

# Credit Card Risk Classification: Logistic Regression vs. Ensemble Learning

Aleksandar Dimitrov  
October 12, 2023

# Table of Contents:

- Introduction
- Ensemble Learning: A look at the Landscape
- Ensemble Learning: Definition and Types
- Ensemble Methods: Structure and Operating Principles
- Exploring Ensemble Methods
- Logistic Regression: Definition and Limitation
- Model Outcome and Risk Segmentation
- Application Models vs. Behavioral Models for Credit Risk
- Models Performance Evaluating: Key Measures
- Application
- Questions and Discussion

## Introduction:

- In the context of our analysis, we have conducted an empirical comparison based on the following GitHub project: [Credit Card Risk Classification: Logistic Regression vs. Ensemble Learning](#)
- The empirical comparison is structured around the specific GitHub project, aiming to evaluate its performance and functionality
- The empirical comparison serves as a critical component of our evaluation process

# Ensemble Learning: A look at the Landscape:

## Multiple Perspectives

- Ensemble Methods utilize multiple models, not just a single one, to make predictions
- They lead to better generalization performance
- By aggregating predictions from diverse models, they can capture underlying patterns in the data more effectively

### Increased Security and Reliability

- Achieving greater security through assessment from various angles

### Reduced Subjectivity

- Avoiding biases associated with a single model

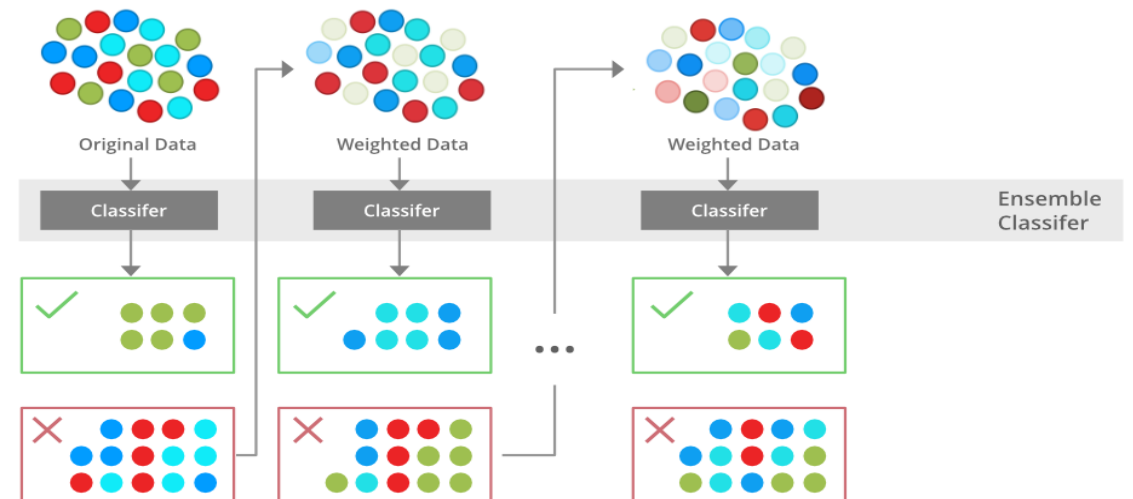
### Robustness to Model Variability

- Even if some of the individual models have variations in their performance, the ensemble as a whole can still provide reliable predictions

## Error Reduction (Improving Accuracy)

- Ensemble Methods reduce errors in predictions by combining multiple models
- Improvement in the accuracy of predictions
- Diversity of models compensates for common errors in individual models
- Ensemble Methods can reduce the impact of noise or outliers in the data, leading to more stable and reliable predictions

## Boosting in Machine Learning

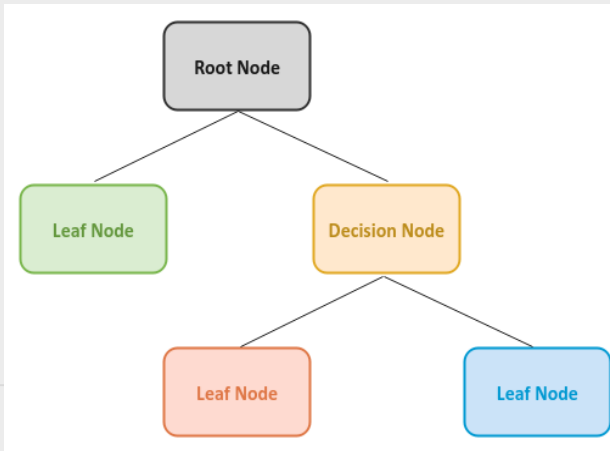


# Ensemble Learning:

## Definition and Types

### Decision Trees

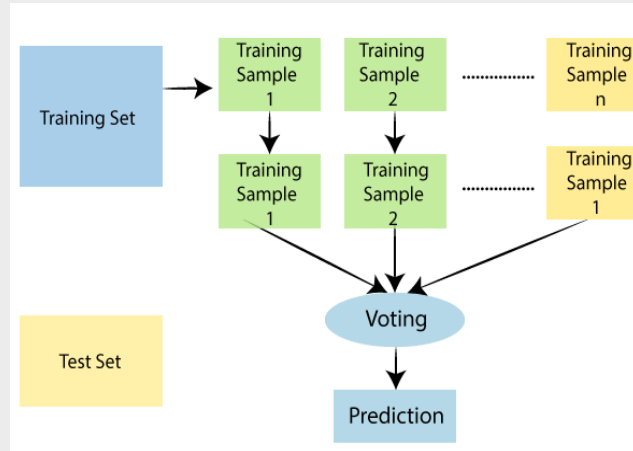
#### Simple Path to Classification



- **Hierarchical Questioning:** Forms a tree structure with input-related questions
- **Classification Focus:** Primarily for classifying data
- **Visual Clarity:** Easy-to-visualize structure
- **Interactive Decision Path:** Users answer questions for classification

### Random Forest

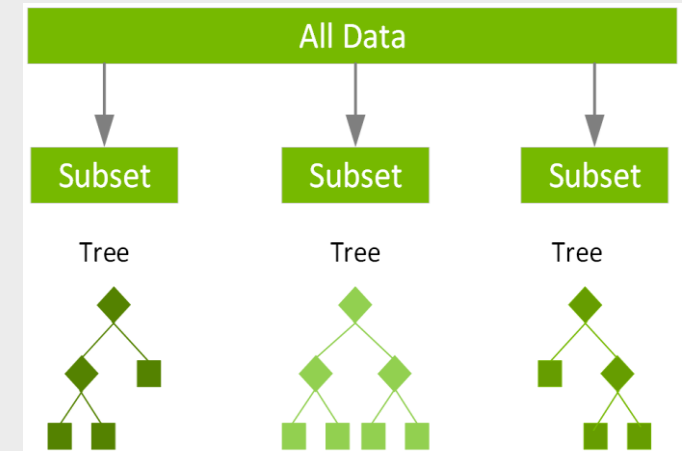
#### Forested Accuracy



- **Decision Tree Ensemble:** It's a combination of multiple decision trees
- **Result Combining:** The final result is based on averaging or majority voting from the trees
- **Bootstrap Aggregating:** It creates diverse datasets and assigns them to trees
- **Accuracy Enhancement:** Its aim is to improve accuracy and mitigate overfitting

### XGBoost

#### Extreme Gradient Precision

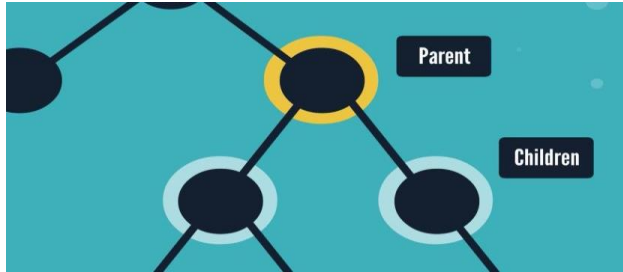


- **Sparse Data Handling:** XGBoost excels with sparse data
- **Weighted Trees:** It employs weighted trees for better predictions
- **Scalable:** Efficient with large datasets
- **Auto Feature Selection:** Automatically selects essential features

# Decision Trees:

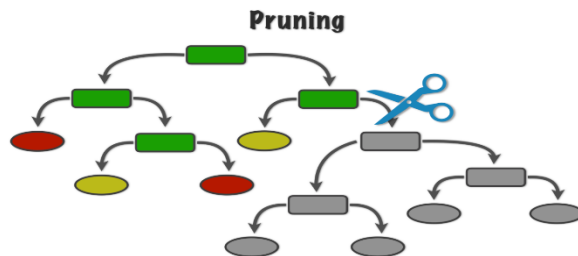
## Structure and Operating Principles

### Decision Trees Model



$$\text{Gini Impurity} = \sum p_i(1 - p_i) = 1 - \sum(p_i^2)$$

$$\text{Entropy} = -\sum p_i \log_2(p_i)$$



### Structure and Operating Principles

#### Construction:

- Build Node by Node
- Nodes selected for the best data split

#### Gini Impurity:

- Measures misclassification in a node
- Lower Gini impurity indicates better split

#### Information Gain (Entropy):

- Measures information in child node
- DT aim to maximize information gain

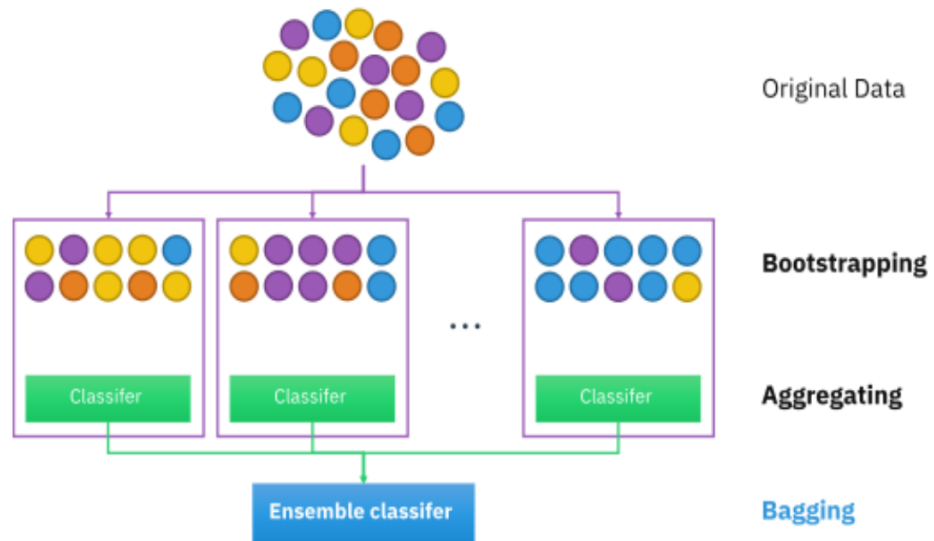
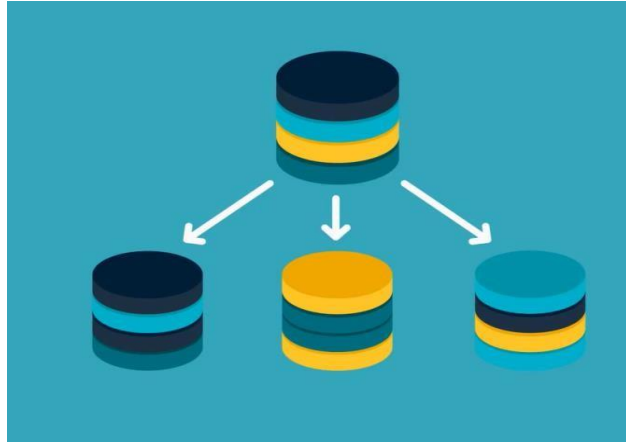
#### Pruning to prevent Overfitting:

- DT tend to overfit
- Pruning reduces complexity by removing branches

# Random Forest:

## Structure and Operating Principles

### Random Forest Model



### Structure and Operating Principles

#### Collective wisdom:

- Ensemble of multiple decision trees
- Combining different learning algorithms

#### Majority Voting:

- With multiple DT, we get multiple answers
- Selecting the most common outcome

#### Bootstrapping:

- Randomly sampling the original dataset with replacement
- Feature subsets are chosen randomly for each tree

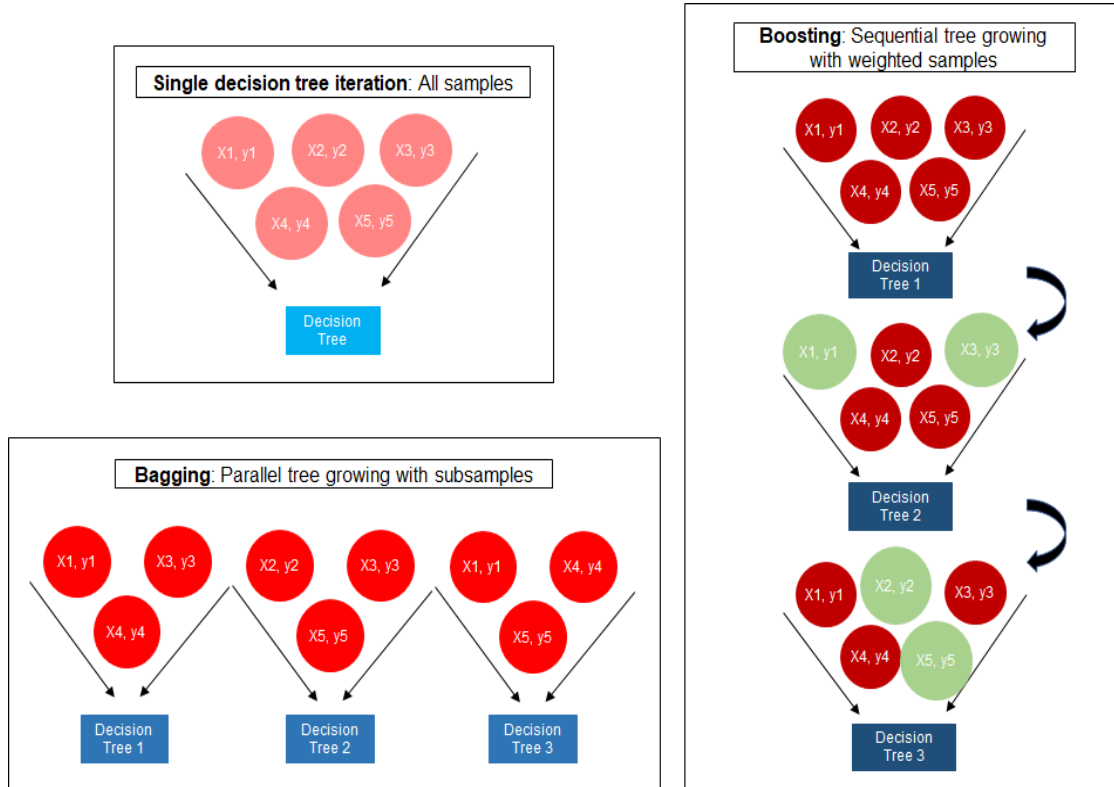
#### Random Forest:

- Bootstrap Aggregated decision trees
- Bootstrapping is type of regularization

# XGBoost:

## Structure and Operating Principles

### XGBoost Model



### Structure and Operating Principles

#### Boosting:

- Predictive accuracy by combining different models
- Model optimizing by trial and error

#### Gradient Boosting:

- Gradient boosting algorithm implementation
- Employing DT with limited depth as weak learners

#### Iterative Model Construction:

- XGBoost is build iteratively
- Each iteration improving upon the errors of the previous ones

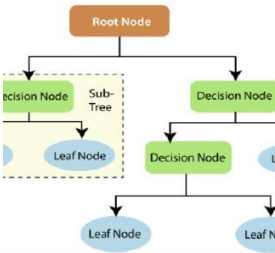
#### Objective Function and Regularization (for Regression only):

- Objective function defining
- L1(Lasso) and L2(Ridge) incorporating



# Exploring Ensemble Methods:

## Decision Trees

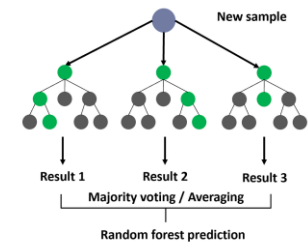


### Applications and Use Cases

Most recent characteristics	<div>Used For</div> <div></div> <div><ul style="list-style-type: none"><li>Regression</li><li>Classification</li></ul></div>	<div>Avoid Overfitting ?</div> <div></div> <div><ul style="list-style-type: none"><li>Perform pruning (during or after the process)</li></ul></div>
	<div>Input</div> <div></div> <div><ul style="list-style-type: none"><li>Numerical</li><li>Categorical</li></ul></div>	<div>Pros</div> <div></div> <div><ul style="list-style-type: none"><li>Simple to understand and interpret</li><li>In build feature selection</li><li>Fast testing on new data</li></ul></div>
	<div>Handles</div> <div></div> <div><ul style="list-style-type: none"><li>Small Dataset</li><li>Large Dataset</li><li>Sparse data</li><li>Hight Dimensions</li></ul></div>	<div>Cons</div> <div></div> <div><ul style="list-style-type: none"><li>Sensitivity to Small Data Variations</li><li>Prone to Overfitting</li><li>Limited Optimization Options</li></ul></div>
	<div>Preprocess</div> <div></div> <div><ul style="list-style-type: none"><li>No preprocessing required</li></ul></div>	<div>Application</div> <div></div> <div><ul style="list-style-type: none"><li>Credit card fraud detection</li><li>Credit risk application and behave models</li><li>Medical diagnoses</li></ul></div>
	<div>Algorithm Speed</div> <div></div> <div><ul style="list-style-type: none"><li>Training: Fast</li><li>Testing: Really Fast</li></ul></div>	

# Exploring Ensemble Methods:

## Random Forest

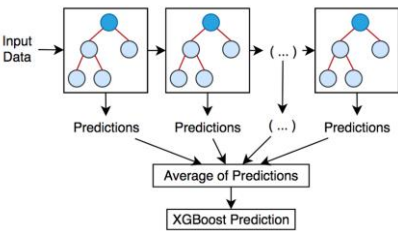


### Applications and Use Cases

Most recent characteristics	<div>Used For</div> <div></div> <div><ul style="list-style-type: none"><li>Regression</li><li>Classification</li></ul></div>	<div>Avoid Overfitting ?</div> <div></div> <div><ul style="list-style-type: none"><li>Prune the individual trees</li></ul></div>
	<div>Input</div> <div></div> <div><ul style="list-style-type: none"><li>Numerical</li><li>Categorical</li></ul></div>	
	<div>Handles</div> <div></div> <div><ul style="list-style-type: none"><li>Small Dataset</li><li>Large Dataset</li><li>Sparse data</li><li>Hight Dimensions</li></ul></div>	
	<div>Preprocess</div> <div></div> <div><ul style="list-style-type: none"><li>No preprocessing required</li></ul></div>	
	<div>Algorithm Speed</div> <div></div> <div><ul style="list-style-type: none"><li>Training: Fast</li><li>Testing: Fast</li></ul></div>	
	<div>Pros</div> <div></div> <div><ul style="list-style-type: none"><li>Performs well with large dataset</li><li>Lots of hyperparameters to control</li><li>Relatively better results than decision trees</li></ul></div>	
	<div>Cons</div> <div></div> <div><ul style="list-style-type: none"><li>Less interpretable than decision trees</li><li>Doesn't solve regression problems well</li><li>Outperformed by gradient-boosted trees</li></ul></div>	
	<div>Application</div> <div></div> <div><ul style="list-style-type: none"><li>Credit card fraud detection</li><li>Credit card application and behave models</li><li>Medical diagnoses</li></ul></div>	

# Exploring Ensemble Methods:

## XGBoost



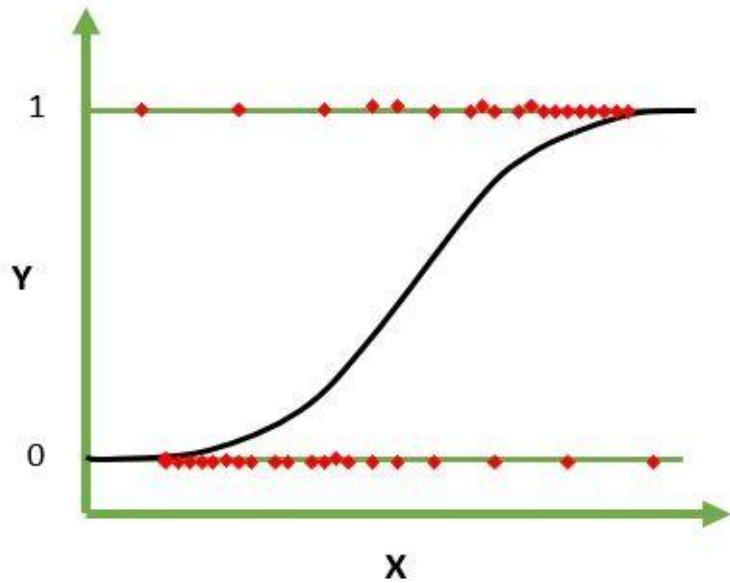
### Applications and Use Cases

Most recent characteristics	<div>Used For</div> <div></div> <div><ul style="list-style-type: none"><li>Regression</li><li>Classification</li></ul></div>		<div>Avoid Overfitting ?</div> <div></div> <div><ul style="list-style-type: none"><li>Resistant to overfitting</li></ul></div>
	<div>Input</div> <div></div> <div><ul style="list-style-type: none"><li>Numerical</li><li>Categorical</li></ul></div>		<div>Pros</div> <div></div> <div><ul style="list-style-type: none"><li>Easy to parallelize</li><li>Smart trees penalize</li><li>Easily scalable</li></ul></div>
	<div>Handles</div> <div></div> <div><ul style="list-style-type: none"><li>Small Dataset</li><li>Large Dataset</li><li>Sparse data</li><li>Hight Dimensions</li></ul></div>		<div>Cons</div> <div></div> <div><ul style="list-style-type: none"><li>Lack of interpretability</li><li>Difficult to optimize the many different parameters</li></ul></div>
	<div>Preprocess</div> <div></div> <div><ul style="list-style-type: none"><li>No preprocessing required</li></ul></div>		<div>Application</div> <div></div> <div><ul style="list-style-type: none"><li>Credit card fraud detection</li><li>Credit card application and behave models</li><li>Store sales predictions</li></ul></div>
	<div>Algorithm Speed</div> <div></div> <div><ul style="list-style-type: none"><li>Training: Fast</li><li>Testing: Fast</li></ul></div>		

# Logistic Regression:

## Definition and Limitation

### Logistic Model



The logistic regression predicts the probability of an event occurring

Logistic regression model

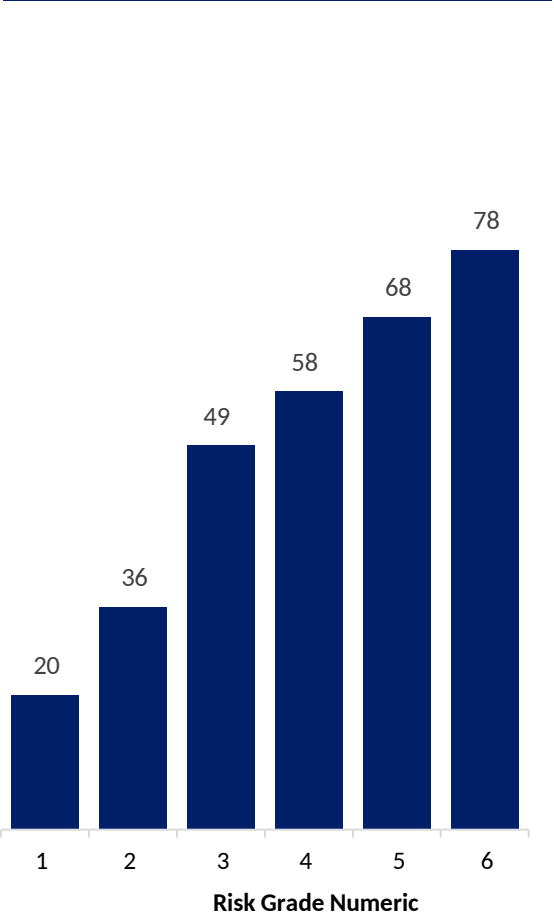
$$\log(\text{odds}) = \beta_0 + \beta_1 x + \cdots \beta_k x_k$$

### Logistic Regression Constraints

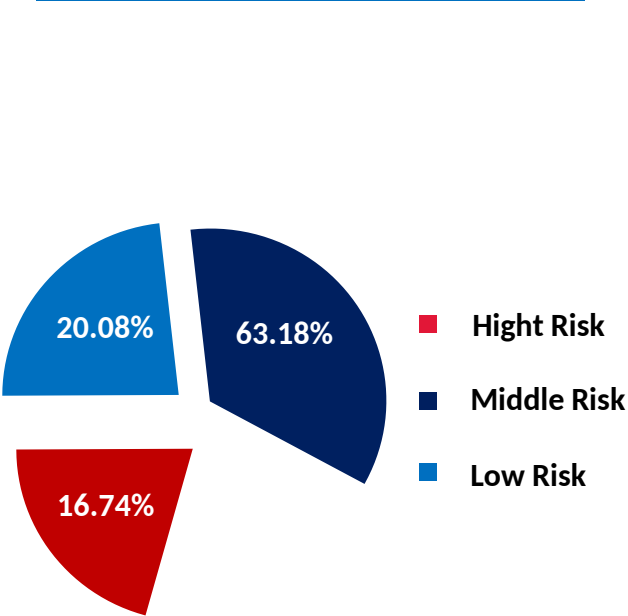
- **Assumes Linearity:** Logistic regression assumes a linear relationship, which can't capture complex non-linear patterns
- **Binary Classification Bias:** May need modifications for multi-class problems
- **Overfitting Risk:** Susceptible to overfitting
- **Independence Assumption:** Assumes observations are independent, which may not hold in time series or spatial data
- **Outlier Sensitivity:** Outliers can strongly influence results
- **Handling Missing Data:** Logistic regression struggles with missing data and imputation can introduce bias

# Model Outcome and Risk Segmentation:

Bad Rate after activation



Risk Class Segmentation

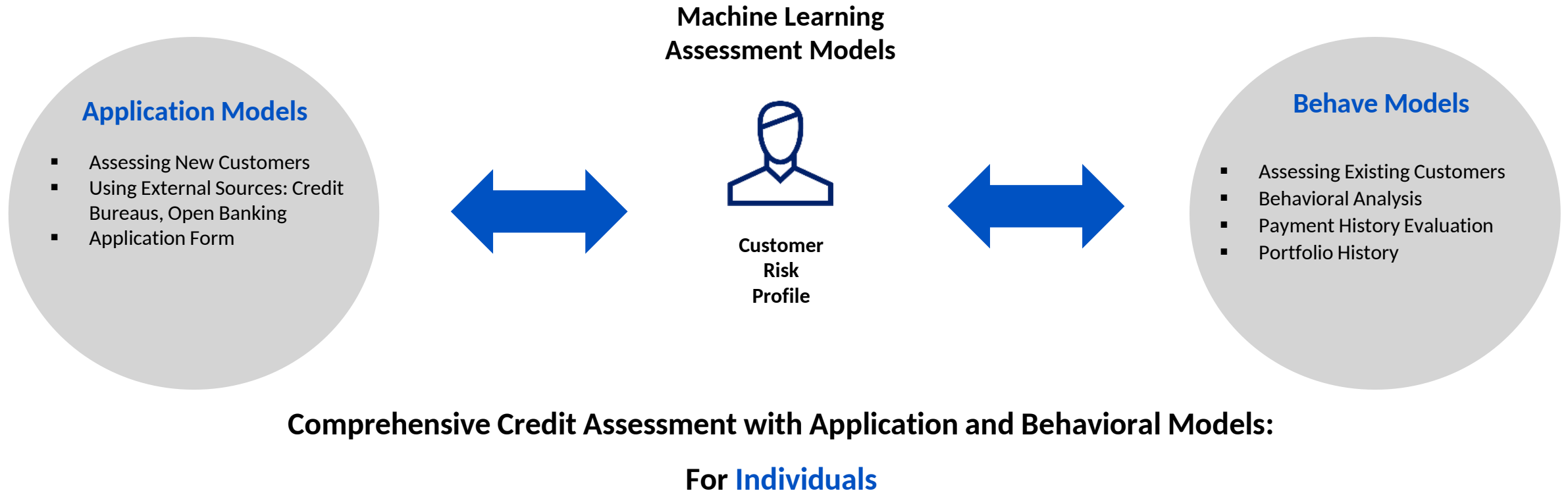


Analyzing Model Output

- **Model Results:** An overview of model outcomes and output probabilities
- **Risk Segmentation:** How the model divides data into segments based on default probability
- **Customer Classification :** Categorizing customers from low-risk (marked as 1) to high-risk (marked as 6-15 classes)
- **Segment Assessment:** Evaluating risk segments(Low, Middle and Hight Risk ), including calculating the percentage of bad customers
- **Model Application:** How to apply the model's output for new data assessment and decision-making

# Application Models vs. Behavioral Models for Credit Risk:

ML models leverage data and automation to improve decision-making and efficiency in various contexts



# Models Performance Evaluation:

## KEY Measures

### What numbers tell us...

#### Confusion Matrix and Classification Report

01

CM is **critical tool** in model evaluation. CR includes metrics like precision, recall, F1-score, and support for each class

#### AUC (Area under the Curve)

02

**Intuitive Interpretation:** A higher AUC means better client classification by the model

#### Gini Coefficient


03

Gini is calculated as  **$Gini = 2 * AUC - 1$**

### It's Working...


For Credit Risk models, higher Gini values are preferred

75%+




**For Application models**  
Gini > 50% implies effective model, correctly identifies ~75% clients

↑ 50%



**For Application models**  
Gini > 50% is desirable

30%



**For Behave models**  
Gini > 60% is desirable, for even higher accuracy

15

## Application:

### Performance Measure

$$\text{Accuracy} = (\text{TN} + \text{TP}) / \text{TN} + \text{FN} + \text{FP} + \text{TP}$$

$$\text{Precision} = \text{TP} / \text{TP} + \text{FP}$$

$$\text{Recall} = \text{TP} / \text{TP} + \text{FN}$$

$$\text{F}_1\text{Score} = 2 / (1/\text{precision} + 1/\text{recall})$$

### Description

- The ratio between the number of all correctly predicted samples and the number of all samples
- The ratio between the number of true positives and the number of all samples classified as positive
- The ratio between the number of true positives and the number of all samples whose true class is the positive one
- The harmonic mean of precision and recall



**THANK YOU !**