dataset. The logistic regression, which is a linear algorithm, achieved an accuracy of approximately 80.96%. In contrast, the random forest, which is a non-linear algorithm, achieved a higher accuracy of about 91.44%.

- The non-linear random forest algorithm outperforms the linear logistic regression algorithm by around 10 percentage points in accuracy, indicating that the random forest model is better at capturing the underlying patterns and relationships in the

dataset. This performance difference suggests that a non-linear algorithm might be more suitable for this particular dataset than a linear one.

**11. Does linear algorithm is easy to interpret as compared to non linear algorithms**

- Yes, generally, linear algorithms like logistic regression are easier to interpret than non-linear algorithms like random forest. The reason for this is that linear algorithms have a simpler structure and fewer parameters, making it easier to understand how the algorithm is making predictions. The coefficients of the features in a logistic regression model, for example, can provide insight into how each feature affects the target variable.

- In contrast, non-linear algorithms like random forest have more complex structures and parameters that interact with each other in more intricate ways. This makes it more difficult to understand how the algorithm is making its predictions. However, even though non-linear algorithms are generally harder to interpret, they can often achieve higher accuracy than linear algorithms on complex datasets

**12. How to choose the parameters to optimise the algorithm**

- Choosing the optimal parameters for a machine learning algorithm is an essential step in achieving the best possible performance. There are several techniques you can use to tune the hyperparameters of your algorithm, including:

- Grid Search: In grid search, you specify a set of possible values for each hyperparameter, and the algorithm tries all possible combinations of these values to determine the optimal set of hyperparameters that maximize performance.

- Random Search: In random search, you randomly select values for each hyperparameter from a specified range of possible values. This technique can be more efficient than grid search when there are many hyperparameters to search over.

**What is the metric (Measure of Success) for this business problem?**

- The measure of success for this business problem is the accuracy of the predictive model in identifying clients who are likely to default on their home equity loans.

The ultimate goal is to minimize the number of defaulters and reduce the financial loss associated with defaults.

- In this case, accuracy can be defined as the proportion of correctly identified defaulters and non-defaulters in the test dataset.

**Should we consider the trade-off between the interpretability and model performance to choose the best model**

- Yes, we should consider the trade-off between the interpretability and model performance when choosing the best model. A more interpretable model is often easier to understand and explain to stakeholders, but it may sacrifice some level of performance compared to a more complex, less interpretable model. On the other hand, a more complex and less interpretable model may achieve better performance on the prediction task but may be difficult to understand and explain to stakeholders.

- The choice between interpretability and performance should be based on the specific goals and requirements of the business problem. For example, if the goal is to understand the factors that contribute to loan defaults and to take action to mitigate these risks, a more interpretable model such as logistic regression may be a better choice, even if it sacrifices some level of performance. Conversely, if the goal is to achieve the highest possible prediction accuracy, a more complex model such as random forest or neural networks may be a better choice, even if it sacrifices some level of interpretability.

**What are the refined insights from EDA and model building?**

The refined insights from exploratory data analysis (EDA) and model building for the loan default prediction problem can be summarized as follows:

**EDA insights:**

- There are missing values in the dataset, which were imputed using appropriate methods such as mean imputation for numerical variables and mode imputation for categorical variables.
- The dataset is imbalanced, with only 20% of clients experiencing a loan default. This should be taken into account when selecting evaluation metrics and building the model.
- The mortgage amount appears to be a significant factor in predicting loan defaults, with defaulted clients having a lower average mortgage amount compared to those who repaid their loans.
- The distribution of years at present job (YOJ) varies across the dataset, with a majority of clients having less than 10 years of experience at their current job.

**Model building insights:**

- Random forest model achieved a higher accuracy compared to logistic regression model in predicting loan defaults.
- The precision and recall for identifying loan defaults were higher in random forest compared to logistic regression.
- There is a trade-off between model interpretability and performance, with logistic regression being more interpretable but less accurate compared to random forest.
- Overall, the refined insights from EDA and model building suggest that the mortgage amount is an important factor in predicting loan defaults, and that random forest is a better choice for building a predictive model for this problem. The insights also suggest that there is a trade-off between model interpretability and performance, and that the choice between these factors should be made based on the specific requirements and constraints of the business problem.

**What observations and insights can be drawn from the confusion matrix and classification report?**

The confusion matrix and classification report provide insights into the performance of the models in predicting loan defaults. The classification report shows the precision, recall, and F1-score for each class (0 and 1), as well as the overall accuracy of the model. The confusion matrix shows the number of true positives, true negatives, false positives, and false negatives for each model.

- From the classification reports, we can see that both models have higher precision and recall for loans that will be repaid (class 0) compared to loans that will default (class 1). This suggests that the models are better at predicting loans that will be repaid compared to loans that will default.
- The accuracy of the logistic regression model is 0.81, which is lower than the accuracy of the random forest model (0.91). This suggests that the random forest model is better at predicting the target variable than the logistic regression model.
- The F1-score is a weighted harmonic mean of precision and recall, and is often used as a measure of the overall performance of a classification model. In this case, the F1-score for class 1 is lower than the F1-score for class 0 in both models, suggesting that the models are more accurate at predicting loans that will be repaid compared to loans that will default.
- Looking at the confusion matrices, we can see that the random forest model has more true positives and fewer false negatives than the logistic regression model. This suggests that the random forest model is better at predicting loans that will default than the logistic regression model.

In summary, both models have higher precision and recall for loans that will be repaid, and are less accurate at predicting loans that will default. The random forest model has a higher accuracy than the logistic regression model and is better at predicting loan defaults, with a higher number of true positives and lower number of false negatives.

**Is the model performance good enough for deployment in production?**

- Whether or not a model is good enough for deployment in production depends on the specific requirements of the business problem at hand, as well as the acceptable level of risk and error for the task.
- In this case, the accuracy of the random forest model is 0.91, which is relatively high and suggests that the model can make accurate predictions for the target variable. However, it's important to note that the model is more accurate at predicting loans that will be repaid compared to loans that will default, which may be a concern depending on the specific business problem.
- Additionally, it's important to consider the potential consequences of false positives and false negatives, as well as the potential impact of model errors on the overall business objective. For example, if the business objective is to minimize losses from loan defaults, false negatives (where the model incorrectly predicts that a loan will be repaid when it actually defaults) could have a significant impact on the business outcome.
- Therefore, before deploying the model in production, it's important to conduct a thorough evaluation of the model performance and consider the potential consequences of model errors in the specific business context. It may also be necessary to continue to monitor and refine the model over time to ensure that it continues to perform well and meets the evolving needs of the business problem.

**What is proposal for final solution design? What are expected benefits and costs (assume numbers) of this solution design?**
- Exploratory Data Analysis (EDA): Analyze the data to gain insights into the relationships between variables and the overall patterns of the data. This can be done through visualizations and statistical analysis.
- Data Preprocessing: Clean and preprocess the data, including handling missing values and outliers, encoding categorical variables, and scaling numerical variables.
- Model Building: Train and evaluate different machine learning models such as logistic regression, decision trees, random forests, or neural networks to identify the best model for the task. This will involve tuning hyperparameters and optimizing the model for performance.
- Model Evaluation: Evaluate the performance of the selected model using appropriate evaluation metrics such as accuracy, precision, recall, and F1 score.