# LOAN STATUS PREDICTION

## Introduction :

This report provides an overview of the performance of the Loan Status Prediction Model developed to assist in the decision-making process for approving or rejecting loan applications. The model was trained using a dataset comprising various applicant and property metrics.

## Data Overview :

The dataset includes information on applicants who have previously applied for property loans. Key features include Applicant Income, Loan Amount, Credit History, Co-applicant Income, and others.

## Model Details :

Several classification models were evaluated, including Logistic Regression, Random Forest, Support Vector Classifier , KNN  and XGBoost. The models were trained on a subset of the data and tested on a separate set to assess their performance.

## Model Performance :

| Model | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|
| XGBClassifier | 0.75 | 0.78 | 0.89 | 0.83 |
| Random Forest | 0.81 | 0.79 | 0.99 | 0.88 |
| Logistic Regression | 0.78 | 0.77 | 0.96 | 0.86 |
| KNeighborsClassifier | 0.65 | 0.70 | 0.84 | 0.77 |
| Support Vector Classifier | 0.68 | 0.38 | 1.00 | 0.81 |

## Key Insights

1. High Recall Models:
   - The **Random Forest** and **Logistic Regression** models have very high recall scores of **0.99** and **0.96**, respectively. This indicates that these models are particularly good at identifying applicants who are likely to get their loans approved, with fewer false negatives.
   - The **Support Vector Classifier** has a perfect recall score of **1.00**, suggesting it identifies all positive cases, but its low precision score of **0.38** indicates a high number of false positives.

2. Balanced Performance:
    ○ The **XGBClassifier** shows a balanced performance with an F1 score of **0.83**, indicating a good balance between precision and recall. This model is less likely to make mistakes on either side of the prediction spectrum.
3. Precision vs. Recall Trade-off:
    ○ There is a trade-off between precision and recall observed in the models. For instance, while the **Support Vector Classifier** has a high recall, its precision is quite low, which might not be ideal for the bank if the cost of a false positive is high.
4. Overall Accuracy:
    ○ The **Random Forest** model has the highest overall accuracy of **0.81**, making it the most accurate model among those tested. However, accuracy alone doesn't tell the full story, as it doesn't account for the class imbalance that is often present in loan datasets.
5. Potential for Model Improvement:
    ○ The **KNeighborsClassifier** has the lowest scores across all metrics, which suggests that this model may not be capturing the complexities of the dataset well. Feature engineering, more complex models, or additional data might be needed to improve its performance.
6. Model Selection for Deployment:
    ○ The choice of model for deployment should consider the bank's tolerance for false positives versus false negatives. If avoiding false negatives is crucial (i.e., not denying a loan to someone who should be approved), a model with a high recall like the **Random Forest** might be preferred.
    ○ If the bank wants to minimize the risk of default (false positives), then a model with higher precision would be more suitable.
7. Feature Importance:
    ○ An analysis of feature importance could provide further insights into which factors are most predictive of loan approval. This could also help in refining the models for better performance.

## Conclusion:

The Random Forest model stands out with the highest accuracy and an excellent recall rate, suggesting it as a strong candidate for deployment. However, the final decision should also consider the bank's specific cost-benefit analysis and business objectives. Continuous model evaluation and updates are recommended to adapt to new patterns in applicant data over time.