

การทำนายการผิดนัดชำระหนี้ของสินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพ :
กรณีศึกษาสถาบันการเงินแห่งหนึ่งในประเทศไทย

PREDICTING DEFAULTS ON NANO FINANCE LOANS: A CASE
STUDY OF A FINANCIAL INSTITUTION IN THAILAND



การค้นคว้าอิสระนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตร
ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาสถิติและการวิเคราะห์ธุรกิจ
ภาควิชาสถิติ คณะวิทยาศาสตร์
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง
พ.ศ. 2568

KMITL-2025-SC-M-050-020

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

PREDICTING DEFAULTS ON NANO FINANCE LOANS: A CASE
STUDY OF A FINANCIAL INSTITUTION IN THAILAND



PREMSUDA PATSAYO

AN INDEPENDENT STUDY SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENT FOR THE DEGREE OF MASTER OF SCIENCE IN STATISTICS AND
BUSINESS ANALYTICS

DEPARTMENT OF STATISTICS SCHOOL OF SCIENCE
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

2025

KMITL-2025-SC-M-050-020

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



COPYRIGHT 2025

SCHOOL OF SCIENCE

KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หัวข้อวิทยานิพนธ์	การทำนายการผิมนัดชำระหนี้ของสินเชื่อรายย่อยเพื่อผู้ประกอบ อาชีพ : กรณีศึกษาสถาบันการเงินแห่งหนึ่งในประเทศไทย
ชื่อนักศึกษา	เปรมสุดา ปัดสาโย
รหัสประจำตัว	65056106
ปริญญา	วิทยาศาสตรมหาบัณฑิต (สถิติและการวิเคราะห์ธุรกิจ)
ภาควิชา	สถิติ
พ.ศ.	2568
อาจารย์ที่ปรึกษาวิทยานิพนธ์	ผู้ช่วยศาสตราจารย์ ดร.กนกกรรณ์ ลีโรจนาประภา

บทคัดย่อ

งานวิจัยนี้มีวัตถุประสงค์เพื่อสร้างแบบจำลองที่เหมาะสมสำหรับการทำนายการผิมนัดชำระหนี้ของสินเชื่อรายย่อยเพื่อผู้ประกอบอาชีพ (Nano Finance) และวิเคราะห์ปัจจัยที่มีอิทธิพลต่อการผิมนัดชำระหนี้ ในการศึกษาครั้งนี้ ผู้วิจัยได้แบ่งกลุ่มลูกหนี้ออกเป็น 5 กลุ่มตามพฤติกรรมการผิมนัดชำระหนี้ ได้แก่ ลูกหนี้ทั่วไป กลุ่มความเสี่ยงต่ำ กลุ่มความเสี่ยงปานกลาง กลุ่มความเสี่ยงสูง และกลุ่มความเสี่ยงวิกฤต โดยมีวัตถุประสงค์เพื่อศึกษาความแม่นยำของแบบจำลองในแต่ละกลุ่มความเสี่ยงและออกแบบแนวทางการบริหารจัดการที่เหมาะสม ผู้วิจัยนำเทคนิคการเรียนรู้ของเครื่องและการเรียนรู้เชิงลึก 5 วิธี ได้แก่ การถดถอยโลจิสติก (Logistic Regression) ต้นไม้ตัดสินใจ (Decision Tree) โครงข่ายประสาทเทียม (Neural Network) LSTM (Long Short-Term Memory) และ Bidirectional LSTM มาทำการเปรียบเทียบโดยใช้ค่าประสิทธิภาพโดยรวม (F1-score) และ ค่าความระลึก (Recall) ผลการวิจัยพบว่า ในกลุ่มลูกหนี้ความเสี่ยงในภาพรวม แบบจำลอง Bidirectional LSTM ให้ค่าประสิทธิภาพโดยรวม (F1-score) สูงสุดที่ร้อยละ 80.56 และค่าความระลึก (Recall) เท่ากับร้อยละ 80.68 แบบจำลองดังกล่าวยังแสดงศักยภาพในการทำนายได้แม่นยำในกลุ่มลูกหนี้ที่มีประวัติผิมนัดต่อเนื่อง โดยเฉพาะกลุ่มความเสี่ยงต่ำ กลาง และสูง ส่วนกลุ่มความเสี่ยงวิกฤต พบว่า LSTM มีประสิทธิภาพสูงสุด นอกจากนี้ยังพบว่าปัจจัยสำคัญที่มีอิทธิพลต่อการผิมนัดชำระหนี้ในทุกแบบจำลอง ได้แก่ ปัจจัยที่เกี่ยวข้องกับประวัติการผิมนัดชำระหนี้ย้อนหลังของลูกหนี้

คำสำคัญ : การถดถอยโลจิสติก, การทำนายการผิมนัดชำระหนี้, การแบ่งกลุ่มลูกหนี้, ต้นไม้ตัดสินใจ, โครงข่ายประสาทเทียม, ปัจจัยเสี่ยง, Bidirectional LSTM, LSTM, สินเชื่อรายย่อย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Independent Study Title	Predicting Defaults on Nano Finance Loans: A Case Study of a Financial Institution in Thailand
Student Name	Premsuda Patsayo
Student ID	65056106
Degree	Master of Science (Statistics And Business Analytics)
Department	Statistics
Year	2025
Independent Study Advisor	Assistant Professor Dr. Kanogkan Leerojanaprapa

Abstract

This study aims to develop suitable models for predicting loan defaults among nano finance borrowers and analyze the factors influencing these defaults. The study categorizes borrowers into five distinct groups based on their delinquency behaviors: general borrowers, low-risk, moderate-risk, high-risk, and critical-risk groups. The primary objective is to evaluate model accuracy within each risk segment and propose effective management strategies tailored to each category. Five machine learning and deep learning methodologies—Logistic Regression, Decision Tree, Neural Network, Long Short-Term Memory (LSTM), and Bidirectional LSTM, evaluated using F1-score and Recall metrics. Results indicate that the Bidirectional LSTM achieved the highest overall performance in the general risk group, with an F1-score of 80.56% and Recall of 80.68%. Additionally, the Bidirectional LSTM showed superior predictive accuracy for customers with histories of continuous default, particularly within the Low, Moderate, and High-risk groups. However, the LSTM model performed best for the Critical Risk group. Across all models, historical default-related factors consistently emerged as the most influential predictors of loan default.

Keywords: Logistic Regression, Default Prediction, Customer Segmentation, Decision Tree, Neural Network, Risk Factors, Bidirectional LSTM, LSTM, Nano Finance

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

กิตติกรรมประกาศ

ข้าพเจ้าขอขอบคุณอาจารย์ที่ปรึกษา ผศ.ดร.กนกวรรณ ลีโรจนาประภา ที่ได้กรุณาให้คำแนะนำ ชี้แนะแนวทาง ตลอดจนให้ความช่วยเหลือในทุกขั้นตอนของการดำเนินงานวิจัยฉบับนี้ด้วยความเอาใจใส่และอดทนเสมอมา

ขอขอบคุณกรรมการสอบ ผศ.ดร.สิทธิชัย เจริญเศรษฐศิลป์ และผศ.ดร.พรพิมล ชัยวุฒิศักดิ์ ที่ได้ให้ข้อเสนอแนะและคำแนะนำอันมีคุณค่า ทำให้งานวิจัยฉบับนี้มีความสมบูรณ์ยิ่งขึ้น

ขอขอบคุณสถาบันการเงินที่ให้ข้อมูลและสนับสนุนทรัพยากรอันสำคัญในการดำเนินการวิจัยนี้ รวมถึงขอขอบคุณเพื่อนร่วมงาน หัวหน้างาน และผู้ที่เกี่ยวข้องทุกท่านที่ทำให้กำลังใจ ให้คำปรึกษาและความร่วมมือด้วยดีตลอดมา

ท้ายที่สุดนี้ขอขอบคุณ ภรรยา ชินคำวงศ์ เพื่อนที่ชักชวนมาให้มาศึกษา ขอขอบคุณคุณนพพร แยมสุวรรณ รวมถึงครอบครัวที่เป็นกำลังใจและส่งเสริมสนับสนุนสำคัญในทุกช่วงเวลาของการศึกษา และการดำเนินงานวิจัยฉบับนี้ให้สำเร็จลุล่วงไปด้วยดี ขอขอบพระคุณทุกท่านมา ณ โอกาสนี้

นางสาวเปรมสุตา ปัดสาโย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ

	หน้า
บทคัดย่อภาษาไทย	ก
บทคัดย่อภาษาอังกฤษ	ข
กิตติกรรมประกาศ	ค
สารบัญ.....	ง
สารบัญตาราง.....	ช
สารบัญรูป.....	ญ
บทที่ 1 บทนำ.....	1
1.1 ความสำคัญและที่มาของปัญหา	1
1.2 วัตถุประสงค์ของงานวิจัย	3
1.3 ขอบเขตการวิจัย.....	4
1.4 ประโยชน์ที่คาดว่าจะได้รับ	4
1.5 กรอบแนวคิดการวิจัย.....	4
1.6 นิยามศัพท์เฉพาะ	6
บทที่ 2 ทฤษฎีและงานวิจัยที่เกี่ยวข้อง	8
2.1 แนวคิดและทฤษฎีเกี่ยวกับปัจจัยส่วนบุคคล.....	8
2.2 แนวคิดและทฤษฎีเกี่ยวกับปัจจัยด้านลักษณะสินเชื่อ.....	9
2.2.1 ประเภทของสินเชื่อ	9
2.2.2 ปัจจัยที่ใช้ในการพิจารณาสินเชื่อ.....	10
2.3 แนวคิดและทฤษฎีเกี่ยวข้องกับกระบวนการวิเคราะห์ข้อมูลแบบ CRISP-DM.....	11
2.3.1 การทำความเข้าใจธุรกิจ (Business Understanding).....	12
2.3.2 การทำความเข้าใจข้อมูล (Data Understanding)	13
2.3.3 การเตรียมข้อมูล (Data Preparation)	14
2.3.4 การสร้างแบบจำลอง (Modeling).....	14
2.3.5 การประเมินผล (Evaluation).....	15
2.3.6 การนำไปใช้ (Deployment).....	16
2.4 แนวคิดและทฤษฎีเกี่ยวกับการแบ่งข้อมูลเพื่อทำการทดสอบและการสุ่มตัวอย่างข้อมูล	16
2.4.1 การแบ่งข้อมูลเพื่อประเมินแบบจำลอง.....	17
2.4.2 เทคนิคการสุ่มตัวอย่างข้อมูล	17

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ (ต่อ)

	หน้า
2.5 แนวคิดและทฤษฎีเกี่ยวข้องกับการจำแนกประเภทของข้อมูล (Classification).....	18
2.5.1 แบบจำลอง Logistic Regression	18
2.5.2 แบบจำลอง Decision Tree	20
2.5.3 แบบจำลอง Neural Networks	22
2.5.4 แบบจำลอง Long Short-Term Memory (LSTM).....	24
2.5.5 แบบจำลอง Bidirectional Long Short-Term Memory (BiLSTM).....	26
2.6 แนวคิดและทฤษฎีเกี่ยวข้องกับการเลือกไฮเปอร์พารามิเตอร์แบบกริด (Grid Search Cross-Validation).....	27
2.7 แนวคิดและทฤษฎีเกี่ยวข้องกับการประเมินประสิทธิภาพในการทำนายของแบบจำลอง.....	29
2.7.1 เมทริกซ์ความสับสน (Confusion matrix).....	29
2.7.2 การจำแนกประเภท (Classification Report).....	30
2.7.3 การทดสอบสมมติฐานด้วย McNemar's Test.....	32
2.8 แนวคิดและทฤษฎีเกี่ยวข้องกับความสำคัญของคุณลักษณะ (Feature Importance).....	33
2.8.1 ความสำคัญของคุณลักษณะแบบ Permutation Importance.....	34
2.8.2 SHAP (SHapley Additive exPlanations).....	35
2.9 งานวิจัยที่เกี่ยวข้อง.....	37
บทที่ 3 วิธีดำเนินการวิจัย.....	40
3.1 การทำความเข้าใจธุรกิจ (Business Understanding).....	40
3.2 การทำความเข้าใจข้อมูล (Data Understanding)	40
3.2.1 การเก็บข้อมูล (Data Collection).....	41
3.2.2 การอธิบายข้อมูล.....	41
3.3 การเตรียมข้อมูล (Data Preparation)	44
3.3.1 การเลือกข้อมูล (Data Selection)	44
3.3.2 การทำความสะอาดข้อมูล (Data Cleansing).....	44
3.3.3 การเตรียมข้อมูล (Data Preparation)	46
3.3.4 การแปลงข้อมูล (Feature Encoding)	47
3.4 การสร้างแบบจำลอง (Modeling).....	48
3.4.1 การเลือกวิธีแบ่งข้อมูลเป็นข้อมูลฝึกฝนและข้อมูลทดสอบ (Training/Testing).....	48
3.4.2 การเลือกไฮเปอร์พารามิเตอร์ด้วยการค้นหาแบบกริด (Grid Search)	49

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ (ต่อ)

	หน้า
3.4.3 ขั้นตอนการสร้างแบบจำลอง.....	53
3.5 การประเมินผล (Evaluation).....	55
3.5.1 เมทริกซ์ความสับสน (Confusion Matrix).....	55
3.5.2 การจำแนกประเภท (Classification Report).....	55
3.6 การนำแบบจำลองไปใช้งาน (Deployment).....	56
บทที่ 4 ผลการวิจัยและอภิปรายผล	57
4.1 สถิติเชิงพรรณนา (Descriptive Statistics).....	59
4.2 ผลการคัดเลือกการแบ่งข้อมูลชุดทดสอบด้วยวิธีแบบแยกชุด (Hold-out) และวิธีแบบไขว้ (K-Fold Cross Validation).....	67
4.2.1 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Logistic Regression.....	70
4.2.2 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Decision Tree	71
4.2.3 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Neural Network.....	72
4.2.4 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง LSTM.....	73
4.2.5 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Bidirectional LSTM.....	74
4.3 ผลการทำนายการผิदनัดชำระหนี้ในกลุ่มที่ 1 ลูกค้าทั่วไป (General Customer).....	75
4.3.1 ผลการทำนายการผิदनัดชำระหนี้ของแบบจำลอง Logistic Regression	76
4.3.2 ผลการทำนายการผิदनัดชำระหนี้ของแบบจำลอง Decision Tree.....	77
4.3.3 ผลการทำนายการผิदनัดชำระหนี้ของแบบจำลอง Neural Network.....	79
4.3.4 ผลการทำนายการผิदनัดชำระหนี้ของแบบจำลอง LSTM.....	81
4.3.5 ผลการทำนายการผิदनัดชำระหนี้ของแบบจำลอง Bidirectional LSTM	83
4.3.6 ผลการเปรียบเทียบประสิทธิภาพของ 5 เทคนิคในการทำนายการผิदनัดชำระหนี้.....	84
4.3.7 การทดสอบสมมติฐานด้วย McNemar's Test.....	85
4.3.8 ผลการวิเคราะห์ความสำคัญของตัวแปร (Feature Importance)	86
4.4 ผลการทำนายการผิदनัดชำระหนี้ในกลุ่มที่ 2 ความเสี่ยงต่ำ (Low Risk)	95
4.5 ผลการทำนายการผิदनัดชำระหนี้ในกลุ่มที่ 3 ความเสี่ยงปานกลาง (Moderate Risk).....	97
4.6 ผลการทำนายการผิदनัดชำระหนี้ในกลุ่มที่ 4 ความเสี่ยงสูง (High Risk)	100
4.7 ผลการทำนายการผิदनัดชำระหนี้ในกลุ่มที่ 5 ความเสี่ยงวิกฤต (Critical Risk).....	102
4.8 อภิปรายผล	105

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญ (ต่อ)

	หน้า
บทที่ 5 สรุปผลการวิจัยและข้อเสนอแนะ.....	106
5.1 สรุปผลการวิจัย	106
5.2 ข้อเสนอแนะ.....	108
เอกสารอ้างอิง	110
ประวัติผู้เขียน.....	115



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญตาราง

ตารางที่	หน้า
ตารางที่ 1.1	6
ตารางที่ 2.1	32
ตารางที่ 3.1	41
ตารางที่ 3.2	42
ตารางที่ 3.3	45
ตารางที่ 3.4	47
ตารางที่ 3.5	47
ตารางที่ 3.6	50
ตารางที่ 3.7	50
ตารางที่ 3.8	51
ตารางที่ 3.9	51
ตารางที่ 3.10	52
ตารางที่ 3.11	53
ตารางที่ 4.1	66
ตารางที่ 4.2	68
ตารางที่ 4.3	69
ตารางที่ 4.4	71
ตารางที่ 4.5	72
ตารางที่ 4.6	73
ตารางที่ 4.7	74
ตารางที่ 4.8	75

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
ตารางที่ 4.9 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Logistic Regression	77
ตารางที่ 4.10 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Decision Tree.....	79
ตารางที่ 4.11 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Neural Network.....	80
ตารางที่ 4.12 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง LSTM.....	82
ตารางที่ 4.13 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Bidirectional LSTM.....	84
ตารางที่ 4.14 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 1 ลูกค้าทั่วไป (General Customer).....	85
ตารางที่ 4.15 ผลการเปรียบเทียบระหว่าง Bidirectional LSTM กับแบบจำลองอื่นๆ.....	86
ตารางที่ 4.16 ผลวิเคราะห์ Permutation Importance และ SHAP Value ของแต่ละแบบจำลอง	95
ตารางที่ 4.17 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองในกลุ่มที่ 2 ความเสี่ยงต่ำ (Low Risk).....	96
ตารางที่ 4.18 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 2 ความเสี่ยงต่ำ (Low Risk)	97
ตารางที่ 4.19 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองในกลุ่มที่ 3 ความเสี่ยงปานกลาง (Moderate Risk)	98
ตารางที่ 4.20 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 3 ความเสี่ยงปานกลาง (Moderate Risk)	99
ตารางที่ 4.21 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองในกลุ่มที่ 4 ความเสี่ยงสูง (High Risk).....	101
ตารางที่ 4.22 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 4 ความเสี่ยงสูง (High Risk)	102
ตารางที่ 4.23 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองในกลุ่มที่ 5 ความเสี่ยงวิกฤต (Critical Risk)	103
ตารางที่ 4.24 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 5 ความเสี่ยงวิกฤต (Critical Risk).....	104
ตารางที่ 5.1 ผลเปรียบเทียบค่าประสิทธิภาพโดยรวม (F1-Score) ตามกลุ่มลูกค้าของทั้ง 5 เทคนิค	106
ตารางที่ 5.2 ผลเปรียบเทียบค่าความระลึก (Recall) ตามกลุ่มลูกค้าของทั้ง 5 เทคนิค.....	106

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สารบัญรูป

รูปที่	หน้า
รูปที่ 1.1 สัดส่วนหนี้ครัวเรือนต่อ GDP ในช่วงไตรมาส 1 ปี พ.ศ. 2566	1
รูปที่ 1.2 กรอบแนวคิดในการวิจัย	5
รูปที่ 2.1 กระบวนการวิเคราะห์ข้อมูลแบบ CRISP-DM.....	12
รูปที่ 2.2 ตัวอย่างการแยกกลุ่มข้อมูลสองกลุ่มด้วย Logistic Regression	19
รูปที่ 2.3 แสดงการแบ่งข้อมูลแต่ละโหนดและการจำแนกข้อมูลที่ใบของต้นไม้	21
รูปที่ 2.4 โครงสร้างโครงข่ายประสาทเทียมแบบง่าย	23
รูปที่ 2.5 โครงสร้างเซลล์ LSTM.....	24
รูปที่ 2.6 โครงสร้างของแบบจำลอง Bidirectional LSTM	26
รูปที่ 2.7 ขั้นตอนการทำงานของ Grid Search Cross-Validation	28
รูปที่ 2.8 เมทริกซ์ความสับสน (Confusion Matrix) แบบค่าที่ทำนายกับค่าจริง (Predict-Actual) 30	
รูปที่ 2.9 เมทริกซ์ความสับสน (Confusion Matrix) แบบค่าจริงกับค่าที่ทำนาย (Actual-Predict) 30	
รูปที่ 2.10 ขั้นตอนการคำนวณ Permutation Importance	35
รูปที่ 2.11 ตัวอย่าง Beeswarm Plot ของค่า SHAP	36
รูปที่ 3.1 การแปลงข้อมูลเชิงกลุ่มแบบไม่มีลำดับด้วยวิธี One Hot Encoding	48
รูปที่ 3.2 วิธีแบ่งข้อมูลแบบแยกชุด (Hold-Out).....	48
รูปที่ 3.3 วิธีแบ่งข้อมูลแบบไขว้ (K-Fold Cross Validation).....	49
รูปที่ 3.4 การประเมินผลด้วยเมทริกซ์ความสับสน (Confusion Matrix) ของแบบจำลอง.....	55
รูปที่ 3.5 การประเมินผลข้อมูลโดยการวิเคราะห์ความแม่นยำของแบบจำลอง	56
รูปที่ 4.1 วิธีการแบ่งกลุ่มลูกค้า	58
รูปที่ 4.2 สถานะผิมนัดชำระเงินของลูกค้านี้ ณ เดือน มกราคม พ.ศ. 2567	59
รูปที่ 4.3 เพศของลูกค้านี้รายบัญชี	60
รูปที่ 4.4 การเข้าโปรแกรมมาตรการช่วยเหลือลูกค้านี้รายบัญชี	61
รูปที่ 4.5 สถานภาพการสมรสของลูกค้านี้รายบัญชี	62
รูปที่ 4.6 ภูมิภาคที่อยู่อาศัยของบัญชีลูกค้านี้รายบัญชี.....	63
รูปที่ 4.7 ระดับความเสี่ยงของลูกค้านี้รายบัญชี.....	64
รูปที่ 4.8 ชั้นหนี้ของลูกค้านี้รายบัญชี	65
รูปที่ 4.9 เมทริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Logistic Regression	76

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

รูปที่ 4.10 เมตริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Decision Tree..... 78

รูปที่ 4.11 เมตริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Neural Network.....80

รูปที่ 4.12 เมตริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง LSTM.....81

รูปที่ 4.13 เมตริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Bidirectional LSTM.....83

รูปที่ 4.14 ผลจากการทำ Permutation Importance ของแบบจำลอง Logistic Regression87

รูปที่ 4.15 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Logistic Regression.....88

รูปที่ 4.16 ผลจากการทำ Permutation Importance ของแบบจำลอง Decision Tree.....89

รูปที่ 4.17 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Decision Tree89

รูปที่ 4.18 ผลจากการทำ Permutation Importance ของแบบจำลอง Neural Network.....90

รูปที่ 4.19 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Neural Network.....91

รูปที่ 4.20 ผลจากการทำ Permutation Importance ของแบบจำลอง LSTM92

รูปที่ 4.21 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง LSTM.....92

รูปที่ 4.22 ผลจากการทำ Permutation Importance ของแบบจำลอง Bidirectional LSTM93

รูปที่ 4.23 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Bidirectional LSTM.....94

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 1

บทนำ

1.1 ความสำคัญและที่มาของปัญหา

สินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพเป็นแหล่งเงินทุนสำคัญสำหรับผู้ประกอบการอาชีพอิสระ และประชาชนรายย่อยที่ไม่สามารถเข้าถึงสินเชื่อจากสถาบันการเงินหลักได้ สินเชื่อประเภทนี้มีบทบาทสำคัญในการสนับสนุนการเติบโตของเศรษฐกิจของประเทศ โดยเฉพาะประชาชนรายย่อยที่ต้องการเงินทุนสำหรับผู้ประกอบการอาชีพ แต่ไม่สามารถเข้าถึงแหล่งสินเชื่อในระบบของสถาบันการเงินได้ อย่างไรก็ตาม สินเชื่อรายย่อยมักมีความเสี่ยงสูงกว่าสินเชื่อประเภทอื่น เนื่องจากผู้กู้ส่วนใหญ่มักขาดเอกสารที่แสดงแหล่งที่มาของรายได้ มีประวัติทางการเงินที่ไม่ชัดเจน หรือไม่มีทรัพย์สินเป็นหลักประกันในการกู้ยืมเงิน (อุไรพรรณ เจริญรัฐ และ ภาสกร ตาปสนันท์, 2558) ด้วยเหตุนี้ การผิมนัดชำระหนี้จึงกลายเป็นปัญหาสำคัญที่สถาบันการเงินต้องบริหารจัดการอย่างมีประสิทธิภาพเพื่อลดความเสี่ยงที่อาจเกิดขึ้น

ปัญหาหนี้ครัวเรือนในประเทศไทยยังคงเป็นประเด็นสำคัญที่ส่งผลกระทบต่อทั้งระบบเศรษฐกิจและความมั่นคงทางการเงินของประเทศ จากข้อมูล ณ เดือนกรกฎาคม พ.ศ. 2566 พบว่าคนไทยมากกว่า 1 ใน 3 แบกรับภาระหนี้ครัวเรือน คิดเป็นร้อยละ 90.6 ของ GDP (ธนาคารแห่งประเทศไทย, 2566)



รูปที่ 1.1 สัดส่วนหนี้ครัวเรือนต่อ GDP ในช่วงไตรมาส 1 ปี พ.ศ. 2566

ที่มา : ธนาคารแห่งประเทศไทย (2566)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยเฉพาะกลุ่มผู้มีรายได้น้อยซึ่งมีข้อจำกัดในการเข้าถึงสินเชื่อในระบบและมักต้องพึ่งพาสินเชื่อรายย่อยหรือแหล่งเงินกู้นอกระบบ ภาระหนี้ที่สูงนี้นำไปสู่ปัญหาการชำระหนี้ล่าช้าหรือผิดนัดชำระ ส่งผลให้ลูกหนี้จำนวนมากเผชิญความเสี่ยงทางการเงินที่ทวีความรุนแรงขึ้นอย่างต่อเนื่อง

การเพิ่มขึ้นของอัตราการผิดนัดชำระหนี้ในกลุ่มสินเชื่อรายย่อยไม่เพียงส่งผลต่อฐานะทางการเงินของลูกหนี้เท่านั้น แต่ยังก่อให้เกิดผลกระทบต่อสถาบันการเงินโดยตรง กล่าวคือ เมื่อหนี้เสีย (Non-Performing Loan: NPL) เพิ่มขึ้น สถาบันการเงินต้องเผชิญกับความเสียหายด้านสภาพคล่องและต้นทุนในการบริหารจัดการหนี้ที่สูงขึ้น อาจส่งผลต่อความเชื่อมั่นของนักลงทุนและเสถียรภาพของระบบการเงินโดยรวม หากไม่สามารถประเมินและบริหารความเสี่ยงเหล่านี้ได้อย่างมีประสิทธิภาพย่อมสร้างความเปราะบางให้กับเศรษฐกิจในระยะยาว ด้วยความท้าทายดังกล่าว การประเมินความเสี่ยงของลูกหนี้และการคาดการณ์โอกาสผิดนัดชำระหนี้จึงเป็นภารกิจสำคัญของสถาบันการเงิน อย่างไรก็ตาม วิธีการประเมินแบบดั้งเดิม เช่น การพิจารณาเพียงประวัติการเงินหรือการใช้แบบจำลองเชิงสถิติบางประเภท อาจไม่เพียงพอในการรองรับข้อมูลที่มีปริมาณมากและมีความซับซ้อน

ในยุคดิจิทัลที่ข้อมูลขนาดใหญ่ (Big Data) มีบทบาทสำคัญในการพัฒนาระบบการเงิน การนำปัญญาประดิษฐ์ (AI) และการเรียนรู้ของเครื่อง (Machine Learning: ML) เข้ามาช่วยวิเคราะห์และทำนายความเสี่ยงสินเชื่อจึงเป็นแนวทางที่มีประสิทธิภาพและได้รับความสนใจอย่างกว้างขวาง (Noriega et al., 2023) โดยเฉพาะในกรณีของสินเชื่อรายย่อยที่ผู้กู้มักมีพฤติกรรมและความเสี่ยงที่ซับซ้อน แบบจำลองดั้งเดิม เช่น การถดถอยโลจิสติก (Logistic Regression) และต้นไม้ตัดสินใจ (Decision Trees) อาจไม่เพียงพอในการรองรับข้อมูลที่ซับซ้อนมากขึ้น ดังนั้นมีการนำเทคนิคขั้นสูง เช่น เครือข่ายประสาทเทียม (Neural Networks) และโครงข่ายหน่วยความจำระยะยาว (LSTM) มาใช้เพิ่มประสิทธิภาพในการวิเคราะห์ข้อมูล โดย Neural Networks สามารถจำลองความสัมพันธ์ที่ซับซ้อนระหว่างตัวแปร ส่วน LSTM และ Bidirectional LSTM มีประสิทธิภาพในการวิเคราะห์ข้อมูลที่มีลำดับเวลา ซึ่งเป็นประโยชน์ในการระบุพฤติกรรมผู้กู้และคาดการณ์การผิดนัดชำระหนี้ (Goodfellow et al., 2016; Ala'raj et al., 2021)

งานวิจัยนี้มีเป้าหมายเพื่อเปรียบเทียบประสิทธิภาพของแบบจำลองการเรียนรู้ของเครื่อง (Machine Learning: ML) และแบบจำลองการเรียนรู้เชิงลึก (Deep Learning: DL) รวมทั้งสิ้น 5 เทคนิค ได้แก่ Logistic Regression, Decision Trees, Neural Networks, Long Short-Term Memory (LSTM) และ Bidirectional LSTM ในการทำนายความเสี่ยงของการผิดนัดชำระของลูกหนี้ในกลุ่มสินเชื่อรายย่อยสำหรับผู้ประกอบอาชีพ โดยมุ่งหวังให้การวิเคราะห์มีความครอบคลุมทั้งในระบภาพรวมและเชิงลึก เพื่อให้การวิเคราะห์มีความครอบคลุมและสามารถนำไปสู่การบริหารจัดการความเสี่ยงที่เฉพาะเจาะจงยิ่งขึ้น งานวิจัยนี้ได้แบ่งกลุ่มลูกหนี้ออกเป็น 5 กลุ่มตามพฤติกรรมผิดนัดชำระหนี้ ได้แก่ **กลุ่มที่ 1 คือ กลุ่มลูกค้าทั่วไป (General Customer)** ซึ่งครอบคลุมลูกหนี้ทั้งหมดโดยไม่จำกัดพฤติกรรมหรือประวัติการผิดนัดชำระ การวิเคราะห์กลุ่มนี้มีเป้าหมายเพื่อเป็นฐานอ้างอิง

(Baseline) สำหรับการเปรียบเทียบกับแบบจำลองเฉพาะกลุ่ม และเพื่อประเมินภาพรวมของพอร์ตเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สินเชื่อ **กลุ่มที่ 2** คือ **กลุ่มความเสี่ยงต่ำ (Low Risk)** ซึ่งประกอบด้วยลูกหนี้ที่ผิดนัดชำระหนี้เพียง 1 ครั้งในช่วง 2 เดือน สะท้อนถึงความล่าช้าหรือความไม่สม่ำเสมอที่อาจเกิดจากปัจจัยชั่วคราว เช่น การลืมนำชำระหนี้หรือรายได้ขาดช่วง เหมาะสำหรับการจัดการเชิงเตือนหรือส่งเสริมวินัยทางการเงิน **กลุ่มที่ 3** คือ **กลุ่มความเสี่ยงปานกลาง (Moderate Risk)** ประกอบด้วยลูกหนี้ที่ผิดนัดในเดือนล่าสุด ซึ่งอาจเป็นสัญญาณเริ่มต้นของเปลี่ยนแปลงด้านรายได้หรือปัญหาทางการเงินเฉียบพลัน เป็นกลุ่มที่แสดงสัญญาณความเปราะบางทางการเงินที่เริ่มต้นและมีแนวโน้มพัฒนาไปในทางลบ หากไม่ได้รับการติดตามอย่างใกล้ชิด จึงเหมาะสำหรับการใช้แบบจำลองที่สามารถตรวจจับ Early Warning Signal ได้อย่างแม่นยำ **กลุ่มที่ 4** คือ **กลุ่มความเสี่ยงสูง (High Risk)** ซึ่งประกอบด้วยลูกหนี้ที่ผิดนัดติดต่อกันเป็นเวลา 2 เดือน สะท้อนถึงปัญหาทางการเงินที่เริ่มมีความต่อเนื่องและมีแนวโน้มพัฒนาไปสู่สถานะหนี้เสียหากไม่ได้รับการจัดการอย่างทันทั่วถึง จำเป็นต้องพิจารณามาตรการเชิงลึก เช่น การปรับโครงสร้างหนี้ และกลุ่มสุดท้ายคือ **กลุ่มที่ 5** ซึ่งเป็น**กลุ่มความเสี่ยงวิกฤต (Critical Risk)** ประกอบด้วย ลูกหนี้ที่ผิดนัดชำระหนี้ติดต่อกันเป็นเวลา 3 เดือน ซึ่งถือว่าอยู่ในภาวะที่มีโอกาสสูงในการกลายเป็นหนี้เสีย (Non-Performing Loan: NPL) และส่งผลกระทบต่อเสถียรภาพทางการเงินของสถาบันการเงิน และอาจต้องใช้มาตรการติดตามเข้มข้นหรือการดำเนินการทางกฎหมาย การจำแนกกลุ่มในลักษณะนี้จึงไม่ได้เป็นเพียงการจัดกลุ่มเพื่อวิเคราะห์โมเดลเชิงเทคนิคเท่านั้น แต่ยังเป็นแนวทางสำคัญในการกำหนดกลยุทธ์การบริหารความเสี่ยงเชิงรุกที่แตกต่างกันในแต่ละระดับความรุนแรงของปัญหา ทั้งในแง่ของการติดตาม การเจรจา การปรับโครงสร้างหนี้ หรือแม้กระทั่งการดำเนินการทางกฎหมายในกลุ่มที่มีความเสี่ยงสูงสุด

ผลลัพธ์จากการศึกษานี้จะช่วยสนับสนุนการตัดสินใจของสถาบันการเงินให้สามารถนำไปใช้ในการประเมินความเสี่ยงได้อย่างแม่นยำมากขึ้น สามารถนำไปประยุกต์ใช้ในการกำหนดนโยบายสินเชื่อที่เหมาะสมกับระดับความเสี่ยงของลูกหนี้แต่ละกลุ่ม เป็นเครื่องมือสำคัญที่ช่วยในการบริหารความเสี่ยงสินเชื่อในเชิงรุก โดยเฉพาะอย่างยิ่งในการติดตามสถานการณ์หนี้เสีย (NPL) และการวางแผนมาตรการรับมือได้ทันทั่วถึง การนำผลลัพธ์เหล่านี้ไปใช้งานจะช่วยลดความสูญเสียทางการเงินที่เกิดจากหนี้เสีย ทั้งยังช่วยปรับปรุงประสิทธิภาพการจัดสรรทรัพยากรและลดต้นทุนในการบริหารหนี้นอกจากนี้ยังมีส่วนสนับสนุนการทำงานของธนาคารแห่งประเทศไทยในการตรวจสอบและกำกับดูแลสถาบันการเงินในเชิงรุกยิ่งขึ้น ช่วยส่งเสริมให้เกิดระบบการเงินที่มีความโปร่งใส ตรวจสอบได้ และมีประสิทธิภาพ อันนำไปสู่การสร้างเสถียรภาพและความมั่นคงทางการเงินที่ยั่งยืนต่อเศรษฐกิจไทยในระยะยาว (อโณทัย พุทธิสารี และคณะ, 2561)

1.2 วัตถุประสงค์ของงานวิจัย

- 1) สร้างแบบจำลองที่เหมาะสมสำหรับกลุ่มลูกค้าที่มีพฤติกรรมผิดนัดชำระหนี้ที่แตกต่างกัน
- 2) วิเคราะห์ปัจจัยที่มีอิทธิพลต่อการผิดนัดชำระหนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.3 ขอบเขตการวิจัย

1) ขอบเขตของแผนงาน

งานวิจัยนี้ศึกษาเกี่ยวกับการเปรียบเทียบแบบจำลองทำนายการผิดนัดชำระหนี้ โดยใช้ข้อมูลสินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพจากสถาบันการเงินแห่งหนึ่ง ข้อมูลที่ใช้จะครอบคลุมเฉพาะบัญชีสินเชื่อที่ยังมีการเคลื่อนไหวในเดือนมกราคม พ.ศ. 2567

2) ขอบเขตของข้อมูล

ข้อมูลที่ใช้วิเคราะห์ประกอบด้วย 2 ส่วนหลัก ได้แก่ ข้อมูลทั่วไปของลูกค้าหนี้ที่บัญชียังมีการเคลื่อนไหวในเดือนมกราคม พ.ศ. 2567 โดยเลือกบัญชีที่เปิดมาแล้วไม่น้อยกว่า 1 ปีนับจากเดือนธันวาคม พ.ศ. 2566 และข้อมูลด้านพฤติกรรมการชำระหนี้ย้อนหลังระหว่างเดือนมกราคม พ.ศ. 2566 ถึงเดือนธันวาคม พ.ศ. 2566 ทั้งนี้แบบจำลองถูกพัฒนาเพื่อทำนายสถานะของลูกค้าหนี้ในเดือนมกราคม พ.ศ. 2567 โดย 0 หมายถึง ลูกค้าหนี้ไม่ผิดนัดชำระหนี้ (Non-Default) และ 1 หมายถึง ลูกค้าหนี้ผิดนัดชำระหนี้ (Default) ณ สิ้นเดือนมกราคม พ.ศ. 2567

3) ขอบเขตด้านเวลา

นำข้อมูลพฤติกรรมของลูกค้าหนี้ย้อนหลังตลอดปี พ.ศ. 2566 มาเพื่อสร้างตัวแปรอิสระที่เกี่ยวข้องกับพฤติกรรมการชำระหนี้ในช่วง 1 ปีที่ผ่านมา

4) เครื่องมือที่ใช้ในการสร้างแบบจำลอง

การพัฒนาแบบจำลองการทำนายการผิดนัดชำระหนี้ดำเนินการผ่านโปรแกรมภาษาไพธอน (Python Programming) โดยใช้ Jupyter Notebook เป็นแพลตฟอร์มหลักในการพัฒนาแบบจำลอง และใช้ SAS Program สำหรับการจัดเก็บและเตรียมข้อมูลที่เกี่ยวข้อง

1.4 ประโยชน์ที่คาดว่าจะได้รับ

1) สร้างแบบจำลองที่ช่วยให้สถาบันการเงินใช้ในการวางแผนและกำหนดกลยุทธ์ให้สอดคล้องกับนโยบายการช่วยเหลือลูกหนี้ของภาครัฐและธนาคารแห่งประเทศไทย

2) สถาบันการเงินสามารถนำผลจากแบบจำลองไปใช้ในการประเมินความเสี่ยงและคัดกรองลูกค้าได้อย่างมีประสิทธิภาพมากขึ้น ซึ่งช่วยให้สามารถวางแผนบริหารจัดการและรับมือกับปัญหาหนี้เสีย (NPL) ได้อย่างเหมาะสมมากยิ่งขึ้น

1.5 กรอบแนวคิดการวิจัย

จากการศึกษาข้อมูลสินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพ ประกอบกับการทบทวนวรรณกรรมและงานวิจัยที่เกี่ยวข้อง ทำให้สามารถระบุปัจจัยที่เกี่ยวข้องกับพฤติกรรมการผิดนัดชำระหนี้ได้อย่างหลากหลาย ทั้งในด้านคุณลักษณะส่วนบุคคล ลักษณะของสินเชื่อ ซึ่งสามารถนำมาใช้เป็นตัวแปรอิสระ (Independent Variables) นอกจากนี้ เพื่อให้การประเมินแบบจำลองมีความครอบคลุมทั้งในเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ระดับภาพรวมและเฉพาะกลุ่ม งานวิจัยนี้ได้ดำเนินการแบ่งกลุ่มลูกหนี้ตามประวัติการผิดนัดชำระหนี้ ออกเป็น 5 กลุ่ม ดังนี้

ตัวแปรต้น (Independent Variables)



รูปที่ 1.2 กรอบแนวคิดในการวิจัย

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1.6 นิยามศัพท์เฉพาะ

การผลิตนัดชำระหนี้และไม่ผลิตนัดชำระหนี้ ในงานวิจัยนี้ คือ **ผลิตนัดชำระหนี้** (1) หมายถึง กรณีที่ลูกหนี้ไม่ชำระหนี้ภายในสิ้นเดือนที่ถึงกำหนด หรือมีการชำระหนี้บางส่วนซึ่งไม่ครบตามจำนวนค่างวด ในทางกลับกัน **“ไม่ผลิตนัดชำระหนี้”** (0) หมายถึง ลูกหนี้ที่ชำระหนี้ครบจำนวนภายในเดือนเดียวกันกับวันที่ครบกำหนด ไม่ว่าจะชำระตรงเวลา ก่อนกำหนด หรือภายหลังวันที่กำหนดเล็กน้อย แต่ยังคงอยู่ในเดือนเดียวกัน

ตารางที่ 1.1 นิยามของสถานการณ์การผลิตนัดชำระหนี้

รหัสบัญชี	วันกำหนดชำระหนี้ในแต่ละงวด	วันที่ลูกค้าชำระจริงในงวดนั้น	รูปแบบการชำระ	สถานะการผลิตนัดชำระหนี้ ณ สิ้นเดือน (Y)	วันกำหนดชำระค่างวดครั้งถัดไป	หมายเหตุ
C01	02/01/67	02/01/67	ชำระเต็ม	ไม่ผลิตนัด	02/02/67	ชำระตรงเวลา
C02	05/01/67	30/01/67	ชำระเต็ม	ไม่ผลิตนัด	05/02/67	ชำระทันเดือน
C03	10/01/67	31/01/67	ชำระบางส่วน	ผลิตนัด	10/01/67	ชำระไม่ครบยอด
C04	10/01/67	-	ไม่ชำระ	ผลิตนัด	10/01/67	ไม่มีการชำระ
C05	15/01/67	12/01/67	ชำระเต็ม	ไม่ผลิตนัด	15/02/67	ชำระก่อนกำหนด
C06	15/01/67	29/01/67	ชำระบางส่วน	ผลิตนัด	15/01/67	ชำระไม่ครบยอด
C07	20/01/67	28/01/67	ชำระเต็ม	ไม่ผลิตนัด	20/02/67	ชำระทันเดือน
C08	25/01/67	05/02/67	ชำระเต็ม	ผลิตนัด	25/01/67	ชำระหลังสิ้นเดือน
C09	28/01/67	30/01/67	ชำระบางส่วน	ผลิตนัด	28/01/67	ชำระไม่ครบยอด
C10	30/01/67	-	ไม่ชำระ	ผลิตนัด	30/01/67	ไม่มีการชำระ

ตารางที่ 1.1 แสดงตัวอย่างสถานการณ์การชำระหนี้เพื่อประกอบการพิจารณานิยามดังกล่าว ตัวอย่างเช่น กรณีบัญชี C08 ลูกหนี้มีกำหนดชำระหนี้วันที่ 25 มกราคม พ.ศ. 2567 แต่ชำระเต็มจำนวน ในวันที่ 5 กุมภาพันธ์ พ.ศ. 2567 ซึ่งเป็นการชำระข้ามเดือน แม้จะชำระครบถ้วน แต่เนื่องจากไม่ได้ชำระภายในเดือนที่ครบกำหนด จึงถือเป็น ผลิตนัดชำระหนี้ ณ สิ้นเดือนมกราคม พ.ศ. 2567 สะท้อนให้เห็นเกณฑ์การจัดสถานะที่ให้ความสำคัญกับ ช่วงเวลาในการชำระ มากกว่าจำนวนที่ชำระเพียงอย่างเดียว

สินเชื่อรายย่อยเพื่อการประกอบอาชีพภายใต้การกำกับ คือ สินเชื่อที่ภาครัฐต้องการให้การสนับสนุน เนื่องจากเห็นปัญหาของประชาชนรายย่อยที่ต้องการใช้เงินทุนเพื่อการประกอบอาชีพ แต่ไม่สามารถเข้าถึงสินเชื่อในระบบสถาบันการเงินได้ (ธนาคารแห่งประเทศไทย, 2558)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

หนี้เสีย หรือ NPL (Non-performing Loan) คือ หนี้ที่ค้างชำระเกิน 90 วัน หรือไม่ได้ชำระคืนตามข้อกำหนดตกลงกันตามสัญญา (ธนาคารกรุงไทย, 2566)

กลุ่มที่ 1 ในงานวิจัยนี้ คือ กลุ่มลูกค้าทั่วไป (General Customer) ซึ่งครอบคลุมลูกหนี้ทั้งหมดโดยไม่จำกัดพฤติกรรมหรือประวัติการผิดนัดชำระหนี้

กลุ่มที่ 2 ในงานวิจัยนี้ คือ กลุ่มความเสี่ยงต่ำ (Low Risk) ซึ่งประกอบด้วยลูกหนี้ที่ผิดนัดชำระหนี้เพียง 1 ครั้งในช่วง 2 เดือน

กลุ่มที่ 3 ในงานวิจัยนี้ คือ กลุ่มความเสี่ยงปานกลาง (Moderate Risk) ประกอบด้วยลูกหนี้ที่ผิดนัดชำระหนี้ในเดือนล่าสุด

กลุ่มที่ 4 ในงานวิจัยนี้ คือ กลุ่มความเสี่ยงสูง (High Risk) ซึ่งประกอบด้วยลูกหนี้ที่ผิดนัดชำระหนี้ต่อเนื่องกันเป็นเวลา 2 เดือน

กลุ่มที่ 5 ในงานวิจัยนี้ คือ กลุ่มความเสี่ยงวิกฤต (Critical Risk) ประกอบด้วยลูกหนี้ที่ผิดนัดชำระหนี้ต่อเนื่องกันเป็นเวลา 3 เดือน



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 2

ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

ในการศึกษาเรื่อง การทำนายการผิวน้ำของสินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพ ผู้วิจัยได้ศึกษาค้นคว้าแนวคิด ทฤษฎีและงานวิจัยที่เกี่ยวข้องต่างๆ ดังต่อไปนี้

- 2.1 แนวคิดและทฤษฎีเกี่ยวข้องกับปัจจัยส่วนบุคคล
- 2.2 แนวคิดและทฤษฎีเกี่ยวกับปัจจัยด้านลักษณะสินเชื่อ
- 2.3 แนวคิดและทฤษฎีเกี่ยวข้องกับกระบวนการวิเคราะห์ข้อมูลแบบ CRISP-DM
- 2.4 แนวคิดและทฤษฎีเกี่ยวกับการแบ่งข้อมูลเพื่อทำการทดสอบและการสุ่มตัวอย่างข้อมูล
- 2.5 แนวคิดและทฤษฎีเกี่ยวข้องกับการจำแนกประเภทของข้อมูล (Classification)
- 2.6 แนวคิดและทฤษฎีเกี่ยวข้องกับการเลือกไฮเปอร์พารามิเตอร์
- 2.7 แนวคิดและทฤษฎีเกี่ยวข้องกับการประเมินประสิทธิภาพในการทำนายของแบบจำลอง
- 2.8 แนวคิดและทฤษฎีเกี่ยวข้องกับความสำคัญของคุณลักษณะ (Feature Importance)
- 2.9 งานวิจัยที่เกี่ยวข้อง

2.1 แนวคิดและทฤษฎีเกี่ยวกับปัจจัยส่วนบุคคล

โคมลักษณะ สุวรรณกาญจน์ และคณะ (2021) ได้กล่าวถึงปัจจัยส่วนบุคคลว่า ปัจจัยด้านประชากรศาสตร์หรือปัจจัยส่วนบุคคลเป็นปัจจัยหนึ่งที่ส่งผลต่อพฤติกรรมของ ผู้บริโภค การตัดสินใจก่อนพฤติกรรมต่างๆ มักจะมีความแตกต่างกันเมื่อผู้ก่อพฤติกรรมมีปัจจัยส่วนบุคคลที่แตกต่างกัน

เปรมกมล หงษ์ยนต์ (2019) ได้กล่าวว่า ลักษณะประชากรศาสตร์ คือ ประชากรศาสตร์เป็นปัจจัยสำคัญที่นักการตลาดนิยมนำมาใช้เป็นเกณฑ์ในการแบ่งส่วนตลาด เช่น เพศ อายุ ระดับการศึกษา รายได้ เป็นต้น บุคคลที่มีลักษณะทางประชากรศาสตร์ที่แตกต่างกัน ย่อมมีความคิดทัศนคติ และพฤติกรรมที่แตกต่างกัน

ปัญญารัตน์ หนูสิงค์ (2019) ปัจจัยที่มีผลต่อการเกิดหนี้ที่ไม่ก่อให้เกิดรายได้ของลูกหนี้ ธนาคารพัฒนาวิสาหกิจขนาดกลางและขนาดย่อมแห่งประเทศไทย สาขาพัทลุง พบว่าปัจจัยส่วนบุคคลประกอบด้วย เพศ อายุ สถานภาพสมรส ระดับการศึกษา อาชีพ อายุการทำงาน จำนวนสมาชิกในครัวเรือน รายได้เฉลี่ยต่อครัวเรือน รายจ่ายเฉลี่ยต่อเดือนของครัวเรือน ส่งผลต่อการเกิดหนี้ที่ไม่ก่อให้เกิดรายได้

ลูกหนี้ที่มีปัจจัยส่วนบุคคลต่างกัน มีปัจจัยสาเหตุของการเกิดหนี้ที่ไม่ก่อให้เกิดรายได้แตกต่างกัน โดยลูกหนี้ค้ำชำระที่มีระดับการศึกษาต่างกัน มีปัจจัยสาเหตุของการเกิดหนี้ค้ำชำระเนื่องจากลูกหนี้ประสบภาวะการขาดทุนและลูกหนี้มีหนี้สินหลายทางแตกต่างกัน และลูกหนี้ที่มีอาชีพต่างกัน มี

ปัจจัยสาเหตุของการเกิดหนี้ค้ำชำระเนื่องจากลูกหนี้มีภาระค่าใช้จ่ายในครอบครัวสูง เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากแนวคิดและทฤษฎีเกี่ยวกับปัจจัยส่วนบุคคลสรุปได้ว่า ปัจจัยส่วนบุคคลประกอบด้วย เพศ อายุ สถานภาพ ที่อยู่อาศัย ซึ่งเป็นลักษณะของบุคคลที่มีความแตกต่างกันและมีพฤติกรรมในการตอบสนองต่อความต้องการที่แตกต่างกัน

2.2 แนวคิดและทฤษฎีเกี่ยวกับปัจจัยด้านลักษณะสินเชื่อ

สินเชื่อ คือ เงินที่กู้ยืมจากสถาบันการเงินซึ่งจะต้องได้รับอนุญาตตามกฎหมาย โดยจะมีการทำหนังสือสัญญาข้อตกลงเงื่อนไขต่างๆ ซึ่งจะเป็นไปตามที่สถาบันการเงินผู้ให้สินเชื่อเป็นผู้กำหนด โดยมีผลตอบแทนเป็นดอกเบี้ย สินเชื่อที่สถาบันการเงินมีให้บริการมีมากมายหลายประเภท และเงื่อนไขในการกู้ยืม และวัตถุประสงค์ของการให้สินเชื่อ รวมถึงอัตราดอกเบี้ยก็จะแตกต่างกันไปตามแต่ละสถาบันการเงินเจ้าของผลิตภัณฑ์สินเชื่อจะกำหนด (ยูเมะพลัส, 2022)

2.2.1 ประเภทของสินเชื่อ

สินเชื่อสามารถแบ่งแยกประเภทได้หลายรูปแบบ เช่น แบ่งตามประเภทของผู้ขอสินเชื่อ (สินเชื่อส่วนบุคคล สินเชื่อนิติบุคคล และสินเชื่อภาครัฐ) แบ่งตามวัตถุประสงค์ของการขอสินเชื่อ (สินเชื่ออเนกประสงค์ สินเชื่อบ้าน และสินเชื่อรถ) แบ่งตามระยะเวลา (สินเชื่อระยะสั้น สินเชื่อระยะยาว) แต่โดยหลักใหญ่ๆ แล้วเรามักแบ่งตามประเภทของสินเชื่อตามลักษณะของผู้ขอสินเชื่อ โดยสินเชื่อแบ่งตามลักษณะของผู้ขอสินเชื่อได้ดังนี้

1) สินเชื่อส่วนบุคคล หรือ Personal Loan เป็นสินเชื่อเพื่อซื้อสินค้าสำหรับการอุปโภคและบริโภคหรือบริการ ตามความต้องการของผู้ขอสินเชื่อ สามารถแยกย่อยออกตามวัตถุประสงค์ของสินเชื่อได้อีก เช่น สินเชื่ออเนกประสงค์ สินเชื่อรถ สินเชื่อที่อยู่อาศัย สินเชื่อเพื่อการศึกษา รวมไปถึงผลิตภัณฑ์บัตรเครดิต และบัตรกดเงินสด ก็จัดเป็นสินเชื่อประเภทหนึ่งเช่นกัน สินเชื่อส่วนบุคคลยังสามารถแบ่งแยกย่อยออกเป็น สินเชื่อส่วนบุคคลแบบไม่มีบัตร และสินเชื่อส่วนบุคคลแบบมีบัตรได้อีก หลักทรัพย์หรือผู้ค้ำประกันในการยื่นขอสินเชื่ออาจมีหรือไม่มีก็ได้ขึ้นอยู่กับเงื่อนไขของแต่ละสถาบันการเงินซึ่งได้รับอนุญาตจากกระทรวงการคลัง และอยู่ภายใต้การดูแลของธนาคารแห่งประเทศไทยเป็นผู้กำหนด

1.1) สินเชื่อส่วนบุคคลแบบไม่มีบัตร คือ สินเชื่อเงินก้อนอเนกประสงค์ซึ่งผู้ขอสินเชื่อจะต้องยื่นกู้กับทางสถาบันการเงินผู้ให้สินเชื่อโดยตรงและผู้ขอสินเชื่อจะเป็นผู้ระบุจำนวนวงเงินที่ต้องการไว้ในเอกสารตั้งแต่ดำเนินการยื่นขอสินเชื่อ หลังจากสินเชื่อผ่านการอนุมัติแล้วเงินก้อนตามจำนวนที่ได้ระบุตามเอกสารตอนยื่นขอสินเชื่อจะถูกโอนเข้าบัญชีของผู้ขอสินเชื่อโดยตรงภายในครั้งเดียว และผู้ขอสินเชื่อเจ้าของบัญชีธนาคารที่ได้รับโอนเงินก้อนนั้นก็สามารถนำเงินไปใช้จ่ายได้ทันที

1.2) สินเชื่อส่วนบุคคลแบบมีบัตร เช่น บัตรเครดิต และบัตรกดเงินสด

ก) บัตรเครดิต คือ บัตรสินเชื่อที่สถาบันการเงินออกให้กับลูกค้าเพื่อใช้ในการจ่ายค่าสินค้าแทนเงินสดซึ่งเป็นการนำเงินมาใช้จ่ายล่วงหน้า โดยวงเงินใช้จ่ายนั้น เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ต้องไม่เกินวงเงินที่สถาบันการเงินผู้ออกบัตรอนุมัติ บัตรเครดิตยังสามารถใช้เบิกถอนเงินสดจากตู้ ATM แต่อาจมีดอกเบี้ยและค่าธรรมเนียมการเบิกถอนเงินสดเพิ่มเติม ซึ่งมีเงื่อนไขเป็นไปตามที่แต่ละสถาบันการเงินเจ้าของผลิตภัณฑ์กำหนด อีกทั้งบัตรเครดิตยังสามารถใช้ผ่อนชำระสินค้า หรือบริการเป็นรายงวดตามโปรโมชั่นของร้านค้าได้อีกด้วย

ข) บัตรกดเงินสด เป็นอีกรูปแบบหนึ่งของสินเชื่อส่วนบุคคล ที่สถาบันการเงินเป็นผู้ออกให้โดยบัตรกดเงินสดสามารถใช้ได้ 2 รูปแบบ คือ ใช้กดเงินสด และ ใช้ผ่อนชำระสินค้า

2) สินเชื่อนิติบุคคล นิติบุคคล คือบริษัท ห้างหุ้นส่วน สมาคม มูลนิธิ วัด กระทรวง หรือจังหวัด ที่กฎหมายสมมติให้มีเช่นเดียวกับบุคคลธรรมดา สามารถทำนิติกรรมสัญญา มีสิทธิเป็นเจ้าของทรัพย์สินต่างๆ โดยสินเชื่อนิติบุคคลคือ สินเชื่อที่บริษัท ห้าง ร้านที่มีสถานะเป็นนิติบุคคลเป็นผู้ขอสินเชื่อเพื่อลงทุนหรือขยายธุรกิจ มีกำหนดระยะเวลาชำระเงินต้นและดอกเบี้ยที่แน่นอน

3) สินเชื่อภาครัฐ สินเชื่อที่สถาบันทางการเงินปล่อยกู้ให้กับภาครัฐ มักอยู่ในรูปแบบของพันธบัตรชดเชยการขาดดุลงบประมาณ แต่ปัจจุบันไม่ค่อยพบสินเชื่อประเภทนี้แล้ว

เอกกมล เอี่ยม ศรี (2018) ได้กล่าวว่า การวิเคราะห์สินเชื่อเป็นขั้นตอนที่สำคัญซึ่งในการพิจารณาจะทำการวิเคราะห์จากข้อมูลทางการเงิน เพื่อช่วยให้สถาบันการเงินหรือผู้ให้สินเชื่อประมาณการถึงความสามารถในการชำระหนี้คืนของผู้ขอสินเชื่อ ช่วยลดความเสี่ยงหนี้ที่ไม่ก่อให้เกิดรายได้และปัญหาหนี้สูญที่จะตามมา ดังนั้นงบกระแสเงินสดจะเป็นข้อมูลสนับสนุนการประมาณการทางการเงินของผู้ขอสินเชื่อ เนื่องจากข้อมูลแสดงถึงศักยภาพทางการเงิน และความยั่งยืนของกิจการ

2.2.2 ปัจจัยที่ใช้ในการพิจารณาสินเชื่อ

บริษัท โสมบายเออร์ไกด์ จำกัด (2017) กล่าวว่า การพิจารณาอนุมัติสินเชื่อขึ้นอยู่กับนโยบายและหลักเกณฑ์ของแต่ละธนาคาร โดยทั่วไปจะมีปัจจัยหลักๆ ที่ใช้ประกอบการพิจารณา เช่น นโยบายสินเชื่อของธนาคาร ซึ่งธนาคารบางแห่งอาจกำหนดว่าผู้ยื่นกู้ต้องไม่มีประวัติการค้างชำระในช่วง 1-2 เดือนย้อนหลัง หรืองดให้สินเชื่อแก่ลูกค้าใหม่ในกลุ่มอาชีพอุตสาหกรรมที่มีความเสี่ยงสูง หรือกำหนดวัตถุประสงค์ในการขอสินเชื่อ เช่น ใช้เป็นเงินหมุนเวียนในการทำธุรกิจ นอกจากนี้ยังมีหลัก 5Cs ที่ธนาคารให้ความสำคัญและนำมาพิจารณาซึ่งประกอบด้วย

1) Character อุปนิสัยของลูกค้า ถือเป็นคุณลักษณะและความน่าเชื่อถือของผู้ขอสินเชื่อ เพื่อดูว่าตัวผู้กู้ มีวินัยการใช้เงิน และประวัติการชำระสินเชื่อในอดีตเป็นอย่างไร ซึ่งจะบอกถึงความสามารถในการใช้หนี้ เช่น ในกรณีบุคคลธรรมดาอาจพิจารณาจากอายุ อาชีพ สถานภาพสมรส ส่วนกรณีผู้ขอสินเชื่อประกอบธุรกิจอาจพิจารณาจากประเภทของธุรกิจ ประวัติของผู้บริหาร และระยะเวลาในการดำเนินธุรกิจ เป็นต้น

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- 2) Capacity ความสามารถในการชำระหนี้ ถือเป็นหัวใจสำคัญในการพิจารณาสินเชื่อ โดยความสามารถในการจ่ายชำระหนี้คืนได้ตามระยะเวลาที่กำหนดไว้ บวกความมั่นคงของรายได้หลังหักค่าใช้จ่ายต่างๆ แล้ว ต้องเหลือเพียงพอชำระหนี้ทั้งหมดในแต่ละเดือนด้วย
- 3) Capital เงินทุน นับว่าเงินทุน สินทรัพย์ เงินฝากของผู้กู้ หรือ ผู้ขอสินเชื่อถือเป็นหลักประกันให้กู้ยืม ซึ่งจำเป็นอย่างยิ่งในสินเชื่อธุรกิจ แม้ว่าสินทรัพย์เหล่านี้จะไม่ใช่ว่าแหล่งเงินสำหรับชำระหนี้ แต่จะเป็นแหล่งเงินสำรองสำหรับการชำระหนี้ของผู้ขอสินเชื่อ ในกรณีเกิดปัญหาไม่สามารถชำระหนี้ได้
- 4) Collateral หลักประกัน คือผู้ค้ำประกันหรือหลักประกันที่ผู้ขอสินเชื่อนำมาจำนำ หรือ จำนองเพื่อให้สถาบันการเงินเกิดความมั่นใจ และลดความเสี่ยงหากผู้ขอสินเชื่อไม่ชำระหนี้ตามกำหนด ซึ่งสามารถให้ผู้ค้ำประกันชำระหนี้แทน หรือนำหลักประกันมาขายทอดตลาดได้ตามที่กฎหมายกำหนด
- 5) Condition เงื่อนไข คือ ปัจจัยภายนอกที่มีผลกระทบต่อรายได้ของผู้ขอสินเชื่อ เช่น เศรษฐกิจ เงินเฟ้อ ความมั่นคงในรายได้ รวมถึงการงาน ปัญหาสังคม และสิ่งแวดล้อม ที่มีผลกระทบต่อความเป็นไปได้ของโครงการลงทุนหรือรายได้ของผู้ขอสินเชื่อ ซึ่งจะมีผลต่อความสามารถในการชำระหนี้

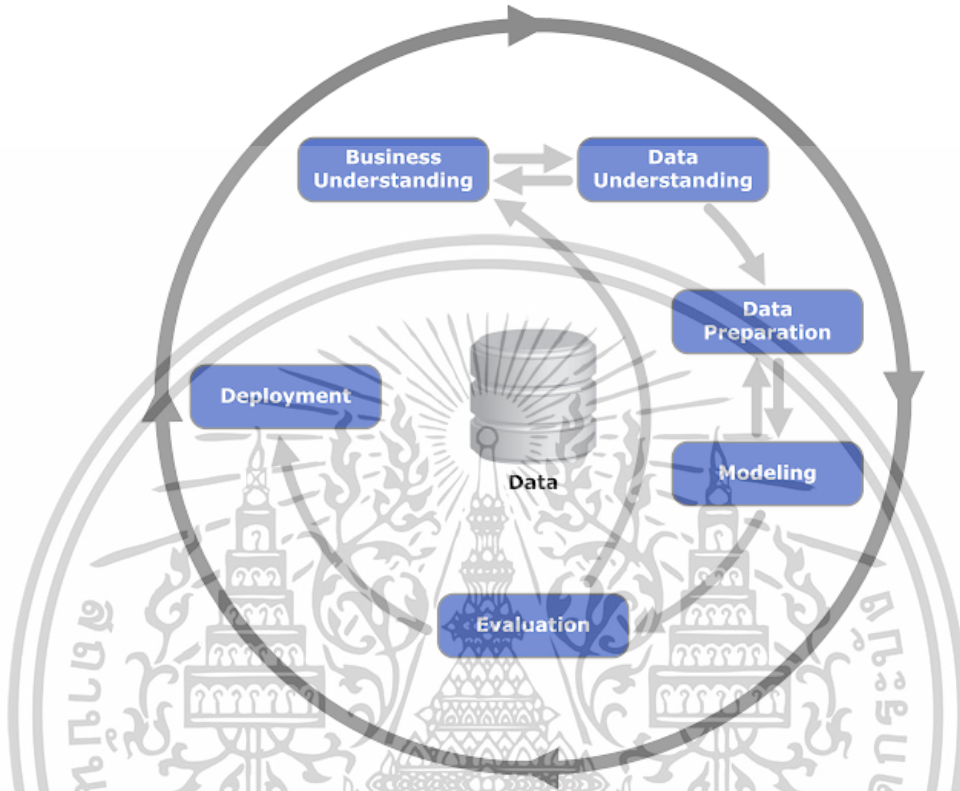
จาก 5Cs ข้างต้น ปัจจุบันเกือบทุกธนาคารจะเน้นไปที่ความสามารถในการชำระหนี้ (Ability of Repay) และความตั้งใจในการชำระหนี้ (Willingness to repay) เป็นหลัก หมายถึง รายได้ของผู้กู้สินเชื่อ โดยผู้กู้อาจเป็นได้ทั้งมนุษย์เงินเดือนที่มี “รายได้ประจำ” กับ “ผู้ประกอบการอาชีพอิสระ” ที่แสดงหลักฐานรายได้แตกต่างกัน ซึ่งการชำระหนี้คืนให้กับธนาคารส่วนใหญ่จะชำระกันเป็นงวดเป็นประจำอย่างสม่ำเสมอ หรือเดือนละครั้ง ดังนั้นรายได้เป็นประจำมีความหมายมากกว่ารายได้ที่ไม่แน่นอน ยิ่งถ้าผู้กู้มีอาชีพการงาน มีรายได้ที่มั่นคง ไม่เคยมีหนี้สินล้นพ้นตัว และมีความสามารถในการชำระหนี้ การขอสินเชื่อคงไม่มีปัญหาอะไร แต่ถ้าไม่มีหน้าที่การงานที่มั่นคง ไม่มีเอกสารใบรับรองเงินเดือน โอกาสได้เงินกู้ยังคงมีอยู่ เพียงแต่ต้องใช้หลักฐานทางการเงินอื่นแทน เช่น ใบแจ้งยอดคงเหลือ หรือ statement ของบัญชีเงินฝาก

จากแนวคิดและทฤษฎีเกี่ยวกับปัจจัยด้านลักษณะสินเชื่อสรุปได้ว่า ปัจจัยด้านลักษณะสินเชื่อเกือบทุกสถาบันการเงินจะเน้นไปที่ความสามารถในการชำระหนี้ และความตั้งใจในการชำระหนี้ รวมทั้งระยะเวลาการชำระหนี้ ประเภทการขอสินเชื่อ วงเงินสินเชื่อ ซึ่งเป็นลักษณะด้านสินเชื่อที่มีลักษณะเฉพาะของในแต่ละบุคคล

2.3 แนวคิดและทฤษฎีเกี่ยวข้องกับกระบวนการวิเคราะห์ข้อมูลแบบ CRISP-DM

CRISP-DM (Cross-Industry Standard Process for Data Mining) เป็นกรอบการทำงานที่พัฒนาขึ้นในช่วงปลายปี 1990 โดยมีเป้าหมายเพื่อเป็นมาตรฐานสำหรับการทำเหมืองข้อมูล (Data Mining) เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Mining) และการวิเคราะห์ข้อมูลในอุตสาหกรรมหลากหลายประเภท กรอบการทำงานนี้ให้แนวทางที่ชัดเจนและมีโครงสร้างสำหรับการดำเนินโครงการข้อมูล ตั้งแต่การเข้าใจปัญหาทางธุรกิจไปจนถึงการนำผลลัพธ์ไปใช้งานจริง (Chapman et al., 2000; Shearer, 2000)



รูปที่ 2.1 กระบวนการวิเคราะห์ข้อมูลแบบ CRISP-DM
ที่มา : Chapman, et al. (2000)

CRISP-DM ถือเป็นมาตรฐานกลางสำหรับการดำเนินโครงการเหมืองข้อมูลและการวิเคราะห์ข้อมูลที่สามารถนำไปประยุกต์ใช้ได้กับหลากหลายอุตสาหกรรม ไม่ว่าจะเป็นธุรกิจการเงิน การตลาด พาณิชยอิเล็กทรอนิกส์ การแพทย์ และอุตสาหกรรมอื่นๆ ที่ต้องพึ่งพิงการวิเคราะห์ข้อมูลเพื่อสร้างคุณค่าและสนับสนุนการตัดสินใจเชิงกลยุทธ์ (Wirth & Hipp, 2000) กรอบแนวคิดนี้ได้ถูกออกแบบมาเพื่อให้ทุกฝ่ายที่เกี่ยวข้องในโครงการข้อมูล ตั้งแต่ผู้บริหาร นักวิเคราะห์ข้อมูล ไปจนถึงผู้พัฒนาแบบจำลอง สามารถเข้าใจขั้นตอนการทำงานในภาพรวมได้ตรงกัน และสามารถดำเนินงานอย่างเป็นระบบ โดยจากรูปกระบวนการวิเคราะห์ข้อมูลแบบ CRISP-DM มีทั้งหมด 6 ขั้นตอนดังนี้

2.3.1 การทำความเข้าใจธุรกิจ (Business Understanding)

การทำความเข้าใจธุรกิจ (Business Understanding) เป็นขั้นตอนแรกและเป็นรากฐานสำคัญในกระบวนการวิเคราะห์ข้อมูลตามแนวทาง CRISP-DM โดยมีวัตถุประสงค์เพื่อให้ผู้ดำเนินโครงการสามารถระบุปัญหาทางธุรกิจ วิเคราะห์ความต้องการขององค์กร และกำหนดเป้าหมายเชิงกลยุทธ์ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ยุทธ์ที่สอดคล้องกับภารกิจและวิสัยทัศน์ขององค์กรอย่างชัดเจน (Chapman et al., 2000; Shearer, 2000) กระบวนการในขั้นตอนนี้เริ่มต้นด้วยการศึกษาและทำความเข้าใจบริบททางธุรกิจ อาทิเช่น สถานการณ์การแข่งขันในตลาด กลยุทธ์การดำเนินงาน จุดแข็ง จุดอ่อน และโอกาสขององค์กร รวมถึง การระบุผู้มีส่วนได้ส่วนเสียหลัก (Key Stakeholders) และข้อจำกัดที่เกี่ยวข้อง (Wirth & Hipp, 2000) ทั้งนี้การนิยามปัญหาทางธุรกิจให้ชัดเจนจะช่วยให้การวิเคราะห์ข้อมูลมีทิศทางที่ตรงกับวัตถุประสงค์จริงขององค์กร จากนั้นจะต้องทำการกำหนดวัตถุประสงค์ของโครงการวิเคราะห์ข้อมูล (Data Mining Goals) โดยแปลงเป้าหมายทางธุรกิจที่มีลักษณะกว้างและเป็นนามธรรม ให้กลายเป็น เป้าหมายเชิงปฏิบัติที่สามารถประเมินผลได้ เช่น การทำนายโอกาสผิदनัดชำระหนี้ของลูกค้า การคัดกรองกลุ่มลูกค้าที่มีความเสี่ยงสูง หรือการเพิ่มอัตราการตอบกลับของแคมเปญการตลาด ซึ่งต้องมีการ กำหนดตัวชี้วัด (Key Performance Indicators: KPIs) อย่างชัดเจน (Chapman et al., 2000)

โดยสรุป การทำความเข้าใจธุรกิจเป็นขั้นตอนสำคัญที่ช่วยวางรากฐานและทิศทางของ โครงการวิเคราะห์ข้อมูล ทำให้การดำเนินงานในแต่ละขั้นตอนถัดไปเป็นไปอย่างมีประสิทธิภาพและ ตอบโจทย์ทางธุรกิจอย่างแท้จริง (Shearer, 2000; Wirth & Hipp, 2000)

2.3.2 การทำความเข้าใจข้อมูล (Data Understanding)

ถือเป็นขั้นตอนสำคัญลำดับต้นในกระบวนการ CRISP-DM ที่มีบทบาทในการเชื่อมโยง เป้าหมายทางธุรกิจกับข้อมูลที่มีอยู่จริง โดยมุ่งเน้นให้ผู้วิจัยหรือนักวิเคราะห์ข้อมูลสามารถเข้าถึงและ เข้าใจข้อมูลอย่างลึกซึ้งซึ่งก่อนที่จะดำเนินการวิเคราะห์หรือสร้างแบบจำลองในขั้นตอนถัดไป (Chapman et al., 2000) กระบวนการทำความเข้าใจข้อมูลนั้นเริ่มจากการรวบรวมข้อมูลเบื้องต้น จากแหล่งข้อมูลต่างๆ ที่เกี่ยวข้องกับปัญหาทางธุรกิจ ไม่ว่าจะเป็นข้อมูลภายในองค์กร เช่น ฐานข้อมูล ลูกค้า ประวัติธุรกรรม หรือข้อมูลภายนอก เช่น ข้อมูลจากสถาบันภายนอก หรือแหล่งข้อมูล สาธารณะ (Shearer, 2000) หลังจากนั้นจะเข้าสู่การอธิบายคุณลักษณะของข้อมูลแต่ละตัวแปร ทั้งใน ด้านประเภทข้อมูล (เชิงปริมาณ, เชิงกลุ่ม) ช่วงของข้อมูล การกระจายตัว และแนวโน้มของข้อมูล (Wirth & Hipp, 2000) ในขั้นตอนนี้ยังรวมถึงการสำรวจข้อมูลเชิงลึกผ่านเทคนิค Visualization เช่น การสร้างกราฟแท่ง กราฟกระจาย หรือการใช้สถิติเชิงพรรณนา เพื่อช่วยให้เห็นโครงสร้างและรูปแบบ ข้อมูลที่สำคัญ อีกทั้งยังสามารถตรวจสอบความสัมพันธ์เบื้องต้นระหว่างตัวแปรต่างๆ ได้ชัดเจนยิ่งขึ้น นอกจากนี้ การตรวจสอบคุณภาพข้อมูล (Data Quality Verification) ก็เป็นส่วนสำคัญ อาทิเช่น การค้นหาค่าขาดหาย (Missing Values) การตรวจสอบค่าผิดปกติ (Outliers) ข้อมูลซ้ำซ้อน (Duplicates) หรือความไม่สมเหตุสมผลของข้อมูล ซึ่งปัญหาเหล่านี้หากไม่ได้รับการจัดการอาจส่งผล ให้ผลลัพธ์ของการวิเคราะห์คลาดเคลื่อนได้ (Chapman et al., 2000; Wirth & Hipp, 2000)

โดยสรุป การทำความเข้าใจข้อมูลเป็นขั้นตอนที่มีเป้าหมายเพื่อให้ผู้วิเคราะห์สามารถตัดสินใจ ได้อย่างถูกต้องในการเตรียมข้อมูลและเลือกแนวทางการวิเคราะห์ในขั้นตอนถัดไป อีกทั้งยังช่วยลด ความเสี่ยงของข้อผิดพลาดที่อาจเกิดขึ้นจากข้อมูลที่มีคุณภาพต่ำหรือไม่เหมาะสม (Shearer, 2000)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.3.3 การเตรียมข้อมูล (Data Preparation)

การเตรียมข้อมูล (Data Preparation) เป็นขั้นตอนสำคัญและใช้เวลามากที่สุดในกระบวนการวิเคราะห์ข้อมูลตามแนวทาง CRISP-DM เนื่องจากคุณภาพของข้อมูลที่เตรียมไว้อย่างเหมาะสมจะส่งผลโดยตรงต่อความแม่นยำและประสิทธิภาพของแบบจำลองวิเคราะห์ข้อมูลในขั้นตอนต่อไป (Chapman et al., 2000) หากข้อมูลขาดความสมบูรณ์ หรือมีปัญหาด้านคุณภาพ จะทำให้ผลการวิเคราะห์คลาดเคลื่อนและนำไปสู่ข้อสรุปที่ไม่ถูกต้อง ขั้นตอนการเตรียมข้อมูลประกอบด้วยกิจกรรมหลักหลายประการ ได้แก่

- 1) การเลือกข้อมูล (Data Selection) คัดเลือกเฉพาะชุดข้อมูล ตัวแปร หรือคุณลักษณะที่เกี่ยวข้องกับวัตถุประสงค์ของการวิเคราะห์ ทั้งนี้เพื่อขจัดข้อมูลที่ไม่จำเป็น ลดมิติของข้อมูล และป้องกันปัญหาข้อมูลซ้ำซ้อน
- 2) การทำความสะอาดข้อมูล (Data Cleaning) ตรวจสอบและจัดการกับค่าขาดหาย (Missing Values) ข้อมูลผิดปกติ (Outliers) และข้อมูลที่ขัดแย้งกัน (Inconsistencies) เช่น การเติมค่า การลบข้อมูล หรือการแก้ไขข้อผิดพลาด
- 3) การแปลงข้อมูล (Data Transformation) ดำเนินการปรับเปลี่ยนรูปแบบข้อมูลให้อยู่ในรูปแบบที่เหมาะสมต่อการวิเคราะห์ เช่น การแปลงข้อมูลเชิงกลุ่มเป็นตัวเลข (Encoding) การสเกลข้อมูล (Scaling/Normalization) หรือการแปลงข้อมูลเป็นช่วง (Binning)
- 4) การสร้างตัวแปรใหม่ (Feature Engineering) สร้างตัวแปรใหม่หรือดัดแปลงตัวแปรเดิม เพื่อเพิ่มศักยภาพในการทำนายของแบบจำลอง ตัวอย่างเช่น การรวมข้อมูลจากหลายฟิลด์ หรือการสร้างตัวแปรเชิงอนุพันธ์จากข้อมูลเดิม
- 5) การผสานข้อมูล (Data Integration) รวมข้อมูลจากหลายแหล่งเพื่อให้ได้ชุดข้อมูลที่สมบูรณ์และครอบคลุมมากยิ่งขึ้น

การดำเนินการเตรียมข้อมูลเหล่านี้จำเป็นต้องใช้ความรอบคอบ ความรู้เชิงธุรกิจ และทักษะด้านสถิติเพื่อให้มั่นใจว่าข้อมูลที่ได้มีคุณภาพเหมาะสมสำหรับการนำไปสร้างแบบจำลองในขั้นตอน Modeling ได้อย่างมีประสิทธิภาพ (Shearer, 2000) นอกจากนี้ การเตรียมข้อมูลที่ดียังช่วยลดความเสี่ยงของปัญหา Overfitting หรือการตีความผลลัพธ์ที่ผิดพลาด

โดยสรุป การเตรียมข้อมูลจึงถือเป็นรากฐานที่สำคัญของความสำเร็จในการวิเคราะห์ข้อมูล และเป็นปัจจัยที่มีผลโดยตรงต่อความน่าเชื่อถือและความแม่นยำของผลลัพธ์ในโครงการวิเคราะห์ข้อมูล (Wirth & Hipp, 2000)

2.3.4 การสร้างแบบจำลอง (Modeling)

เป็นขั้นตอนสำคัญที่ช่วยแปลงข้อมูลที่ได้ผ่านการเตรียมไว้อย่างเหมาะสมให้กลายเป็นแบบจำลองทางคณิตศาสตร์หรือเชิงสถิติ เพื่อใช้ในการทำนายหรือจำแนกข้อมูลตามเป้าหมายของโครงการ

โดยในขั้นตอนนี้ผู้เชี่ยวชาญต้องเลือกเทคนิคหรืออัลกอริทึมการวิเคราะห์ข้อมูลที่เหมาะสมกับลักษณะเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นิยามให้หน้าไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ของปัญหาและรูปแบบของข้อมูล เช่น การจำแนกประเภท (Classification) การถดถอย (Regression) การจัดกลุ่ม (Clustering) หรือเทคนิคเชิงลึก (Deep Learning) (Chapman et al., 2000; Shearer, 2000) การเลือกเทคนิคที่เหมาะสมจำเป็นต้องพิจารณาปัจจัยต่างๆ เช่น ลักษณะของตัวแปรเป้าหมาย (Target Variable) ขนาดและมิติของข้อมูล ความสมดุลของคลาส และข้อจำกัดด้านทรัพยากรของระบบ จากนั้นจะดำเนินการสร้างแบบจำลองโดยใช้ข้อมูลฝึกฝน (Training Data) ที่ได้จากขั้นตอนการเตรียมข้อมูล พร้อมทั้งมีการปรับแต่งค่าพารามิเตอร์ของแบบจำลอง (Hyperparameter Tuning) เพื่อให้ได้ประสิทธิภาพที่ดีที่สุด เช่น การกำหนดค่า Learning Rate จำนวนชั้นของ Neural Network หรือค่าความลึกของ Decision Tree ซึ่งอาจดำเนินการผ่านเทคนิคการค้นหาแบบ Grid Search หรือ Random Search (Wirth & Hipp, 2000)

นอกจากนี้ ในกระบวนการสร้างแบบจำลอง ยังมีการกำหนดวิธีการทดสอบประสิทธิภาพแบบจำลอง เช่น การแบ่งข้อมูลแบบแบ่งชุด (Hold-Out) หรือ Cross-Validation เพื่อประเมินศักยภาพของแบบจำลองในการคาดการณ์ข้อมูลใหม่อย่างถูกต้อง และเพื่อลดความเสี่ยงของปัญหา Overfitting หรือ Underfitting ที่อาจเกิดขึ้น (Chapman et al., 2000)

โดยสรุป ขั้นตอนการสร้างแบบจำลองถือเป็นหัวใจสำคัญของการวิเคราะห์ข้อมูล เนื่องจากผลลัพธ์ที่ได้จากแบบจำลองจะถูกนำไปใช้สนับสนุนการตัดสินใจเชิงธุรกิจในขั้นตอนถัดไป ดังนั้น การเลือกเทคนิคที่เหมาะสม การปรับแต่งพารามิเตอร์อย่างรอบคอบ และการประเมินผลที่ถูกต้องจึงเป็นสิ่งสำคัญในการสร้างความน่าเชื่อถือและประสิทธิภาพสูงสุดให้กับโครงการวิเคราะห์ข้อมูล (Shearer, 2000)

2.3.5 การประเมินผล (Evaluation)

การประเมินผล (Evaluation) เป็นขั้นตอนสำคัญในกระบวนการวิเคราะห์ข้อมูลตามกรอบแนวคิด CRISP-DM โดยมีจุดประสงค์เพื่อประเมินและตรวจสอบประสิทธิภาพของแบบจำลองที่สร้างขึ้นว่ามีความถูกต้องและสอดคล้องกับเป้าหมายทางธุรกิจที่กำหนดไว้หรือไม่ (Chapman et al., 2000) กระบวนการนี้ไม่เพียงแต่นำมาซึ่งการวิเคราะห์ค่าทางสถิติหรือเมตริกต่างๆ เช่น Accuracy, Precision, Recall, F1-Score, ROC-AUC แต่ยังรวมถึงการตีความผลลัพธ์และตรวจสอบว่าผลลัพธ์เหล่านั้นตอบโจทย์ความต้องการทางธุรกิจจริงหรือไม่ (Shearer, 2000) ในขั้นตอนนี้จะมีการนำแบบจำลองที่ผ่านการฝึกฝนไปทดสอบกับข้อมูลที่ไม่ได้ใช้ในการฝึกฝน (Test Set) เพื่อตรวจสอบประสิทธิภาพของแบบจำลองบนข้อมูลใหม่ที่ไม่เคยเห็นมาก่อน เพื่อลดความเสี่ยงของ Overfitting หรือ Underfitting นอกจากนี้ ยังมีการประเมินผลการทำงานของแบบจำลองโดยเปรียบเทียบกับตัวชี้วัดหรือเกณฑ์ความสำเร็จที่ได้กำหนดไว้ตั้งแต่ต้นในขั้นตอน Business Understanding เช่น หากเป้าหมายคือการทำนายลูกค้ากลุ่มเสี่ยงสูง ก็อาจต้องให้ความสำคัญกับค่าความระลึก (Recall) ของกลุ่มเป้าหมายเป็นพิเศษ (Wirth & Hipp, 2000) อีกหนึ่งประเด็นสำคัญ คือ การนำเสนอผลการประเมินควรอยู่ในรูปแบบที่เข้าใจง่าย ชัดเจน และเหมาะสมกับกลุ่มผู้มีส่วนได้ส่วนเสีย ไม่ว่าจะเป็นเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ผู้บริหารหรือผู้ใช้งานทางธุรกิจ เช่น การนำเสนอผ่านกราฟ สรุปรายงาน หรือรายงานเชิงสังเคราะห์ เพื่อสนับสนุนการตัดสินใจในขั้นตอนถัดไป หากผลการประเมินไม่เป็นไปตามเป้าหมายที่กำหนด อาจจำเป็นต้องย้อนกลับไปปรับปรุงในขั้นตอนก่อนหน้า เช่น การเลือกตัวแปรหรือปรับแต่งแบบจำลองเพิ่มเติม (Chapman et al., 2000; Shearer, 2000) โดยสรุป การประเมินผลเป็นขั้นตอนที่มีความสำคัญต่อความน่าเชื่อถือของโครงการวิเคราะห์ข้อมูล และมีบทบาทโดยตรงต่อการนำผลลัพธ์ไปใช้ประโยชน์ในทางธุรกิจอย่างแท้จริง

2.3.6 การนำไปใช้ (Deployment)

เป็นขั้นตอนสุดท้ายในกระบวนการ CRISP-DM ซึ่งมีความสำคัญอย่างยิ่งต่อความสำเร็จของโครงการวิเคราะห์ข้อมูลและเหมืองข้อมูล เนื่องจากเป็นกระบวนการที่นำผลลัพธ์หรือแบบจำลองที่ผ่านการประเมินแล้วไปประยุกต์ใช้ในสภาพแวดล้อมจริงขององค์กร (Chapman et al., 2000) จุดมุ่งหมายหลักของขั้นตอนนี้คือการแปลงองค์ความรู้หรือแบบจำลองที่ได้จากการวิเคราะห์ให้กลายเป็นประโยชน์ที่จับต้องได้ ไม่ว่าจะเป็นในรูปแบบของระบบสนับสนุนการตัดสินใจ (Decision Support Systems) ระบบคัดกรองอัตโนมัติ (Automated Scoring) หรือรายงานเชิงวิเคราะห์ที่ผู้บริหารและผู้มีส่วนได้ส่วนเสียสามารถนำไปใช้ตัดสินใจได้อย่างมีประสิทธิภาพ (Shearer, 2000) ลักษณะการนำไปใช้งานอาจแตกต่างกันตามลักษณะของแต่ละโครงการและความต้องการขององค์กร ตัวอย่างเช่น

- 1) การนำแบบจำลองไปฝังในระบบงานจริง (Production Environment) เช่น ระบบให้คะแนนความเสี่ยงสินเชื่อ ระบบแนะนำผลิตภัณฑ์ หรือระบบตรวจจับธุรกรรมผิดปกติ
- 2) การจัดทำรายงานสรุปผลการวิเคราะห์ในรูปแบบที่เข้าใจง่าย เช่น Dashboard หรือ Interactive Reports
- 3) การสร้างระบบแจ้งเตือนอัตโนมัติ หรือการตั้งค่า Workflow อัตโนมัติเพื่อสนับสนุนการดำเนินงานประจำวัน

นอกจากนี้ ขั้นตอนการนำไปใช้งานยังครอบคลุมถึงการวางแผนการบำรุงรักษา (Monitoring & Maintenance) และการติดตามประสิทธิภาพของแบบจำลองอย่างต่อเนื่อง เพื่อให้สามารถตรวจพบและแก้ไขปัญหาที่อาจเกิดขึ้นในอนาคต เช่น ปัญหา Model Drift หรือข้อมูลที่เปลี่ยนแปลงไปตามเวลา (Wirth & Hipp, 2000) การฝึกอบรมบุคลากรที่เกี่ยวข้องก็เป็นอีกหนึ่งภารกิจสำคัญของขั้นตอนนี้ เพื่อให้มั่นใจว่าองค์กรสามารถใช้ประโยชน์จากผลลัพธ์การวิเคราะห์ได้อย่างยั่งยืน

2.4 แนวคิดและทฤษฎีเกี่ยวกับการแบ่งข้อมูลเพื่อทำการทดสอบและการสุ่มตัวอย่างข้อมูล

การแบ่งข้อมูลและการสุ่มตัวอย่างถือเป็นกระบวนการสำคัญในการพัฒนาแบบจำลองการเรียนรู้ของเครื่อง (Machine Learning) โดยมีวัตถุประสงค์หลักเพื่อประเมินประสิทธิภาพของเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

แบบจำลองในสภาพแวดล้อมที่ใกล้เคียงกับการใช้งานจริงให้ได้มากที่สุด อีกทั้งยังช่วยลดความเสี่ยงจากปัญหา Overfitting ซึ่งเกิดขึ้นเมื่อแบบจำลองจดจำข้อมูลฝึกฝนมากเกินไปจนไม่สามารถประยุกต์ใช้กับข้อมูลใหม่ได้อย่างมีประสิทธิภาพ การเลือกใช้เทคนิคที่เหมาะสมสำหรับการแบ่งข้อมูลและการสุ่มตัวอย่าง เป็นปัจจัยสำคัญที่ส่งผลโดยตรงต่อความน่าเชื่อถือของการประเมินผลแบบจำลอง

2.4.1 การแบ่งข้อมูลเพื่อทำการทดสอบ

การแบ่งข้อมูล (Data Partitioning หรือ Data Splitting) เป็นวิธีการพื้นฐานที่ใช้ในการจัดสรรข้อมูลออกเป็นชุดฝึกฝน (Training Set) และชุดทดสอบ (Testing Set) โดยมีเป้าหมายเพื่อวัดประสิทธิภาพของแบบจำลองบนข้อมูลที่ไม่เคยถูกใช้ในการฝึก เพื่อประเมินความสามารถในการทำนายข้อมูลใหม่ เทคนิคที่ได้รับความนิยมมีดังนี้

1) การแบ่งข้อมูลแบบแยกชุด (Hold-out)

การแบ่งข้อมูลแบบแยกชุด (Hold-out) เป็นวิธีการที่เรียบง่ายและได้รับความนิยมอย่างแพร่หลาย โดยจะทำการสุ่มแบ่งชุดข้อมูลทั้งหมดออกเป็น 2 หรือ 3 ส่วน ได้แก่ ชุดฝึกฝน (Train Set) สำหรับใช้สร้างและฝึกแบบจำลอง ชุดทดสอบ (Test Set) สำหรับประเมินผลการทำนายของแบบจำลองบนข้อมูลที่ไม่เคยใช้ในการฝึก และในบางกรณีอาจมีชุดตรวจสอบ (Validation Set) ด้วย (Kohavi, 1995) อัตราส่วนที่นิยมใช้ได้แก่ 70:30 หรือ 80:20 ตามลำดับ ข้อดีของวิธีนี้ คือมีความรวดเร็วและใช้งานง่าย เหมาะสำหรับชุดข้อมูลขนาดใหญ่ที่มีความหลากหลาย อย่างไรก็ตาม วิธีนี้มีข้อจำกัดในเรื่องความผันผวนของผลลัพธ์ หากชุดข้อมูลมีขนาดเล็กหรือเกิดการสุ่มแบ่งที่ไม่สมดุล อาจทำให้ผลการประเมินแบบจำลองคลาดเคลื่อนจากความเป็นจริงได้

2) การแบ่งข้อมูลแบบไขว้ (K-Fold Cross-Validation)

การแบ่งข้อมูลแบบไขว้ (K-Fold Cross-Validation) เป็นเทคนิคที่มีความยืดหยุ่นและเหมาะสมสำหรับชุดข้อมูลขนาดกลางถึงเล็ก โดยเฉพาะการแบ่งข้อมูลแบบไขว้ (K-Fold Cross-Validation) ซึ่งจะทำการแบ่งข้อมูลทั้งหมดออกเป็น k ส่วนเท่าๆ กัน แล้วสลับนำแต่ละส่วนมาใช้เป็นชุดทดสอบสลับกับชุดฝึกฝนจนครบทุกส่วน จากนั้นจะนำค่าประเมินผลแต่ละรอบมาหาค่าเฉลี่ยเพื่อความน่าเชื่อถือ (Hastie, Tibshirani, & Friedman, 2009) วิธีนี้ช่วยลดอคติ (Bias) และลดความแปรปรวน (Variance) ที่อาจเกิดจากการสุ่มแบ่งข้อมูลเพียงครั้งเดียว

2.4.2 เทคนิคการสุ่มตัวอย่างข้อมูล

การสุ่มตัวอย่างเป็นอีกหนึ่งแนวทางที่ช่วยให้การแบ่งข้อมูลมีความแม่นยำและเหมาะสมมากยิ่งขึ้น โดยเฉพาะในกรณีที่ข้อมูลมีลักษณะไม่สมดุล หรือมีปริมาณจำกัด เทคนิคที่นิยมใช้มีดังนี้

1) การสุ่มแบบแบ่งชั้น (Stratified Sampling)

Stratified Sampling เป็นเทคนิคการสุ่มข้อมูลที่รักษาสัดส่วนของแต่ละกลุ่มเป้าหมาย (Class) ให้ใกล้เคียงกันระหว่างชุดฝึกฝนและชุดทดสอบ เหมาะสมอย่างยิ่งในกรณีที่ข้อมูลมีความไม่

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

สมดุระหว่างคลาส เช่น กรณีศึกษาเรื่องการทำนายลูกค้าที่จะผิดนัดชำระหนี้ ซึ่งกลุ่มผู้ผิดนัดมีจำนวนน้อยกว่ากลุ่มปกติ (Lemaître, Nogueira, & Aridas, 2017) การใช้ Stratified Sampling ช่วยให้แบบจำลองได้รับข้อมูลจากทุกกลุ่มเป้าหมายอย่างเพียงพอทั้งในชั้นฝึกและทดสอบส่งผลให้ผลการประเมินมีความน่าเชื่อถือยิ่งขึ้น

2) การสุ่มแบบบูตสตรอป (Bootstrapping)

Bootstrapping เป็นเทคนิคการสุ่มตัวอย่างซ้ำโดยมีการคืนข้อมูล (Sampling with Replacement) เพื่อสร้างชุดข้อมูลย่อยจำนวนมากสำหรับฝึกและทดสอบแบบจำลอง นำผลที่ได้มาหาค่าเฉลี่ยและประมาณค่าความแปรปรวน (Hastie, Tibshirani, & Friedman, 2009) ข้อดีของ Bootstrapping คือเหมาะกับข้อมูลขนาดเล็กและใช้ประเมินความไม่แน่นอนของแบบจำลองได้ดี อย่างไรก็ตาม อาจซับซ้อนในการตีความผลและใช้เวลาในการประมวลผลสูง

จากแนวคิดและทฤษฎีเกี่ยวกับการแบ่งข้อมูลเพื่อทำการทดสอบและการสุ่มตัวอย่างข้อมูลสรุปได้ว่า การเลือกใช้วิธีการแบ่งข้อมูลที่เหมาะสมในแต่ละงานวิจัยขึ้นอยู่กับขนาดข้อมูล ลักษณะของปัญหา ความสมดุลของกลุ่มเป้าหมาย และข้อจำกัดด้านคำนวณ วิธีการแต่ละแบบล้วนมีจุดแข็งและข้อจำกัดเฉพาะตัว (Kuhn & Johnson, 2013) สำหรับการวิจัยครั้งนี้ ผู้วิจัยได้นำเทคนิคการแบ่งข้อมูลแบบแยกชุด (Hold-out) และการแบ่งข้อมูลแบบไขว้ (K-Fold Cross-Validation) มาเปรียบเทียบเพื่อประเมินประสิทธิภาพของแบบจำลองในการทำนายการผิดนัดชำระหนี้ นอกจากนี้ เพื่อให้การแบ่งข้อมูลมีความสมดุลระหว่างกลุ่มเป้าหมาย ได้ใช้เทคนิคการสุ่มแบบแบ่งชั้น (Stratified Sampling) ร่วมกับกระบวนการแบ่งข้อมูล เพื่อรักษาสัดส่วนของกลุ่มลูกค้าที่มีประวัติผิดนัดและไม่ผิดนัดให้ใกล้เคียงกันทั้งในชุดฝึกและชุดทดสอบ ซึ่งช่วยลดอคติจากข้อมูลไม่สมดุล

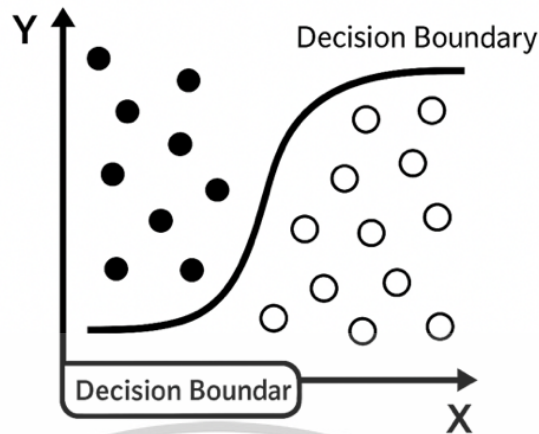
2.5 แนวคิดและทฤษฎีเกี่ยวข้องกับการจำแนกประเภทของข้อมูล (Classification)

การจำแนกประเภทของข้อมูล (Classification) หมายถึง กระบวนการเรียนรู้จากข้อมูลตัวอย่างที่มีป้ายกำกับ เพื่อสร้างแบบจำลองสำหรับทำนายหรือจัดหมวดหมู่ข้อมูลใหม่ให้อยู่ในกลุ่มที่เหมาะสมโดยอัตโนมัติ ไม่ว่าจะเป็นการวิเคราะห์ความเสี่ยงลูกค้า การวินิจฉัยโรค การตรวจจับการฉ้อโกง หรือการจำแนกข้อความภาษา (Kotsiantis, 2007; Hastie, Tibshirani, & Friedman, 2009)

2.5.1 แบบจำลอง Logistic Regression

Logistic Regression เป็นแบบจำลองเชิงสถิติที่นิยมอย่างมากในงานจำแนกประเภท (Classification) โดยเฉพาะเมื่อข้อมูลเป้าหมายเป็นแบบเชิงกลุ่ม (Categorical) เช่น กลุ่มผิดนัด/ไม่ผิดนัดชำระหนี้ จุดแข็งของแบบจำลองนี้คือความเรียบง่าย การตีความได้ตรงไปตรงมา และประสิทธิภาพที่ดีเมื่อตัวแปรต้นมีความสัมพันธ์เชิงเส้นกับ Logit ของความน่าจะเป็น (Han et al., 2011; Hosmer et al., 2013) Logistic Regression จึงกลายเป็นแบบจำลองมาตรฐานที่นักวิจัยและผู้พัฒนาด้านข้อมูลเลือกใช้เป็น Baseline ในการเปรียบเทียบกับเทคนิค Machine Learning อื่นๆ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 2.2 ตัวอย่างการแยกกลุ่มข้อมูลสองกลุ่มด้วย Logistic Regression

แบบจำลอง Logistic Regression พัฒนาค่อยๆมาจากแนวคิดของ Linear Regression ด้วยการนำผลรวมเชิงเส้นของตัวแปรต้น (Feature) และสัมประสิทธิ์ (Coefficient) มาแปลงค่าด้วย ฟังก์ชันซิกมอยด์ (Sigmoid Function) เพื่อให้ผลลัพธ์อยู่ในช่วงค่าความน่าจะเป็นระหว่าง 0 ถึง 1 เหมาะกับการประเมินโอกาสที่ตัวอย่างข้อมูลหนึ่งจะอยู่ในกลุ่มเป้าหมายหรือไม่ โครงสร้างสมการของ Logistic Regression สามารถเขียนได้ดังนี้

$$\hat{Y}_i = \frac{1}{1 + e^{-z}} \quad \text{เมื่อ } z = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in} \quad (2.1)$$

หรือแปลงให้อยู่ในรูป Logit ได้ดังนี้

$$\ln\left(\frac{\hat{Y}}{1 - \hat{Y}}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n \quad (2.2)$$

ในกระบวนการประมาณค่าสัมประสิทธิ์ (β) จะใช้วิธี Maximum Likelihood Estimation (MLE) ซึ่งเป็นเทคนิคมาตรฐานที่นิยมในงานสถิติและ Machine Learning (Hosmer, Lemeshow & Sturdivant, 2013) ประเภทของ Logistic Regression สามารถแบ่งตามจำนวนกลุ่มของตัวแปรเป้าหมายได้แก่

- 1) การถดถอยโลจิสติกทวิภาค (Binary Logistic Regression) คือตัวแปรตามมีค่าเพียง 2 ค่า เช่น ผิด/ไม่ผิด
- 2) การถดถอยโลจิสติกพหุ (Multinomial Logistic Regression) คือตัวแปรตามมีค่า

มากกว่า 2 ค่าขึ้นไป เช่น การจัดประเภทสินค้าหรือผลสอบหลายระดับ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อผู้เช่าได้เห็นไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

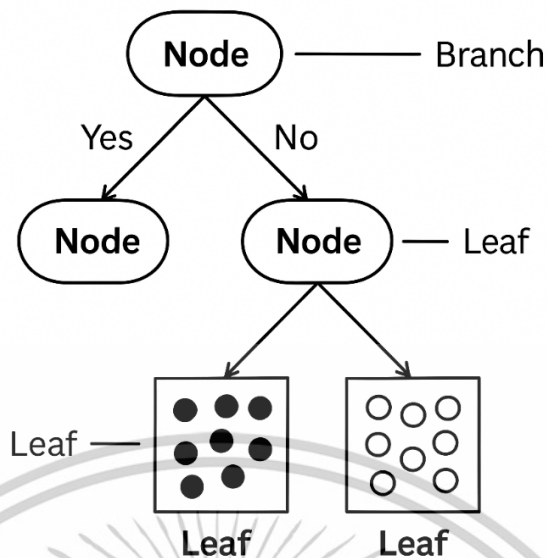
Logistic Regression ได้รับความนิยมอย่างต่อเนื่องเนื่องจากแบบจำลองนี้มีข้อดีในแง่ของความโปร่งใสและความเข้าใจง่าย ผลลัพธ์ที่ได้จากแบบจำลองสามารถตีความเชิงสาเหตุได้โดยตรงผ่านค่าสัมประสิทธิ์ ซึ่งบ่งชี้ถึงทิศทางและขนาดของอิทธิพลของแต่ละตัวแปรต้นต่อโอกาสที่ตัวอย่างจะอยู่ในกลุ่มเป้าหมาย นอกจากนี้ Logistic Regression ยังใช้เวลาฝึกฝนและทดสอบที่รวดเร็วมกเมื่อเทียบกับแบบจำลองที่มีความซับซ้อนสูงกว่า ทั้งยังรองรับการปรับแต่งแบบจำลอง เช่น การเพิ่ม Regularization (L1, L2) เพื่อควบคุมการเกิด Overfitting ในชุดข้อมูลที่มีตัวแปรจำนวนมาก (Ng, 2004) อีกทั้งยังสามารถใช้เป็น Baseline Model สำหรับเปรียบเทียบกับเทคนิค Machine Learning อื่นๆ ได้อย่างเหมาะสมในงานวิจัย

อย่างไรก็ตาม Logistic Regression ก็มีข้อจำกัดสำคัญที่ต้องตระหนัก โดยแบบจำลองนี้ตั้งอยู่บนสมมติฐานว่าความสัมพันธ์ระหว่างตัวแปรต้นกับ Logit ของค่าความน่าจะเป็นจะต้องเป็นเชิงเส้น หากความสัมพันธ์เป็นแบบไม่เชิงเส้นหรือมีความซับซ้อนสูง เช่น มีปฏิสัมพันธ์ระหว่างตัวแปรหรือความสัมพันธ์เชิงโค้ง แบบจำลองจะทำนายได้ไม่ดีเท่ากับเทคนิคที่ยืดหยุ่นกว่า เช่น Tree-based หรือ Neural Network นอกจากนี้ Logistic Regression อาจมีปัญหาเมื่อเจอ Multicollinearity ระหว่างตัวแปรต้น หรือมี Outlier จำนวนมาก ซึ่งจะส่งผลกระทบต่อค่าพารามิเตอร์ให้ผิดเพี้ยนไป ในกรณีที่ข้อมูลกลุ่มเป้าหมายมีความไม่สมดุลสูง (Imbalanced Data) แบบจำลองนี้อาจทำนายกลุ่มที่มีจำนวนน้อยได้ไม่ดีนัก เว้นแต่จะปรับสมดุลข้อมูลก่อนฝึกแบบจำลอง (Han et al., 2011) แม้จะมีข้อจำกัดดังกล่าว Logistic Regression ก็ยังคงเป็นหนึ่งในแบบจำลองที่ได้รับความนิยมสูงสุดในงานวิเคราะห์ข้อมูลที่ต้องการความแม่นยำ ความรวดเร็ว และความสามารถในการตีความผลลัพธ์แบบโปร่งใส เหมาะอย่างยิ่งสำหรับการประเมินความเสี่ยง วิเคราะห์พฤติกรรมลูกค้า และเป็นจุดเริ่มต้นของงานวิจัยเชิงเปรียบเทียบในสาขาต่างๆ (Jon, 2007; Han et al., 2011; Hosmer et al., 2013)

2.5.2 แบบจำลอง Decision Tree

Decision Tree เป็นแบบจำลองเชิงตรรกะที่ได้รับความนิยมอย่างสูงในงานวิเคราะห์ข้อมูลทั้งในแวดวงวิชาการและอุตสาหกรรม โดยเฉพาะในโจทย์ที่ต้องการการจำแนกประเภท (Classification) หรือการทำนายค่าตัวเลข (Regression) แบบจำลองนี้ถูกออกแบบมาให้มีโครงสร้างลำดับชั้นเหมือน “ต้นไม้” ประกอบด้วยโหนดภายใน (Internal Node) สำหรับทดสอบเงื่อนไขของข้อมูล โหนดกิ่ง (Branch) ที่แยกข้อมูลตามผลลัพธ์ และโหนดปลายทาง (Leaf Node) ที่ระบุผลลัพธ์สุดท้ายหรือกลุ่มของข้อมูล (Han, Kamber & Pei, 2011; Safavian & Landgrebe, 1991) ความโดดเด่นของ Decision Tree คือความสามารถในการแปลงเงื่อนไขเชิงตรรกะให้เป็นรูปแบบกราฟิกที่เข้าใจง่าย ผู้ใช้งานสามารถติดตามกระบวนการตัดสินใจของแบบจำลองได้ตลอดเส้นทางจากรากสู่ใบ ทำให้เหมาะสำหรับการอธิบายและนำเสนอผลการวิเคราะห์กับผู้ที่ไม่มีพื้นฐานด้านข้อมูลโดยตรง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 2.3 แสดงการแบ่งข้อมูลแต่ละโหนดและการจำแนกข้อมูลที่ใบของต้นไม้

ในทางทฤษฎี Decision Tree จะทำงานโดยการเลือกตัวแปรและเงื่อนไขที่เหมาะสมที่สุด เพื่อแยกข้อมูลแต่ละกลุ่มให้มีความบริสุทธิ์ (Purity) มากที่สุดในแต่ละระดับของต้นไม้ โดยใช้เกณฑ์ Information Gain, Gini Index หรือ Chi-Square ในการประเมินการแบ่งกลุ่มที่ดีที่สุด (Breiman et al., 1984) ตัวอย่างเช่น หากต้องการทำนายว่าลูกค้าจะผัดนวดชำระหนี้หรือไม่ แบบจำลองจะเลือกตัวแปรที่แบ่งข้อมูลออกเป็นกลุ่ม “เสี่ยง” กับ “ไม่เสี่ยง” ได้ชัดเจนที่สุดเป็นอันดับแรก แล้วจึงแบ่งย่อยลงไป ตามตัวแปรที่มีอิทธิพลรองลงมา โครงสร้างนี้ทำให้แบบจำลองมีความยืดหยุ่น สามารถรองรับข้อมูลที่มีความสัมพันธ์เชิงซ้อนหรือไม่เป็นเชิงเส้น (Non-Linear) ได้ดี และไม่ต้องการการปรับแต่งข้อมูลต้นฉบับ (Data Preprocessing) มากนัก

ประเภทของ Decision Tree สามารถจำแนกได้เป็น

- 1) Classification Tree (ใช้กับตัวแปรเป้าหมายเชิงกลุ่ม)
- 2) Regression Tree (ใช้กับตัวแปรเป้าหมายเชิงปริมาณ)

Decision Tree ได้รับความนิยมสูงในงานวิเคราะห์ข้อมูลเนื่องจากมีข้อดีหลายประการ เช่น ความโปร่งใสและความเข้าใจง่าย ผลลัพธ์ที่ได้สามารถแปลความเป็นกฎหรือเงื่อนไขเชิงตรรกะได้โดยตรง ทำให้สนับสนุนกระบวนการตัดสินใจเชิงนโยบายในภาคธุรกิจหรือหน่วยงานรัฐ นอกจากนี้แบบจำลองนี้ยังสามารถจัดการกับข้อมูลที่ขาดหาย (Missing Value) หรือข้อมูลเชิงหมวดหมู่และตัวเลขได้ในเวลาเดียวกัน โดยไม่ต้องแปลงค่าล่วงหน้า ทั้งยังรองรับการทำ Feature Selection ภายในตัวผ่านกระบวนการเลือกเงื่อนไขที่เหมาะสมที่สุดในแต่ละโหนด และสามารถจัดการกับข้อมูลที่มีความสัมพันธ์ไม่เชิงเส้นได้ดีกว่าแบบจำลองเชิงเส้นอย่าง Logistic Regression (Han et al., 2011; Quinlan, 1986) อย่างไรก็ตาม Decision Tree ก็มีข้อจำกัดที่ควรตระหนัก โดยธรรมชาติของแบบจำลองที่มีความยืดหยุ่นสูงอาจทำให้เกิดปัญหา Overfitting ได้ง่ายในกรณีที่ต้นไม้มีความลึกมาก เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นิยามให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

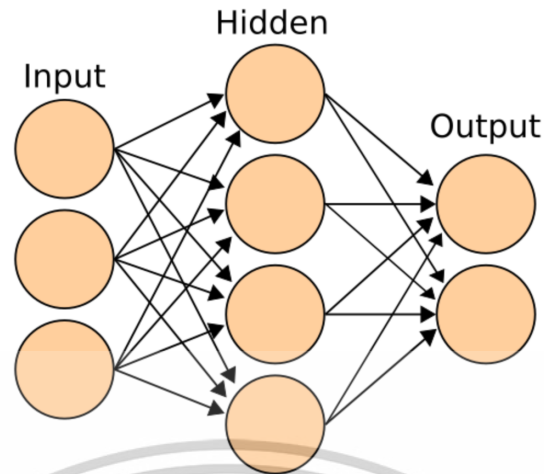
เกินไป หรือตัดสินใจตาม noise ในข้อมูลมากเกินไป นอกจากนี้ผลลัพธ์ที่ได้อาจเปลี่ยนแปลงไปมาก หากข้อมูลฝึกฝนไม่สมดุลหรือมีการสุ่มตัวอย่างที่ต่างกัน แบบจำลองยังไวต่อ Outlier และอาจมีความไม่เสถียร (Variance สูง) เมื่อเทียบกับเทคนิค Ensemble ที่พัฒนาต่อยอด เช่น Random Forest (Breiman, 2001) การควบคุมความลึก (Pruning) หรือการกำหนดจำนวนตัวอย่างขั้นต่ำที่ต้องมีในแต่ละใบจึงเป็นสิ่งสำคัญในการลดความเสี่ยงของ Overfitting

โดยสรุป Decision Tree คือแบบจำลองที่เหมาะสมอย่างยิ่งกับปัญหาที่ต้องการทั้งความแม่นยำ และการตีความผลลัพธ์ได้ง่าย ช่วยสนับสนุนการตัดสินใจเชิงเหตุผลและสามารถนำไปประยุกต์ใช้กับข้อมูลหลากหลายรูปแบบ แม้จะมีข้อจำกัดด้านความเสถียรและปัญหา Overfitting แต่หากได้รับการควบคุมที่เหมาะสม แบบจำลองนี้ยังคงเป็นหนึ่งในเครื่องมือที่ทรงพลังและได้รับความนิยมในสาย Data Science มาจนถึงปัจจุบัน (Han et al., 2011; Breiman, 2001; Quinlan, 1986)

2.5.3 แบบจำลอง Neural Networks

แบบจำลองโครงข่ายประสาทเทียม (Neural Networks) เป็นเทคนิคหนึ่งในกลุ่ม Machine Learning ที่ได้รับความนิยมและมีการพัฒนาอย่างต่อเนื่องตั้งแต่กลางศตวรรษที่ 20 โดยได้รับแรงบันดาลใจจากกลไกการทำงานของสมองมนุษย์ ซึ่งมีหน่วยประมวลผลพื้นฐานเรียกว่า “เซลล์ประสาท” (Neuron) แนวคิดหลักของ Neural Networks คือการเชื่อมโยงหน่วยประมวลผลจำนวนมากเข้าด้วยกันเป็นโครงข่ายหลายชั้น (Multilayer Network) เพื่อถ่ายทอดข้อมูลและเรียนรู้รูปแบบ (Pattern) ที่ซับซ้อนในข้อมูล (Goodfellow, Bengio & Courville, 2016) โครงสร้างพื้นฐานประกอบด้วยชั้นอินพุต (Input Layer) สำหรับรับข้อมูลต้นทาง, ชั้นซ่อน (Hidden Layer) ที่ทำหน้าที่แปลงข้อมูลด้วยฟังก์ชันคณิตศาสตร์แบบไม่เชิงเส้น (Non-linear Activation Function) และชั้นเอาต์พุต (Output Layer) สำหรับให้ผลลัพธ์สุดท้ายของแบบจำลอง

Neural Networks มีจุดเด่นที่สำคัญคือความสามารถในการจำลองความสัมพันธ์ที่ซับซ้อนสูงในข้อมูล สามารถรองรับพีเจอร์จำนวนมากและจับความสัมพันธ์แบบไม่เชิงเส้นได้ดีมากกว่าแบบจำลองคลาสสิกอย่าง Logistic Regression หรือ Decision Tree ในแต่ละโหนดของโครงข่ายจะทำการคำนวณแบบถ่วงน้ำหนัก (Weighted Sum) และนำผ่านฟังก์ชันกระตุ้น (Activation Function) เช่น Sigmoid, ReLU หรือ Tanh เพื่อนำผลลัพธ์ไปยังโหนดชั้นถัดไป กระบวนการเรียนรู้ใน Neural Networks จะอาศัยวิธี Backpropagation ซึ่งเป็นการปรับน้ำหนักภายในโครงข่ายด้วยอัลกอริทึม Gradient Descent เพื่อลดค่าความคลาดเคลื่อน (Loss Function) ระหว่างผลทำนายกับค่าจริง (Rumelhart, Hinton & Williams, 1986)



รูปที่ 2.4 โครงสร้างโครงข่ายประสาทเทียมแบบง่าย

ที่มา : Wikimedia Commons (2006)

ปัจจุบัน Neural Networks แบ่งออกเป็นหลายประเภทตามวัตถุประสงค์การใช้งาน เช่น Feed Forward Neural Network เป็นโครงสร้างพื้นฐานสำหรับการจำแนกประเภท (Classification) หรือการทำนายค่า (Regression), Convolutional Neural Network (CNN) ที่โดดเด่นด้านการวิเคราะห์ภาพและสัญญาณ, และ Recurrent Neural Network (RNN) ซึ่งเหมาะกับข้อมูลลำดับเวลา (Time Series) หรือข้อมูลที่มีลำดับเหตุการณ์ ทั้งนี้ การประยุกต์ใช้งาน Neural Networks มีความหลากหลายมาก ตั้งแต่การจำแนกภาพ เสียง ข้อความ การพยากรณ์พฤติกรรมผู้บริโภค ไปจนถึงการเรียนรู้เชิงลึก (Deep Learning) ในระบบที่มีข้อมูลจำนวนมาก

จุดแข็งของ Neural Networks อยู่ที่ความสามารถในการเรียนรู้จากข้อมูลขนาดใหญ่และความซับซ้อนสูง สามารถประยุกต์ใช้กับโจทย์ที่ต้องการการจับความสัมพันธ์แบบไม่เชิงเส้นหรือข้อมูลหลายมิติได้อย่างมีประสิทธิภาพ แบบจำลองสามารถปรับแต่งโครงสร้างและพารามิเตอร์ให้เหมาะสมกับงานแต่ละประเภท เช่น การเพิ่มจำนวนชั้นซ่อน (Deep Neural Networks) หรือการปรับเปลี่ยนฟังก์ชันกระตุ้น ทั้งนี้ Neural Networks ยังสามารถเรียนรู้คุณลักษณะ (Feature Learning) จากข้อมูลโดยตรงโดยไม่ต้องพึ่งพาการสกัดพีเจอร์แบบเดิมเสมอไป ส่งผลให้มีประสิทธิภาพสูงในงานที่มีความซับซ้อน เช่น การรู้จำใบหน้า การแปลภาษา หรือการวินิจฉัยโรคจากภาพถ่ายทางการแพทย์ (LeCun, Bengio & Hinton, 2015) อย่างไรก็ตาม Neural Networks ก็มีข้อจำกัดที่สำคัญ โดยเฉพาะอย่างยิ่งในแง่ของความโปร่งใส (Interpretability) เนื่องจากโครงสร้างที่ซับซ้อนมากอาจทำให้ผลลัพธ์ที่ได้ยากต่อการตีความว่าเกิดจากอิทธิพลของตัวแปรใด (Black-Box Model) การฝึกแบบจำลองยังต้องการข้อมูลขนาดใหญ่และใช้ทรัพยากรคอมพิวเตอร์สูง ทั้งในแง่ของเวลาและพลังงาน นอกจากนี้แบบจำลองมีความเสี่ยงต่อ Overfitting หากไม่มีการควบคุมที่เหมาะสม เช่น การหยุดการเรียนรู้ก่อน (Early Stopping) Regularization หรือการใช้เทคนิค Dropout อีกทั้งการเลือกพารามิเตอร์และโครงสร้างที่เหมาะสมต้องอาศัยประสบการณ์และการทดลองหลายครั้งจึงจะได้

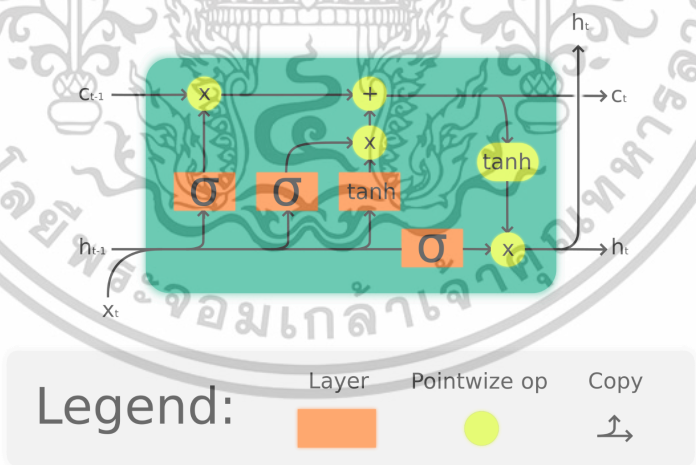
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปเผยแพร่ในสื่ออื่นโดยไม่ได้รับอนุญาต

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ผลลัพธ์ที่ดีที่สุด (Goodfellow et al., 2016) Neural Networks ถือเป็นรากฐานของเทคนิค Deep Learning ตั้งแต่งานวิเคราะห์ข้อมูลพื้นฐานจนถึงปัญญาประดิษฐ์ขั้นสูง อย่างไรก็ตามการเลือกใช้ Neural Networks ควรพิจารณาตามความเหมาะสมของปัญหาและข้อจำกัดของข้อมูลและทรัพยากร

2.5.4 แบบจำลอง Long Short-Term Memory (LSTM)

แบบจำลอง Long Short-Term Memory หรือ LSTM เป็นหนึ่งในโครงข่ายประสาทเทียมแบบลำดับ (Recurrent Neural Network: RNN) ที่ถูกออกแบบมาเพื่อแก้ไขข้อจำกัดของ RNN เดิมในการเรียนรู้ลำดับข้อมูลที่มีระยะเวลายาวหรือมีบริบทข้ามช่วงเวลา (Long-Term Dependency) LSTM ได้รับการพัฒนาขึ้นโดย Hochreiter และ Schmidhuber ในปี 1997 โดยมีกลไกสำคัญคือการสร้าง “หน่วยความจำ” (Memory Cell) ที่สามารถเก็บรักษาข้อมูลสำคัญไว้ได้เป็นเวลานาน พร้อมทั้งใช้ “ประตู” (Gates) หลายประเภทในการควบคุมการไหลของข้อมูลเข้าหรือออกจากหน่วยความจำ (Hochreiter & Schmidhuber, 1997) LSTM ประกอบด้วย Input Gate, Forget Gate และ Output Gate ซึ่งแต่ละประตูจะมีบทบาทในการตัดสินใจว่าข้อมูลส่วนใดควรนำเข้าไปเก็บไว้ หรือส่งต่อไปยังช่วงเวลาถัดไป Input Gate ควบคุมข้อมูลใหม่ที่เข้าสู่เซลล์หน่วยความจำ ส่วน Forget Gate ตัดสินใจว่าข้อมูลใดในหน่วยความจำควรถูกลบออก และ Output Gate ควบคุมว่าข้อมูลใดจะส่งออกไปเป็นผลลัพธ์ของแบบจำลอง LSTM ทำให้สามารถเรียนรู้ข้อมูลลำดับที่มีโครงสร้างซับซ้อนได้อย่างมีประสิทธิภาพ เช่น ลำดับข้อความ ข้อมูลธุรกรรมต่อเนื่อง หรือชุดข้อมูลอนุกรมเวลา (Time Series) ซึ่งแบบจำลองอื่นอาจขาดศักยภาพในการจดจำบริบทระยะยาว



รูปที่ 2.5 โครงสร้างเซลล์ LSTM

ที่มา : Wikimedia Commons (2018)

จากรูปที่ 2.5 แสดงให้เห็นว่าข้อมูลในแต่ละช่วงเวลาจะถูกควบคุมโดยกลไก Gate ต่างๆ ของ LSTM เซลล์ เพื่อกำหนดว่าข้อมูลเดิมควรถูกเก็บหรือลบทิ้ง ข้อมูลใหม่จะเข้าเซลล์หรือไม่ และค่าผลลัพธ์ใดจะถูกส่งต่อไปยังขั้นถัดไป ส่งผลให้ LSTM มีความสามารถในการจดจำข้อมูลที่สำคัญในเอกสารนี้เป็นเอกสารที่ส่งวนเวียนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนุญเตเห็นาเบเซบประะเขชนดานการคว่าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ลำดับเวลาระยะยาวได้ดีกว่า RNN ทั่วไป โดย เส้นลูกศร แสดงทิศทางการไหลของข้อมูลและสถานะต่างๆ ภายในเซลล์ c_{t-1} และ h_{t-1} คือค่าหน่วยความจำและสถานะที่ส่งต่อมาจากช่วงเวลาก่อนหน้า x_t คือข้อมูลอินพุตในช่วงเวลาปัจจุบัน โดยภายในเซลล์จะประกอบด้วยฟังก์ชันหลัก 3 ประเภท ได้แก่

- 1) Sigmoid Layer (σ) ทำหน้าที่เป็น "Gate" หรือประตูควบคุมข้อมูล ได้แก่
 - Forget Gate : ตัดสินใจว่าข้อมูลส่วนใดในสถานะเดิมควรลบทิ้ง
 - Input Gate : กำหนดว่าข้อมูลใหม่ส่วนใดควรเพิ่มเข้าสู่หน่วยความจำ
 - Output Gate : กำหนดว่าข้อมูลส่วนใดควรถูกส่งออกจากเซลล์ไปเป็น Output
- 2) Tanh Layer ใช้สำหรับสร้างเวกเตอร์ใหม่ของข้อมูลที่อาจถูกเพิ่มเข้าสู่ Cell State
- 3) Pointwise Operations (วงกลมสี่เหลี่ยม) : สัญลักษณ์ มีความหมายดังนี้
 - สัญลักษณ์ \times คือ การคูณองค์ประกอบ (Element-Wise Multiplication) เช่น การควบคุมข้อมูลด้วยค่าจากประตู (Gate)
 - สัญลักษณ์ $+$ คือ การบวกองค์ประกอบ (Element-Wise Addition)
 - Copy (สัญลักษณ์ลูกศรโค้ง) คือ การคัดลอกข้อมูล

โดยผลรวมของกระบวนการเหล่านี้คือค่าหน่วยความจำใหม่ (c_t) และสถานะเอาต์พุตใหม่ (h_t) ที่จะถูกส่งไปยังช่วงเวลาถัดไป โครงสร้างของ LSTM ยังเปิดโอกาสให้สามารถเรียนรู้ทั้งข้อมูลระยะสั้น (Short-Term) และระยะยาว (Long-Term) ไปพร้อมกัน โดยการปรับแตงน้ำหนักของแต่ละประตูตามลักษณะข้อมูลและบริบทที่จำเป็น สมการสำคัญของ LSTM ประกอบด้วยการคำนวณค่าฟังก์ชัน Sigmoid และ Tanh สำหรับประตูแต่ละประเภท และมีการอัปเดตค่า Cell State และ Hidden State ในแต่ละช่วงเวลา อัลกอริทึม Backpropagation Through Time (BPTT) ถูกนำมาใช้เพื่อปรับค่าพารามิเตอร์ให้เหมาะสมระหว่างการฝึกฝน (Graves, 2012)

LSTM ได้รับความนิยมอย่างสูงในงานวิเคราะห์ข้อมูลลำดับ เช่น การทำนายข้อมูลอนุกรมเวลา (Time Series Forecasting) การรู้จำเสียงพูด (Speech Recognition) การวิเคราะห์ข้อความ (Text Analysis) และการพยากรณ์เหตุการณ์ในข้อมูลธุรกิจ เพราะสามารถจดจำลำดับหรือ Pattern ที่ยาวและซับซ้อนได้ดีมากกว่า RNN หรือแบบจำลองเชิงเส้นทั่วไป จุดเด่นของ LSTM คือความสามารถในการลดปัญหา Vanishing Gradient ที่พบใน RNN ดั้งเดิม และช่วยให้แบบจำลองสามารถเก็บรักษาข้อมูลสำคัญได้อย่างต่อเนื่องในช่วงเวลาที่ยาวนาน

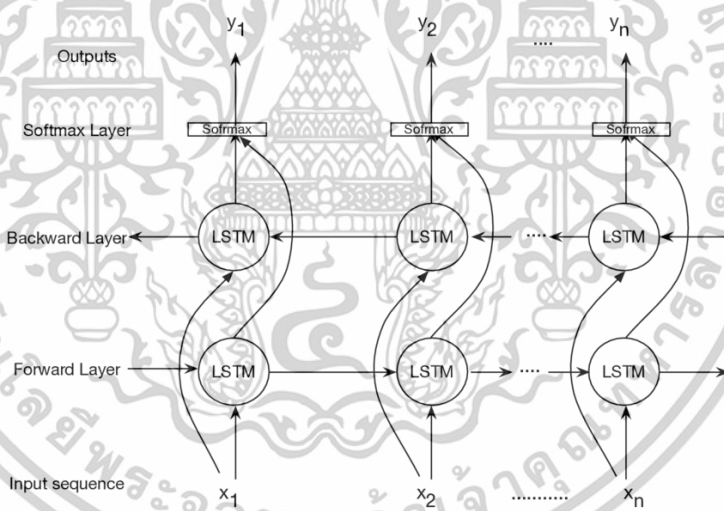
อย่างไรก็ตาม LSTM ก็มีข้อจำกัดในแง่ของการใช้ทรัพยากรคำนวณที่สูงกว่า RNN แบบธรรมดา เนื่องจากโครงสร้างที่ซับซ้อนมากขึ้น นอกจากนี้การฝึกฝนแบบจำลองอาจใช้เวลานาน และต้องมีการปรับแต่งพารามิเตอร์จำนวนมาก เช่น จำนวนหน่วยในแต่ละชั้น (Units) จำนวนชั้น (Layers) หรือ Learning Rate เพื่อให้ได้ประสิทธิภาพสูงสุด อีกทั้งในกรณีที่ข้อมูลไม่ได้มีลำดับหรือ Pattern ที่ยาวมาก แบบจำลองแบบง่ายอย่าง FeedForward Neural Network หรือ RNN แบบดั้งเดิมก็อาจเพียงพอสำหรับการประยุกต์ใช้งาน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.5.5 แบบจำลอง Bidirectional Long Short-Term Memory (BiLSTM)

Bidirectional Long Short-Term Memory หรือ Bidirectional LSTM (BiLSTM) เป็นการขยายขีดความสามารถของโครงข่ายประสาทเทียมแบบ LSTM โดยออกแบบให้สามารถเรียนรู้ข้อมูลลำดับได้ทั้งสองทิศทางพร้อมกัน กล่าวคือ แบบจำลอง BiLSTM จะประมวลผลข้อมูลทั้งจากอดีตสู่ออนาคต (Forward) และจากอนาคตย้อนกลับสู่ออดีต (Backward) ในแต่ละช่วงเวลา ทำให้สามารถดึงข้อมูลบริบท (Context) ที่อยู่ก่อนและหลังตำแหน่งเป้าหมายมาใช้ประกอบการทำนายได้อย่างสมบูรณ์มากยิ่งขึ้น (Schuster & Paliwal, 1997)

หลักการการทำงานของ BiLSTM จะประกอบด้วย LSTM สองชุดซึ่งทำงานขนานกัน โดยชุดหนึ่งรับข้อมูลเรียงลำดับจากซ้ายไปขวา (Forward LSTM) ส่วนอีกชุดรับข้อมูลจากขวาไปซ้าย (Backward LSTM) ผลลัพธ์จากทั้งสองทิศทางจะถูกนำมารวมกันเพื่อสร้างการแทนค่าที่มีข้อมูลบริบทครบถ้วนสำหรับแต่ละจุดในลำดับข้อมูล (Graves & Schmidhuber, 2005) ข้อดีของแนวคิดนี้คือในงานที่ต้องวิเคราะห์ข้อมูลลำดับ เช่น การประมวลผลภาษา การรู้จำเสียง หรือการทำนายอนุกรมเวลา BiLSTM จะสามารถเข้าใจทั้งบริบทก่อนหน้าและหลังจากตำแหน่งปัจจุบัน ส่งผลให้แบบจำลองตัดสินใจได้แม่นยำขึ้น



รูปที่ 2.6 โครงสร้างของแบบจำลอง Bidirectional LSTM

ที่มา : Tavakoli, N. (2020)

จากรูป 2.6 แสดงโครงสร้างของแบบจำลอง Bidirectional LSTM (BiLSTM) ซึ่งประกอบด้วยสองเลเยอร์หลักคือ

Forward Layer: รับข้อมูลจากซ้ายไปขวาเพื่อจับความสัมพันธ์เชิงเวลาแบบลำดับ เช่น จาก $x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_n$

Backward Layer: รับข้อมูลย้อนกลับจากขวาไปซ้ายเพื่อจับความสัมพันธ์จากอนาคตย้อนหลัง เช่น จาก $x_n \rightarrow x_{n-1} \rightarrow \dots \rightarrow x_1$

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ทั้งสองเลเยอร์ใช้ LSTM Units และประมวลผลข้อมูลแบบขนานกัน แล้วผลลัพธ์จากทั้งสองทิศทางจะถูกรวมเข้าด้วยกัน (Concatenated) ที่ Softmax Layer ซึ่งให้ Output สุดท้ายเป็นลำดับ y_1, y_2, \dots, y_n

ตัวอย่างงานที่ BiLSTM ได้รับความนิยมสูง ได้แก่ การวิเคราะห์ลำดับข้อความ (Sequence Labeling) การทำนายพฤติกรรมลูกค้า การวินิจฉัยทางการแพทย์จากข้อมูลลำดับ และการประเมินลำดับเหตุการณ์ในระบบทางธุรกิจ ทั้งนี้ โครงสร้าง BiLSTM ยังคงใช้กลไกเดียวกับ LSTM ประกอบด้วยหน่วยความจำ (Cell State) และประตูควบคุม (Input, Forget, Output Gate) โดยประมวลผลทั้ง Forward และ Backward ในแต่ละ Time Step (Graves et al., 2013)

ข้อดีของ BiLSTM คือความสามารถในการจับ Pattern และบริบททั้งในอดีตและอนาคตในข้อมูลลำดับเดียวกัน ทำให้เหมาะกับปัญหาที่ต้องใช้ข้อมูลบริบทสองด้านเพื่อการทำนายที่แม่นยำ เช่น การทำนายคำในประโยคภาษา หรือการประเมินความเสี่ยงจากอนุกรมข้อมูลทางการเงิน อย่างไรก็ตาม BiLSTM มีข้อจำกัดในแง่ของการใช้ทรัพยากรคำนวณที่สูงขึ้นกว่าระบบ LSTM ปกติ เพราะต้องประมวลผลข้อมูลสองทิศทางพร้อมกัน และไม่เหมาะกับปัญหาที่ต้องการพยากรณ์แบบออนไลน์ (Online Prediction) ซึ่งรู้ข้อมูลเพียงทิศทางเดียว การพัฒนา BiLSTM เป็นรากฐานสำคัญของงานวิจัยด้านประมวลผลภาษา (NLP) การรู้จำเสียงพูด (Speech Recognition) และการทำนายข้อมูลอนุกรมเวลา (Time Series Forecasting) ที่ต้องการเข้าใจข้อมูลทั้งบริบทก่อนหน้าและภายหลังตำแหน่งเป้าหมาย ซึ่งช่วยเพิ่มความแม่นยำให้กับงานเหล่านี้อย่างมีนัยสำคัญ (Schuster & Paliwal, 1997; Graves et al., 2013)

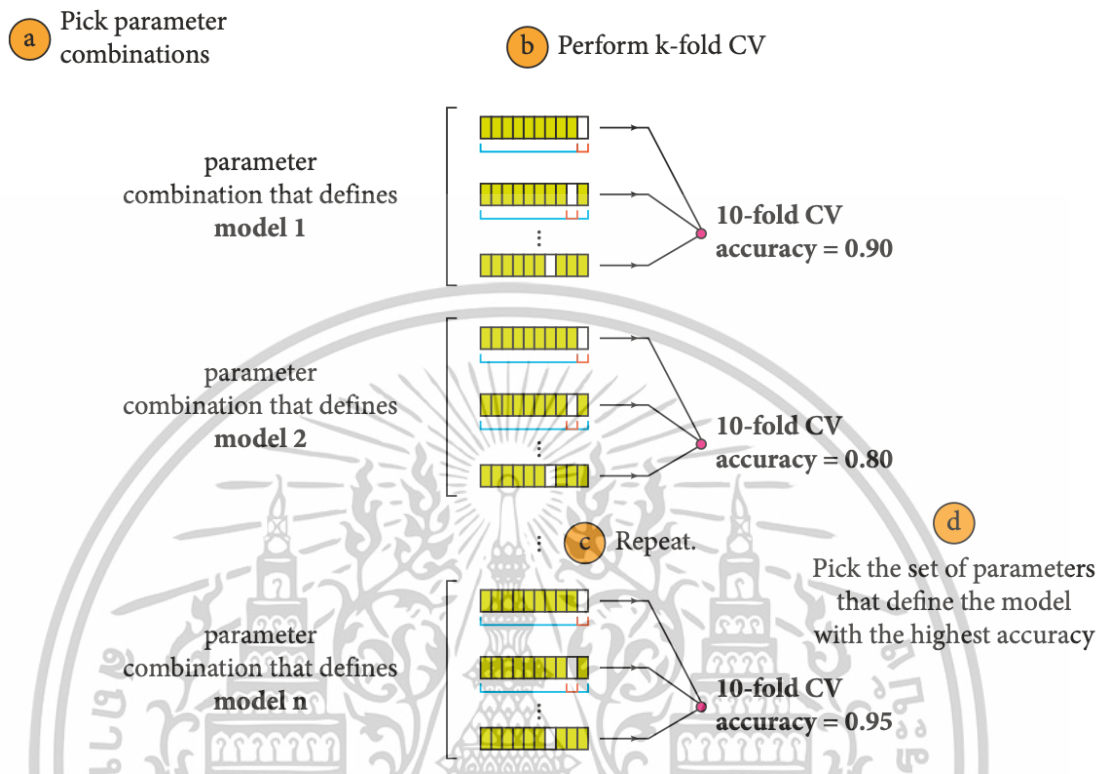
2.6 แนวคิดและทฤษฎีเกี่ยวข้องกับการเลือกไฮเปอร์พารามิเตอร์แบบกริด (Grid Search Cross-Validation)

การเลือกไฮเปอร์พารามิเตอร์ (Hyperparameter Tuning) เป็นขั้นตอนสำคัญในกระบวนการพัฒนา Machine Learning ที่มีผลต่อประสิทธิภาพของแบบจำลองอย่างมีนัยสำคัญ ไฮเปอร์พารามิเตอร์คือพารามิเตอร์ที่ไม่ถูกปรับค่าผ่านกระบวนการเรียนรู้โดยตรง แต่ต้องกำหนดไว้ล่วงหน้า เช่น ค่าความลึกของต้นไม้ (Max_Depth) ใน Decision Tree จำนวนโหนดในชั้นซ่อนของ Neural Network หรือค่า Learning Rate ในการฝึกแบบจำลอง ตัวเลือกที่เหมาะสมของไฮเปอร์พารามิเตอร์จะช่วยเพิ่มความแม่นยำ ลดความซับซ้อนเกินจำเป็น (Overfitting) และทำให้แบบจำลองทั่วไปกับข้อมูลใหม่ได้ดีขึ้น (Bergstra & Bengio, 2012)

Grid Search Cross-Validation เป็นวิธีมาตรฐานและได้รับความนิยมสูงในการเลือกไฮเปอร์พารามิเตอร์ โดยแนวคิดหลักคือการกำหนดชุดค่าของไฮเปอร์พารามิเตอร์ที่ต้องการทดลอง จากนั้นสร้าง "กริด" หรือทุกความเป็นไปได้ของการผสมค่าต่างๆ แล้วนำไปประเมินแบบจำลองแต่ละชุดด้วย

เทคนิค Cross-Validation ซึ่งหมายถึงการแบ่งข้อมูลออกเป็นหลายส่วน (Folds) เพื่อฝึกและทดสอบ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

แบบจำลองซ้ำๆ ในแต่ละชุดค่าพารามิเตอร์ (Kohavi, 1995) วิธีนี้ช่วยให้ได้ผลการประเมินที่มีความเสถียรและลดอคติที่อาจเกิดจากการแบ่งข้อมูลเพียงครั้งเดียว



รูปที่ 2.7 ขั้นตอนการทำงานของ Grid Search Cross-Validation
ที่มา : Fathi et al. (2021)

จากรูปที่ 2.7 เป็นกระบวนการ Grid Search Cross-Validation จะทำงานโดยอัตโนมัติผ่านขั้นตอนต่อเนื่อง ได้แก่

- 1) Pick Parameter Combinations: เริ่มจากการกำหนดชุดของค่าพารามิเตอร์ที่ต้องการทดลอง เช่น Max_Depth หรือ Min_Samples_Split
- 2) Perform K-Fold CV: สำหรับแต่ละชุดค่าพารามิเตอร์ ทำการแบ่งข้อมูลเป็น K ส่วน (เช่น 5 หรือ 10 ส่วน) เพื่อทำการ Cross-Validation โดยในแต่ละรอบจะนำข้อมูลบางส่วนมาใช้ฝึกสอนโมเดล และใช้ส่วนที่เหลือในการทดสอบประสิทธิภาพของโมเดล เพื่อวัดค่าต่างๆ เช่น Accuracy, Recall หรือ F1-score สำหรับชุดพารามิเตอร์นั้นๆ แล้วคำนวณค่าเฉลี่ยประสิทธิภาพจากทุก Fold (เช่น Accuracy เฉลี่ย 0.95, 0.80)
- 3) Repeat: ขั้นตอนนี้จะทำซ้ำกับทุกชุดพารามิเตอร์ที่กำหนดไว้
- 4) Select Best Parameter: เลือกชุดค่าพารามิเตอร์ที่ให้ประสิทธิภาพดีที่สุดเพื่อนำไปใช้กับแบบจำลอง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ข้อดีของ Grid Search Cross-Validation คือความครอบคลุมและความมั่นใจว่าไม่มีการมองข้ามค่าพารามิเตอร์ที่เป็นไปได้ทั้งหมด (Exhaustive Search) ซึ่งช่วยป้องกันปัญหาการเลือกค่าที่ไม่เหมาะสมโดยอาศัยประสบการณ์หรือความรู้สึก อย่างไรก็ตามข้อจำกัดสำคัญคืออาจใช้ทรัพยากรคำนวณสูง เมื่อจำนวนไฮเปอร์พารามิเตอร์หรือขอบเขตของค่าแต่ละตัวมีมาก (Curse of Dimensionality) ทำให้ในกรณีที่มีพารามิเตอร์จำนวนมากหรือชุดข้อมูลขนาดใหญ่ นักวิจัยอาจเลือกใช้วิธีสุ่ม (Random Search) หรือ Bayesian Optimization แทน เพื่อเพิ่มประสิทธิภาพและลดเวลาคำนวณ (Bergstra & Bengio, 2012) Grid Search Cross-Validation จึงเป็นเครื่องมือสำคัญที่ช่วยให้กระบวนการเลือกไฮเปอร์พารามิเตอร์มีความเป็นระบบ เชื่อถือได้ และสามารถปรับให้เหมาะสมกับลักษณะข้อมูลหรือปัญหาได้อย่างเป็นวิทยาศาสตร์ โดยเฉพาะเมื่อนำไปใช้กับงานที่ต้องการผลลัพธ์ที่มั่นคงและโปร่งใส เช่น งานวิจัยหรือการเปรียบเทียบแบบจำลองในเชิงอุตสาหกรรม (Pedregosa et al., 2011)

2.7 แนวคิดและทฤษฎีเกี่ยวข้องกับการประเมินประสิทธิภาพในการทำนายของแบบจำลอง

การประเมินประสิทธิภาพของแบบจำลองเป็นขั้นตอนสำคัญในการวิเคราะห์แบบจำลอง Machine Learning เพื่อให้ทราบถึงความสามารถในการทำนายและข้อจำกัดของแบบจำลอง โดยเฉพาะในปัญหาการจำแนกประเภท (Classification) ซึ่งนิยมใช้เมตริกซ์ต่างๆ เช่น Confusion Matrix และ Classification Report สำหรับการวิเคราะห์เชิงลึก

2.7.1 เมตริกซ์ความสับสน (Confusion Matrix)

Confusion Matrix คือเครื่องมือที่ช่วยประเมินและวิเคราะห์ผลการทำนายของแบบจำลองจำแนกประเภท โดยจะแสดงจำนวนของแต่ละกรณีที่แบบจำลองทำนายถูกหรือผิดเมื่อเปรียบเทียบกับค่าจริง (Fawcett, 2006; Powers, 2011) Confusion Matrix ประกอบด้วย 4 องค์ประกอบหลักดังนี้

- 1) True Positive (TP) คือ จำนวนตัวอย่างที่แบบจำลองสามารถทำนายได้ถูกต้องว่าอยู่ในกลุ่มบวก (Positive Class)
- 2) True Negative (TN) คือ จำนวนตัวอย่างที่แบบจำลองทำนายได้ถูกต้องว่าอยู่ในกลุ่มลบ (Negative Class)
- 3) False Positive (FP) คือ จำนวนตัวอย่างที่แบบจำลองทำนายผิดโดยระบุว่าในกลุ่มบวก ทั้งที่จริงแล้วเป็นกลุ่มลบ (Type I Error)
- 4) False Negative (FN) คือ จำนวนตัวอย่างที่แบบจำลองทำนายผิดโดยระบุว่าในกลุ่มลบ ทั้งที่จริงแล้วเป็นกลุ่มบวก (Type II Error)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เมทริกซ์ความสับสน (Confusion Matrix) นี้ทำให้ผู้ใช้งานเข้าใจประสิทธิภาพของแบบจำลองได้อย่างลึกซึ้ง ทั้งในด้านการวัดความถูกต้องและการวิเคราะห์ข้อผิดพลาดจากการทำนาย Confusion Matrix สามารถนำเสนอได้ 2 รูปแบบ คือ

1) ค่าที่ทำนายกับค่าจริง (Predict-Actual)

Actual \ Predicted	Default	Non default
Default	True Positive (TP)	False Positive (FP)
Non default	False Negative (FN)	True Negative (TN)

รูปที่ 2.8 เมทริกซ์ความสับสน (Confusion Matrix) แบบค่าที่ทำนายกับค่าจริง (Predict-Actual)

2) ค่าจริงกับค่าที่ทำนาย (Actual-Predict)

Actual \ Predicted	Default	Non default
Default	True Positive (TP)	False Negative (FN)
Non default	False Positive (FP)	True Negative (TN)

รูปที่ 2.9 เมทริกซ์ความสับสน (Confusion Matrix) แบบค่าจริงกับค่าที่ทำนาย (Actual-Predict)

Confusion Matrix ทำให้ผู้ใช้งานเข้าใจลักษณะข้อผิดพลาดของแบบจำลองได้อย่างลึกซึ้ง ทั้งด้านการวัดความถูกต้องและการวิเคราะห์ข้อผิดพลาดในแต่ละกลุ่มเป้าหมาย (Géron, 2019)

2.7.2 การจำแนกประเภท (Classification Report)

Classification Report เป็นการแสดงผลวิเคราะห์ประสิทธิภาพของแบบจำลองจำแนกประเภทในระบบ Machine Learning โดยจะรายงานตัวชี้วัดที่สำคัญ เช่น ความถูกต้อง (Accuracy) ความระลึก (Recall) ความแม่นยำ (Precision) และค่าประสิทธิภาพโดยรวม (F1-Score) ซึ่งแต่ละค่ามีความสำคัญและช่วยสะท้อนประสิทธิภาพของแบบจำลองในแง่มุมต่างๆ (Saito & Rehmsmeier, 2015; Powers, 2011) ในการจำแนกประเภท (Classification Report) ประกอบด้วยตัวชี้วัดทางสถิติที่ใช้ประเมินผลลัพธ์ของแบบจำลองจากค่าที่พยากรณ์ไว้เปรียบเทียบกับค่าจริงในชุดข้อมูลทดสอบ โดยมีรายละเอียดดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- 1) ค่าความเที่ยงตรงหรือความถูกต้อง (Accuracy) คือ อัตราส่วนของจำนวนข้อมูลที่ทำนายได้ถูกต้องทั้งหมดต่อจำนวนข้อมูลทั้งหมด คำนวณโดยใช้สูตร

$$Accuracy = \frac{True\ Positive + True\ Negative}{True\ Positive + True\ Negative + False\ Positive + False\ Negative} \quad (2.3)$$

- 2) ค่าความครบถ้วนหรือค่าความระลึก (Recall) คือ อัตราส่วนของจำนวนข้อมูลที่ทำนายได้ถูกต้องในคลาสบวกต่อจำนวนข้อมูลที่ทำนายได้จริงของคลาสบวก คำนวณโดยใช้สูตร

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative} \quad (2.4)$$

- 3) ค่าความแม่นยำ (Precision) คือ อัตราส่วนของจำนวนข้อมูลที่ทำนายได้ถูกต้องในคลาสบวกต่อข้อมูลที่ถูกตั้งในคลาสบวกและจำนวนข้อมูลที่ทำนายผิดพลาดในคลาส คำนวณโดยใช้สูตร

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive} \quad (2.5)$$

- 4) ค่าประสิทธิภาพโดยรวม (F1-Score หรือ F-Measure) เป็นค่าเฉลี่ยแบบ Harmonic ของค่าความระลึก (Recall) และค่าความแม่นยำ (Precision) ให้ค่าระหว่าง 0 ถึง 1 ซึ่งค่าที่สูงขึ้นหมายความว่าประสิทธิภาพของแบบจำลองดีขึ้น คำนวณโดยใช้สูตร

$$F1\ Score = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (2.6)$$

โดยค่า F1-Score สามารถอธิบายได้ดังนี้

- 1) ถ้าค่า F1-Score มีค่าสูง หมายความว่า Precision และ Recall มีค่าดีทั้งคู่
- 2) ถ้าค่า F1-Score ต่ำ แสดงว่ามีอย่างน้อยหนึ่งค่าที่ต่ำ หรือทั้งสองค่าต่ำ (Géron, 2019)

Classification Report ยังมักแสดงค่า “Support” คือ จำนวนตัวอย่างจริงในแต่ละกลุ่ม เพื่อประกอบการวิเคราะห์เชิงปริมาณและเปรียบเทียบประสิทธิภาพของแบบจำลองอย่างเป็นระบบ โดยเฉพาะในปัญหาที่ข้อมูลแต่ละกลุ่มมีขนาดไม่เท่ากัน (Imbalanced Data) (Géron, 2019).

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.7.3 การทดสอบสมมติฐานด้วย McNemar's Test

McNemar's Test เป็นสถิติที่ใช้สำหรับทดสอบความแตกต่างของสัดส่วนระหว่างตัวแปรที่จัดอยู่ในรูปแบบ "จับคู่" (Paired Nominal Data) โดยเฉพาะกรณีที่ต้องการเปรียบเทียบผลลัพธ์ของแบบจำลองสองแบบจำลองหรือสองวิธีการที่นำไปใช้กับตัวอย่างเดียวกัน เช่น การเปรียบเทียบความแม่นยำของแบบจำลอง Machine Learning สองแบบ โดยพิจารณาผลลัพธ์ที่แบบจำลองทำนายผิด - ถูกในชุดข้อมูลเดียวกัน

McNemar's Test จะทำงานกับข้อมูลประเภท 2x2 Contingency Table ที่ประกอบด้วยผลลัพธ์ของแต่ละวิธีการ/แบบจำลองในรูปแบบ binary (เช่น ถูก/ผิด, ใช่/ไม่ใช่) และทดสอบว่าความน่าจะเป็นของการเปลี่ยนแปลงระหว่างสองวิธีการมีความแตกต่างกันหรือไม่ ดังตารางต่อไปนี้

ตารางที่ 2.1 ตัวอย่างการทำนายที่แตกต่างกันระหว่างแบบจำลอง A และแบบจำลอง B สำหรับการทดสอบ McNemar's Test

กลุ่มย่อย	แบบจำลอง B ถูก	แบบจำลอง B ผิด
แบบจำลอง A ถูก	a	b
แบบจำลอง A ผิด	c	d

McNemar's Test จะพิจารณาเฉพาะค่าที่ไม่ตรงกัน (Off-Diagonal) คือ b กับ c

สมมติฐาน

สมมติฐานศูนย์ (H_0): ความน่าจะเป็นที่แบบจำลอง A และแบบจำลอง B ทำนายแตกต่างกันในทิศทางหนึ่ง เท่ากับอีกทิศทางหนึ่ง ($b = c$)

สมมติฐานทางเลือก (H_1): ความน่าจะเป็นที่แบบจำลอง A และแบบจำลอง B ทำนายแตกต่างกันในทิศทางหนึ่ง ไม่เท่ากับอีกทิศทางหนึ่ง ($b \neq c$)

สูตรการคำนวณ

$$\chi^2 = \frac{(b - c)^2}{b + c} \quad (2.7)$$

โดยที่

b คือ จำนวนตัวอย่างที่แบบจำลอง A ผิด แต่แบบจำลอง B ถูก

c คือ จำนวนตัวอย่างที่แบบจำลอง A ถูก แต่แบบจำลอง B ผิด

โดยค่าที่ได้จะนำไปเปรียบเทียบกับตาราง Chi-Square Distribution ที่ระดับอิสระ 1 (df = 1) เพื่อพิจารณาค่า p-value เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2.8 แนวคิดและทฤษฎีเกี่ยวข้องกับความสำคัญของคุณลักษณะ (Feature Importance)

ในการสร้างและวิเคราะห์แบบจำลอง Machine Learning หรือ Data Mining การทำความเข้าใจว่าแต่ละคุณลักษณะ (Feature) มีผลกระทบหรือมีบทบาทสำคัญต่อผลลัพธ์ของแบบจำลองมากน้อยเพียงใด ถือเป็นประเด็นสำคัญที่ช่วยเพิ่มทั้งประสิทธิภาพและความน่าเชื่อถือของการวิเคราะห์ซึ่งในทางทฤษฎี "ความสำคัญของคุณลักษณะ" (Feature Importance) คือกระบวนการประเมินหรือจัดอันดับคุณลักษณะตามระดับอิทธิพลที่แต่ละตัวแปรมีต่อการทำนายหรือการตัดสินใจของแบบจำลอง (Kuhn & Johnson, 2013) การวัดความสำคัญของคุณลักษณะมีจุดประสงค์หลักเพื่อ

- 1) เข้าใจกลไกการทำงานของแบบจำลองและการตีความผลลัพธ์ (Interpretability)
- 2) ช่วยในการเลือกหรือลดจำนวนคุณลักษณะ (Feature Selection) เพื่อให้แบบจำลองง่ายต่อการใช้งานและลดความซับซ้อน
- 3) เพิ่มประสิทธิภาพการทำนาย ลด Overfitting และปรับปรุงเวลาในการประมวลผล

แนวคิดนี้จึงเป็นพื้นฐานสำคัญในการสร้างแบบจำลองที่ แนวทางการประเมินความสำคัญของคุณลักษณะมีหลากหลายวิธี เช่น

Logistic Regression ความสำคัญของคุณลักษณะสามารถวัดได้จากค่าสัมประสิทธิ์ (Coefficient) ของแต่ละตัวแปรต้นในสมการโลจิสติก แบบจำลองนี้จะประมาณสัมประสิทธิ์ผ่านการฝึก (Training) โดยแต่ละค่าจะบอกทิศทางและขนาดของอิทธิพลต่อโอกาสที่ผลลัพธ์เป้าหมายจะเกิดขึ้น หากค่าสัมประสิทธิ์ของคุณลักษณะใดมีค่าสูง (บวกหรือลบ) แสดงว่าคุณลักษณะนั้นมีบทบาทสำคัญต่อการจำแนกประเภท ตัวอย่างเช่น การแปลงค่าสัมประสิทธิ์ให้อยู่ในรูป Odds Ratio จะช่วยตีความอิทธิพลในเชิงปฏิบัติได้ง่ายขึ้น (Kuhn & Johnson, 2013) อย่างไรก็ตาม แบบจำลองนี้อาจได้รับผลกระทบจาก Multicollinearity และอิทธิพลที่คาดเคลื่อนหากข้อมูลไม่ถูกเตรียมไว้อย่างเหมาะสม

Decision Tree สามารถให้คะแนนความสำคัญของแต่ละคุณลักษณะโดยวัดจากการลดค่าความบริสุทธิ์ (Impurity Reduction) เช่น Gini Importance หรือ Information Gain ทุกครั้งที่ตัวแปรต้นถูกใช้แบ่งข้อมูลในแต่ละโหนด หากคุณลักษณะใดถูกเลือกใช้แยกข้อมูลบ่อยและช่วยเพิ่มความบริสุทธิ์ได้มาก จะถือว่ามีความสำคัญสูง (Breiman et al., 1984) การแสดงความสำคัญใน Decision Tree จึงสะท้อนในแง่การแบ่งกลุ่มข้อมูลเป็นหลัก ข้อดีคือสามารถตีความได้ง่ายและรองรับข้อมูลเชิงหมวดหมู่ได้ดี

Neural Network โดยเฉพาะโครงข่ายแบบหลายชั้น (MLP) ความสำคัญของคุณลักษณะไม่สามารถดูได้โดยตรงจากค่าน้ำหนักเหมือนแบบจำลองเชิงเส้น จึงต้องใช้เทคนิคเสริม เช่น Permutation Importance ที่วัดผลกระทบของการสลับค่าคุณลักษณะ หรือใช้เทคนิคเช่น SHAP (SHapley Additive exPlanations) เพื่อประเมินอิทธิพลของแต่ละคุณลักษณะต่อการทำนายในแต่ละเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ละตัวอย่าง (Lundberg & Lee, 2017; Molnar, 2022) เทคนิคเหล่านี้ช่วยให้แบบจำลองซับซ้อนที่เปรียบเสมือน “Black-Box” สามารถตีความผลได้มากขึ้น

Long Short-Term Memory (LSTM) ซึ่งเป็น Neural Network สำหรับข้อมูลลำดับ การวัดความสำคัญของคุณลักษณะโดยตรงทำได้ยากเช่นเดียวกับ Neural Network ทั่วไป จึงนิยมใช้วิธี Permutation Importance, SHAP หรือ Attention Mechanism เพื่อประเมินอิทธิพลของแต่ละฟีเจอร์ ตัวอย่างเช่น การสลับหรือปิดข้อมูลฟีเจอร์หนึ่งขณะประเมินผลลัพธ์ แล้ววัดการเปลี่ยนแปลงของผลลัพธ์ (Molnar, 2022) นอกจากนี้ เทคนิค Attention ยังช่วยชี้ให้เห็นว่าข้อมูลช่วงเวลาใดหรือฟีเจอร์ใดที่ LSTM ให้ความสำคัญในการตัดสินใจ

Bidirectional LSTM (BiLSTM) มีแนวคิดพื้นฐานเดียวกับ LSTM ในการประเมินความสำคัญของคุณลักษณะ แต่เพิ่มความซับซ้อนเนื่องจากพิจารณาบริบททั้งสองทิศทาง (Forward และ Backward) ดังนั้นการใช้เทคนิค SHAP, Permutation Importance หรือ Attention ยังเป็นแนวทางที่นิยมสำหรับ BiLSTM เช่นกัน เพื่อวิเคราะห์ว่าแต่ละฟีเจอร์ส่งผลกระทบต่อผลลัพธ์ในแต่ละตำแหน่งของลำดับข้อมูล (Lundberg & Lee, 2017; Molnar, 2022)

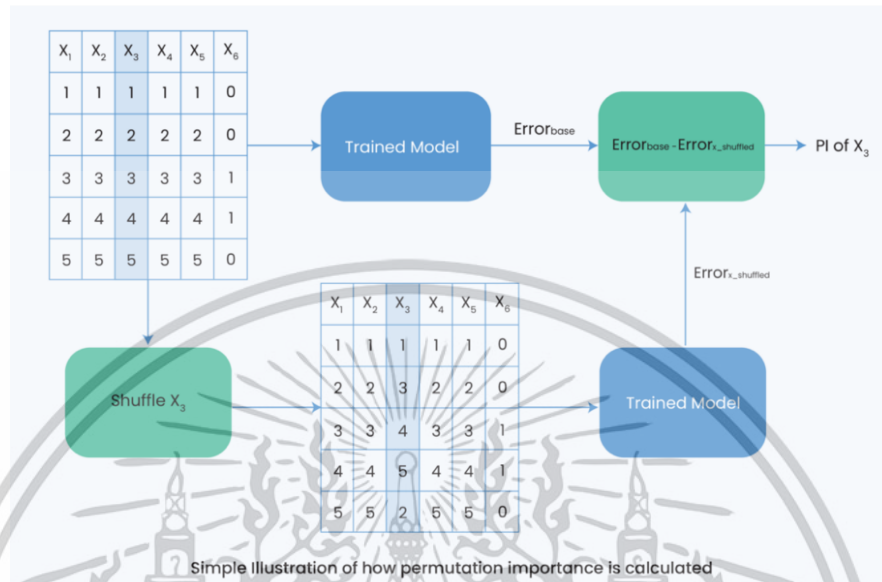
การเข้าใจความสำคัญของคุณลักษณะไม่เพียงช่วยให้แบบจำลองมีประสิทธิภาพสูงขึ้นเท่านั้น แต่ยังช่วยให้การตัดสินใจหรือการแปลผลแบบจำลองเป็นไปอย่างโปร่งใส ตอบโจทย์การนำไปใช้งานจริงในภาคธุรกิจ การแพทย์ การเงิน และงานวิจัยที่ต้องการเหตุผลอ้างอิงประกอบการตัดสินใจ นอกจากนี้การเลือกใช้คุณลักษณะสำคัญ (Feature Selection) ยังช่วยลดความซับซ้อนและต้นทุนข้อมูล รวมถึงช่วยลดโอกาสเกิดปัญหา Multicollinearity ที่อาจส่งผลกระทบต่อแบบจำลองเชิงเส้น (Kuhn & Johnson, 2013; Guyon & Elisseeff, 2003)

การเรียนรู้ของเครื่อง (Machine Learning) ได้กลายมาเป็นเครื่องมือสำคัญในการสร้างแบบจำลอง โดยเฉพาะในด้านการเงินและสินเชื่อ อย่างไรก็ตาม แบบจำลองที่มีความซับซ้อน เช่น Neural Networks แม้จะมีความสามารถในการทำนายที่แม่นยำ แต่กลับมีลักษณะเป็น “กล่องดำ” (Black Box) ที่ยากต่อการตีความว่าการตัดสินใจใดๆ มาจากฟีเจอร์ใดบ้าง การวิเคราะห์ความสำคัญของตัวแปร (Feature Importance) จึงมีบทบาทสำคัญในการ “เปิดกล่องดำ” เพื่อทำความเข้าใจว่าแบบจำลองใช้อะไรเป็นเกณฑ์ในการคาดการณ์ (Molnar, 2022) โดยเทคนิคการวิเคราะห์ความสำคัญที่ได้รับคามนิยมสูงและมีรากฐานทางทฤษฎีที่มั่นคง ได้แก่ Permutation Importance และ SHAP (SHapley Additive exPlanations)

2.8.1 ความสำคัญของคุณลักษณะแบบ Permutation Importance

เป็นวิธีที่ไม่ขึ้นกับชนิดของแบบจำลอง (Model-Agnostic) และสามารถนำไปใช้กับแบบจำลองใดก็ได้ วิธีนี้ถูกนำเสนอและใช้อย่างแพร่หลายในบริบทของการสร้าง Random Forest โดย Breiman (2001) และภายหลังได้ถูกประยุกต์ใช้อย่างกว้างขวางในงาน Machine Learning สมัยใหม่ (Fisher et al., 2019) หลักการของ Permutation Importance คือ การสลับค่าของเอกสารนี้เป็นเอกสารที่สแกนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

พีเจอรไอดีพีเจอรหนึ่งแบบสุมในชุดข้อมูลทดสอบ แล้วประเมินว่าค่าความเที่ยงตรง (Accuracy) หรือค่าประเมินอื่นๆ เปลี่ยนแปลงไปมากน้อยเพียงใด หากการสลับพีเจอรนั้นส่งผลให้ประสิทธิภาพของแบบจำลองลดลงมาก แสดงว่าพีเจอรนั้นมีความสำคัญสูง



รูปที่ 2.10 ขั้นตอนการคำนวณ Permutation Importance

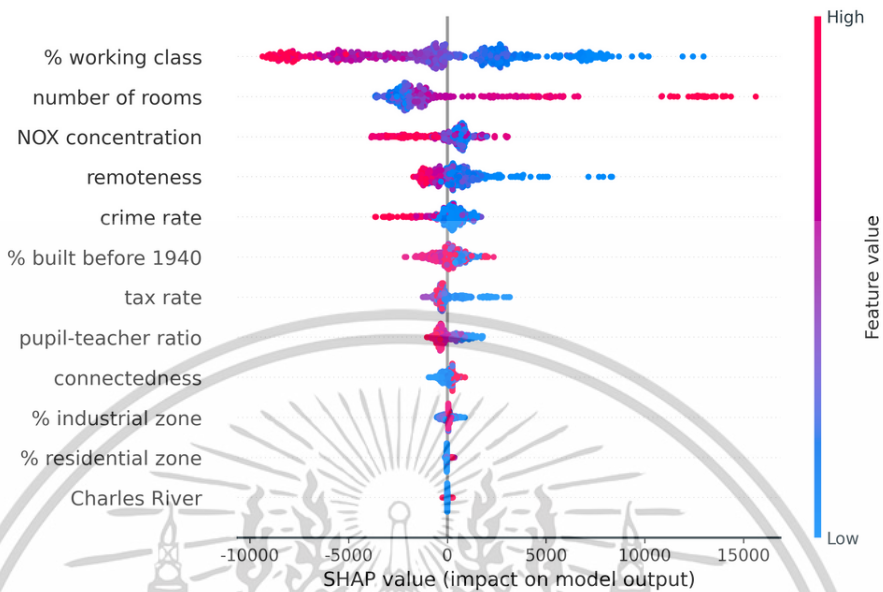
ที่มา: Alon Gubkin (2022)

รูปที่ 2.10 แสดงขั้นตอนการประเมิน Permutation Importance อย่างชัดเจน โดยเริ่มจากการทำนายด้วยข้อมูลเดิม และต่อมาทำการสลับค่าของตัวแปร X_3 ก่อนจะนำไปคำนวณค่าความผิดพลาดใหม่ แล้วจึงหาค่าความแตกต่างเพื่อประเมินระดับความสำคัญของพีเจอรดังกล่าว ข้อดีของวิธีนี้คือความเรียบง่ายและความสามารถในการใช้งานร่วมกับแบบจำลองที่ซับซ้อนทุกชนิด ขณะที่ข้อจำกัดคืออาจใช้เวลานานเมื่อจำนวนพีเจอรมีมาก และอาจเกิดความผิดพลาดหากพีเจอรมีความสัมพันธ์กันสูง (Multicollinearity)

2.8.2 SHAP (SHapley Additive exPlanations)

SHAP เป็นอีกหนึ่งเทคนิคสำคัญที่ถูกพัฒนาขึ้นเพื่อใช้ตีความแบบจำลองโดยเฉพาะในระดับรายบุคคล (Instance-Level Explanation) โดยอ้างอิงจาก Shapley Value ซึ่งมีต้นกำเนิดมาจากทฤษฎีเกม (Cooperative Game Theory) โดยนักคณิตศาสตร์ Lloyd Shapley (1953) หลักการของ SHAP คือการคำนวณค่าผลกระทบของแต่ละพีเจอรโดยเปรียบเสมือนว่าแต่ละพีเจอรเป็นผู้เล่นในเกม และผลลัพธ์ของแบบจำลองเป็นผลรวมของคะแนนที่แต่ละผู้เล่นมีส่วนร่วม (Lundberg & Lee, 2017) โดยการรวมค่าผลกระทบของพีเจอรทั้งหมดจะเท่ากับค่าทำนายของแบบจำลองในกรณีนั้นๆ เสมอ SHAP สามารถอธิบายได้ทั้งในระดับ Global และ Local กล่าวคือสามารถบอกได้ว่าโดยรวมแล้วพีเจอรใดสำคัญที่สุด (Global Importance) และการวิเคราะห์ว่าแบบจำลองให้ผลลัพธ์เอกสารนี้เป็นเอกสารที่สวอนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การทำนายแบบใดกับตัวอย่างข้อมูลเฉพาะราย และฟิเจอร์ไคเป็นปัจจัยสำคัญที่ผลักดันให้แบบจำลองตัดสินใจเช่นนั้น (Local Explanation)



รูปที่ 2.11 ตัวอย่าง Beeswarm Plot ของค่า SHAP

ที่มา: Cooper (2021)

Beeswarm Plot เป็นเครื่องมือที่สำคัญของ SHAP (SHapley Additive exPlanations) สำหรับการอธิบายความสำคัญของฟิเจอร์และการกระจายอิทธิพลของแต่ละตัวแปรที่มีต่อผลลัพธ์ของ Machine Learning โดยแผนภาพนี้สามารถให้ข้อมูลเชิงลึกได้ทั้งในเชิงปริมาณและเชิงคุณภาพ ดังนี้

1) โครงสร้างของ Beeswarm Plot

แกน Y แสดงรายชื่อของตัวแปรอินพุต (Input Variables) ที่เรียงลำดับจากบนลงล่างตามค่าเฉลี่ยสัมบูรณ์ของ SHAP (Mean Absolute SHAP Values) ซึ่งบ่งบอกถึงความสำคัญของแต่ละฟิเจอร์ต่อการทำนายของแบบจำลองในภาพรวม

แกน X แสดงค่า SHAP (SHAP value) ของแต่ละตัวอย่างในชุดข้อมูล ซึ่งสะท้อนถึงผลกระทบของค่าฟิเจอร์นั้นๆ ต่อการทำนายของแบบจำลอง (ค่าเป็นบวกช่วยเพิ่มค่าทำนาย, ค่าเป็นลบช่วยลดค่าทำนาย)

จุดแต่ละจุดแทนตัวอย่างข้อมูลหนึ่งแถว จุดเหล่านี้จะกระจายตามแกน X ตามค่า SHAP ของแต่ละกรณี หากบริเวณใดมีความหนาแน่นสูง จุดจะเรียงซ้อนกันในแนวตั้งคล้ายฝูงผึ้ง

2) ความหมายของสี

สีของแต่ละจุดแสดงถึงค่าดิบของตัวแปรในแต่ละกรณี (ไม่ใช่ค่า SHAP) โดยการใช้การไล่ระดับสีตั้งแต่สีน้ำเงิน (ค่าต่ำ) ถึงสีชมพู/แดง (ค่าสูง)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การพิจารณาการกระจายของสีในแต่ละแถว ช่วยให้เข้าใจทิศทางและความสัมพันธ์ระหว่างค่าตัวแปรกับผลกระทบต่อการทำนาย เช่น หากจุดสีแดงกระจายไปทางด้านขวา แสดงว่าค่าสูงของฟีเจอร์นั้นมีแนวโน้มเพิ่มค่าทำนายของแบบจำลอง

3) การตีความ Beeswarm Plot

Beeswarm plot ไม่ได้แสดงเฉพาะลำดับความสำคัญของฟีเจอร์ (เหมือน Bar Plot) แต่ยังช่วยให้เห็นความสัมพันธ์เชิงลึกระหว่างค่าฟีเจอร์กับการทำนาย เช่น สามารถสังเกตได้ว่าค่าฟีเจอร์ที่สูงหรือต่ำมีแนวโน้มทำให้ผลทำนายเปลี่ยนไปในทิศทางใด

การกระจายตัวของค่า SHAP ตามแกน X สะท้อนถึงระดับอิทธิพลของแต่ละฟีเจอร์ ตัวแปรที่มีการกระจายกว้างแสดงว่ามีผลกระทบต่อผลลัพธ์สูง ส่วนตัวแปรที่มีการกระจายแคบมีผลกระทบต่อผลลัพธ์จำกัด

Beeswarm Plot ให้ข้อมูลทั้งเชิงปริมาณและเชิงคุณภาพในแผนภาพเดียว ช่วยให้เข้าใจการตัดสินใจของแบบจำลองได้ชัดเจนมากขึ้น เหมาะกับการสื่อสารให้ทั้งผู้เชี่ยวชาญและผู้มีส่วนได้ส่วนเสียที่ไม่ใช่สายเทคนิค

2.9 งานวิจัยที่เกี่ยวข้อง

Ala'raj et al. (2021) ได้ศึกษาและพัฒนาแบบจำลอง (เพื่อทำนายความน่าจะเป็นของการขาดชำระหนี้ในเดือนถัดไปสำหรับลูกค้าแต่ละราย โดยใช้ข้อมูลพฤติกรรมชำระหนี้ย้อนหลัง 6 เดือน ผลการทดลองบนชุดข้อมูลบัตรเครดิตของธนาคารได้หวั่น พบว่าแบบจำลอง Bidirectional LSTM ให้ประสิทธิภาพสูงสุดเมื่อเปรียบเทียบกับเทคนิคแบบดั้งเดิม ได้แก่ Logistic Regression, Support Vector Machine (SVM), Random Forest, Multi-layer Perceptron (MLP) และ Gradient Boosting โดยแบบจำลอง Bidirectional LSTM ให้ค่าความเที่ยงตรง (Accuracy) สูงถึง 82.4% และแสดงผลลัพธ์ที่มีนัยสำคัญทางสถิติเมื่อประเมินด้วย McNemar Test

Lessmann et al. (2015) ได้ศึกษาการเปรียบเทียบประสิทธิภาพของแบบจำลอง Machine Learning ในการทำนายความเสี่ยงการผิดนัดชำระหนี้ (Credit Scoring) โดยครอบคลุมทั้งแบบจำลองคลาสสิก เช่น Logistic Regression และ Decision Tree ตลอดจนแบบจำลองสมัยใหม่อย่าง Neural Network, Random Forest, Gradient Boosting และ Support Vector Machine โดยใช้ชุดข้อมูลสินเชื่อผู้บริโภคจากหลากหลายแหล่ง ผลการศึกษาพบว่าแบบจำลองกลุ่ม Ensemble (เช่น Random Forest, Gradient Boosting) และ Neural Network ประสิทธิภาพสูงกว่าวิธีคลาสสิก ในขณะที่เดียวกัน Logistic Regression ยังคงเป็นแบบจำลองที่มีจุดแข็งด้านความโปร่งใสและการตีความ ข้อค้นพบนี้ชี้ให้เห็นถึงความสำคัญของการเลือกเทคนิคที่เหมาะสมกับลักษณะข้อมูลและเป้าหมายการวิเคราะห์ อีกทั้งยังสนับสนุนให้มีการใช้เทคนิค Cross-Validation ในการประเมินเปรียบเทียบแบบจำลองอย่างเป็นระบบ เพื่อให้ได้ผลลัพธ์ที่น่าเชื่อถือและทั่วไปได้ดี

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Refaeilzadeh et al. (2009) ได้ทำการเปรียบเทียบวิธีการแบ่งข้อมูลทั้งสองแบบในบริบทของการสร้างแบบจำลองการจำแนกประเภท พบว่าค่า F1-Score และค่า Accuracy จากการใช้ k-fold Cross-Validation มีความเสถียรและสะท้อนความสามารถของแบบจำลองได้ดีกว่าการแบ่งแบบ Hold-out โดยเฉพาะในชุดข้อมูลที่มีขนาดไม่ใหญ่มาก นอกจากนี้ Cross-Validation ยังช่วยป้องกันปัญหาการเลือกพารามิเตอร์ที่เหมาะสมเกินไปกับชุดข้อมูลใดชุดหนึ่ง (Overfitting) ส่งผลให้แบบจำลองมีความสามารถในการนำไปใช้กับข้อมูลใหม่ได้ดีขึ้น

อโนทัย พุทธาริ และคณะ (2018) ได้ทำการศึกษาเรื่องเครื่องมือในการประเมินคุณภาพสินเชื่อ มีวัตถุประสงค์เพื่อสร้างความรู้ความเข้าใจเกี่ยวกับแนวคิดและการประยุกต์ใช้ Credit risk indicator 2 ประเภท ได้แก่ 1) การใช้ Default rate ในการพิจารณาคุณภาพสินเชื่อควบคู่กับ NPL ratio และ 2) การใช้ Credit Score เพื่อเป็นเครื่องมือในการติดตามและพยากรณ์ความเสี่ยงที่อาจเกิดขึ้นในอนาคตและการบริหารความเสี่ยงหรือกำกับดูแลเชิงรุก โดยบทสรุป คือ การจัดทำแบบจำลอง Credit Score เพื่อใช้ในการแยกแยะลูกหนี้ดี-เสีย และคาดการณ์แนวโน้มการผิดนัดชำระหนี้ภายใน 1 ปี ข้างหน้าจะช่วยเพิ่มประสิทธิภาพและพัฒนาแนวทางการบริหารความเสี่ยงของสถาบันการเงินให้เอื้อต่อแนวทางการตรวจสอบของธนาคารแห่งประเทศไทยที่เน้นกำกับ ดูแลเป็นไปในเชิงรุกมากขึ้น

สำนักงานสถิติแห่งชาติ (2019) ได้ทำการศึกษานี้สินของครัวเรือนเกษตร ปี พ.ศ. 2560 โดยมีจุดประสงค์เพื่อศึกษาปัจจัยที่มีอิทธิพลต่อการเป็นหนี้ของครัวเรือนเกษตร โดยใช้สถิติการวิเคราะห์การถดถอยโลจิสติกแบบสองกลุ่ม (Binary Logistic Regression) พบว่ามี 8 ตัวแปร จากทั้งหมด 13 ตัวแปร ที่มีผลต่อการเป็นหนี้ของครัวเรือนเกษตร ได้แก่ ภาค สถานภาพสมรส เพศ จำนวนสมาชิกในครัวเรือน จำนวนพื้นที่ทำการเกษตร การเป็นเจ้าของที่ดิน ค่าใช้จ่าย และมูลค่าทรัพย์สิน โดยมีรายละเอียดดังนี้ ครัวเรือนเกษตรที่อยู่ในภาคตะวันออกเฉียงเหนือ มีโอกาสเป็นหนี้สูงกว่าครัวเรือนภาคอื่นๆ หัวหน้าครัวเรือนเกษตรที่มีสถานภาพสมรสมีโอกาสเป็นหนี้มากกว่าสถานภาพอื่นๆ หัวหน้าครัวเรือนเพศชายมีโอกาสเป็นหนี้มากกว่าหัวหน้าครัวเรือนหญิง ครัวเรือนเกษตรที่มีจำนวนสมาชิกเพิ่ม 1 คน จะมีโอกาสเป็นหนี้เพิ่มขึ้น ครัวเรือนเกษตรที่มีจำนวนพื้นที่ทำการเกษตรมากขึ้นมีโอกาสเป็นหนี้เพิ่มขึ้น ครัวเรือนเกษตรที่ไม่ใช้ที่ดินในการทำเกษตร (ผู้ทำประมง ป่าไม้ ลำสัตว์ หางของป่า และการบริการทางการเกษตร) มีโอกาสเป็นหนี้ต่ำกว่าครัวเรือนที่มีที่ดินเป็นของตนเอง ส่วนค่าใช้จ่ายและมูลค่าทรัพย์สินถึงแม้จะมีความสัมพันธ์กับการเป็นหนี้ของครัวเรือนเกษตรอย่างมีนัยสำคัญทางสถิติ แต่การเปลี่ยนแปลงของค่าใช้จ่ายและมูลค่าทรัพย์สินจะทำให้โอกาสที่ครัวเรือนจะเป็นหนี้หรือไม่เป็นหนี้เท่ากัน ในขณะที่เขตการปกครอง อายุของหัวหน้าครัวเรือน การศึกษาของหัวหน้าครัวเรือน อัตราพึ่งพิง และรายได้ ไม่พบความสัมพันธ์กับการเป็นหนี้ของครัวเรือนเกษตรอย่างมีนัยสำคัญทางสถิติ

ขวัญฤทัย ฤคตี (2021) ได้ทำการศึกษารวบรวมการวิเคราะห์ปัจจัยกำหนดการผิดนัดชำระค่างวดใน

ธุรกิจเช่าซื้อรถจักรยานยนต์: กรณีศึกษาบริษัทแห่งหนึ่งในกรุงเทพมหานคร ผลการวิเคราะห์จากเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การศึกษานี้พบว่า ปัจจัยกำหนดที่ส่งผลต่อค่าความน่าจะเป็นในการผัดนัดชำระค่าวงวดสูงขึ้น ได้แก่ ลูกหนี้เป็นเพศชาย ลูกหนี้มีสถานภาพสมรส ลูกหนี้ประกอบอาชีพอิสระ ลูกหนี้อายุงานน้อย ลูกหนี้มีรายได้น้อย ลูกหนี้ไม่มีผู้ค้ำประกัน ลูกหนี้มีประวัติการติดเครดิตบูโร วงเงินสินเชื่อสูง อัตราดอกเบี้ยสูง จำนวนเงินดาวน์ต่ำ และมีการผ่อนชำระในระยะสั้นหรือจำนวนงวดในการผ่อนชำระน้อย ตามลำดับ โดยปัจจัยที่ส่งผลต่อโอกาสในการผัดนัดชำระค่าวงวดรุนแรงที่สุด ได้แก่ การที่ลูกหนี้มีประวัติการติดเครดิตบูโรและมีสถานภาพสมรส ตามลำดับ ปัจจัยที่ช่วยลดโอกาสในการผัดนัดชำระค่าวงวดมากที่สุด ได้แก่ การที่ลูกหนี้มีผู้ค้ำประกันและมีการผ่อนชำระในระยะยาว ตามลำดับ โดยลูกหนี้ที่ได้รับผลการประเมินความน่าจะเป็นมาก หมายถึงบุคคลนั้นมีแนวโน้มจะผัดนัดชำระค่าวงวดสูงจัดเป็นลูกหนี้ชั้นแยะ ในขณะที่ลูกหนี้ที่ได้รับผลการประเมินความน่าจะเป็นน้อยหมายถึงบุคคลนั้นมีแนวโน้มจะผัดนัดชำระค่าวงวดต่ำจัดเป็นลูกหนี้ชั้นดี

บทสรุปจากการทบทวนวรรณกรรม

จากการศึกษาทฤษฎีและงานวิจัยที่เกี่ยวข้อง ทำให้ผู้วิจัยสามารถกำหนดตัวแปรอิสระที่มีอิทธิพลต่อการชำระเงินมีอยู่ ดังนี้ 1) ปัจจัยด้านลักษณะส่วนบุคคล ได้แก่ อายุ เพศ อาชีพ สถานภาพสมรส ที่อยู่อาศัย และ 2) ปัจจัยด้านลักษณะสินเชื่อ ได้แก่ ยอดหนี้ ค่าวงวด ระยะเวลาผ่อน วันค้างชำระ การจัดชั้นหนี้ ประวัติการติดต่อ การเบิกใช้วงเงิน ประวัติการชำระหนี้ และการเข้าโปรแกรมมาตรการช่วยเหลือลูกหนี้ เป็นปัจจัยที่มีอิทธิพลต่อการชำระเงินของลูกหนี้

บทที่ 3

วิธีดำเนินการวิจัย

ในบทนี้ผู้วิจัยจะกล่าวถึงวิธีการดำเนินงานวิจัยโดยใช้กระบวนการมาตรฐานอุตสาหกรรมสำหรับการทำเหมืองข้อมูล (Cross-Industry Standard Process for Data Mining: CRISP-DM) ที่ได้กล่าวไว้ในบทที่ 2 ซึ่งมีขั้นตอนการดำเนินงาน 6 ขั้นตอน ดังนี้

1. การทำความเข้าใจธุรกิจ (Business Understanding)
2. การทำความเข้าใจข้อมูล (Data Understanding)
3. การเตรียมข้อมูล (Data Preparation)
4. การสร้างแบบจำลอง (Modeling)
5. การประเมินผล (Evaluation)
6. การนำแบบจำลองไปใช้งาน (Deployment)

3.1 การทำความเข้าใจธุรกิจ (Business Understanding)

การศึกษาครั้งนี้เริ่มต้นจากการวิเคราะห์ความต้องการทางธุรกิจของสถาบันการเงินซึ่งให้บริการสินเชื่อรายย่อย (Nano Finance) มีเป้าหมายสำคัญคือการลดความเสี่ยงจากการผิดนัดชำระหนี้ของลูกค้า และเพิ่มประสิทธิภาพในการบริหารจัดการสินเชื่ออย่างแม่นยำ ผู้วิจัยได้สัมภาษณ์ผู้เกี่ยวข้องในองค์กร เพื่อทำความเข้าใจปัญหาเชิงธุรกิจ ได้แก่ ความท้าทายในการประเมินความเสี่ยงของลูกค้าใหม่ การขาดระบบแจ้งเตือนลูกค้ากลุ่มเสี่ยง และต้นทุนในการติดตามหนี้ที่สูงเมื่อเกิดการผิดนัดชำระ จากการศึกษาข้อมูลเบื้องต้นพบว่า กลุ่มลูกค้าสินเชื่อรายย่อยส่วนใหญ่มักไม่มีหลักประกัน รายได้ไม่แน่นอน และมีพฤติกรรมชำระหนี้ที่หลากหลาย ส่งผลให้การประเมินความเสี่ยงมีความซับซ้อน จึงจำเป็นต้องนำเทคนิคการทำเหมืองข้อมูล (Data Mining) และการเรียนรู้ของเครื่อง (Machine Learning) มาใช้เพื่อพัฒนาแบบจำลองที่สามารถคาดการณ์โอกาสการผิดนัดชำระหนี้ล่วงหน้าได้อย่างมีประสิทธิภาพ ผู้วิจัยจึงได้กำหนดปัญหาเชิงข้อมูลให้เป็นปัญหาการจัดประเภท (Classification) โดยมี "สถานะการผิดนัดชำระหนี้ในเดือนถัดไป" เป็นตัวแปรตาม ซึ่งจะเป็นพื้นฐานในการพัฒนาแบบจำลองทำนายต่อไปในขั้นตอนถัดไปของกระบวนการ CRISP-DM

3.2 การทำความเข้าใจข้อมูล (Data Understanding)

ขั้นตอนการทำความเข้าใจข้อมูลมีเป้าหมายเพื่อให้สามารถทำความเข้าใจข้อมูลที่มีอยู่ได้อย่างลึกซึ้ง ทั้งในเชิงโครงสร้าง คุณภาพ ความสมบูรณ์ และความสัมพันธ์ของตัวแปรต่างๆ โดยในขั้นตอนนี้ ผู้วิจัยได้ดำเนินการดังต่อไปนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.2.1 การเก็บข้อมูล (Data Collection)

ในขั้นตอนการเก็บรวบรวมข้อมูล ผู้วิจัยได้ใช้ข้อมูลทุติภูมิ (Secondary Data) จากฐานข้อมูลของสถาบันการเงินแห่งหนึ่งในประเทศไทย โดยดึงข้อมูลรายบัญชีของลูกค้าสินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพ (Nano Finance) ที่มีสถานะบัญชีเคลื่อนไหว ณ เดือนมกราคม พ.ศ. 2567 รวมทั้งสิ้น 44,251 บัญชี การดึงข้อมูลดังกล่าวดำเนินการด้วยภาษา SQL (Structured Query Language) ซึ่งเหมาะสมสำหรับการคัดเลือกข้อมูลจากฐานข้อมูลขนาดใหญ่ที่มีความซับซ้อน หลังจากนั้น ผู้วิจัยได้นำข้อมูลที่ได้มาผ่านการประมวลผลเบื้องต้น (Preprocessing) ด้วยโปรแกรม SAS โดยใช้ภาษา SAS เพื่อจัดรูปแบบข้อมูล จากนั้นนำมาจัดการกับข้อมูลต่อบน Jupyter Notebook โดยใช้ภาษาไพธอน (Python)

ตารางที่ 3.1 ไลบรารีบน Python ที่ใช้ในงานวิจัย

ไลบรารี (Library)	วัตถุประสงค์การใช้
pandas	จัดการและประมวลผลข้อมูลตาราง
numpy	คำนวณเชิงตัวเลขและอาร์เรย์
scikit-learn	แบบจำลอง ML และเครื่องมือวิเคราะห์
imbalanced-learn	เทคนิคจัดการข้อมูลไม่สมดุล
tensorflow	สร้างและเทรนแบบจำลอง Deep Learning
keras	สร้างแบบจำลอง Neural Network
matplotlib	สร้างกราฟและ visualization
seaborn	Visualization และกราฟสถิติ

3.2.2 การอธิบายข้อมูล

ขั้นตอนการอธิบายข้อมูล (Data Description) ถือเป็นกระบวนการสำคัญที่ช่วยให้ผู้วิจัยสามารถทำความเข้าใจข้อมูลในภาพรวมก่อนดำเนินการวิเคราะห์เชิงลึกในขั้นตอนถัดไป โดยมุ่งเน้นการสำรวจลักษณะของข้อมูลทั้งในเชิงโครงสร้าง (Structure) เชิงปริมาณ (Quantitative) และเชิงคุณภาพ (Qualitative) เพื่อให้สามารถระบุปัญหาที่อาจเกิดขึ้น เช่น ค่าที่ขาดหาย (Missing Values), ค่าผิดปกติ (Outliers), ความไม่สมดุลของกลุ่มข้อมูล (Imbalanced Data), และความสัมพันธ์ระหว่างตัวแปรต่างๆ (Correlation) ซึ่งล้วนเป็นประเด็นที่ส่งผลต่อความแม่นยำและประสิทธิภาพของแบบจำลองทำนายที่พัฒนาขึ้นในภายหลัง งานวิจัยนี้รวบรวมตัวแปรที่คาดว่าเป็นปัจจัยที่ส่งผลต่อการผิดนัดชำระหนี้ จำนวน 38 ตัวแปร ประกอบด้วยตัวแปรอิสระ 37 ตัวแปรและตัวแปรตาม 1 ตัวแปร โดยแบ่งเป็นข้อมูลเชิงกลุ่มจำนวน 7 ตัวแปร และเป็นอัตราส่วน (Ratio) จำนวน 31 ตัวแปร ดังตารางที่ 3.2

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.2 คุณลักษณะของตัวแปรที่ใช้ในงานวิจัย

ลำดับ	ชื่อตัวแปร	คำอธิบาย	ประเภทข้อมูลตามระดับมาตรวัด
1	AGE	อายุ (ปี)	อัตราส่วน (Ratio)
2	GENDER	เพศ	นามบัญญัติ (Nominal)
3	MARITAL	สถานภาพสมรส	นามบัญญัติ (Nominal)
4	REGION	ภูมิภาคที่อยู่อาศัย	นามบัญญัติ (Nominal)
5	DSCR	อัตราส่วนความสามารถในการชำระหนี้ (เดือน)	อัตราส่วน (Ratio)
6	ChildAmount	จำนวนบุตร (คน)	อัตราส่วน (Ratio)
7	MOB	จำนวนเดือนนับจากวันทำ สินเชื่อ	อัตราส่วน (Ratio)
8	TERM	ระยะเวลากู้ยืม (เดือน)	อัตราส่วน (Ratio)
9	RATE	อัตราดอกเบี้ยของเงินกู้	อัตราส่วน (Ratio)
10	DPD	จำนวนวันค้างชำระ (วัน)	อัตราส่วน (Ratio)
11	UTILIZATION	อัตราส่วนยอดเงินกู้คงเหลือต่อ จำนวนเงินกู้	อัตราส่วน (Ratio)
12	INT_PER_VAL	อัตราส่วนยอดดอกเบี้ยคงเหลือ ต่อจำนวนเงินกู้	อัตราส่วน (Ratio)
13	RAT_OST_BAL_INTL_AMT	อัตราส่วนค้างยอดหนี้	อัตราส่วน (Ratio)
14	AVG_UTILIZATION_3	ค่าเฉลี่ยของอัตราส่วนการเบิกใช้ วงเงิน ใน 3 เดือน	อัตราส่วน (Ratio)
15	FLAG_PRG	การเข้าโปรแกรมมาตรการ ช่วยเหลือลูกหนี้	นามบัญญัติ (Nominal)
16	TBAND	ระดับความเสี่ยงลูกหนี้	อันดับ (Ordinal)
17	CLASS	การจัดชั้นหนี้	อันดับ (Ordinal)
18	NUM_FULL_PMT_LAST_3	จำนวนครั้งที่จ่ายเต็ม ใน 3 เดือน ล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
19	NUM_FULL_PMT_LAST_6	จำนวนครั้งที่จ่ายเต็ม ใน 6 เดือน ล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
20	NUM_FULL_PMT_LAST_9	จำนวนครั้งที่จ่ายเต็ม ใน 9 เดือน ล่าสุด (ครั้ง)	อัตราส่วน (Ratio)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.2 (ต่อ) คุณลักษณะของตัวแปรที่ใช้ในงานวิจัย

ลำดับ	ชื่อตัวแปร	คำอธิบาย	ประเภทข้อมูลตามระดับมาตรวัด
21	NUM_FULL_PMT_LAST_12	จำนวนครั้งที่จ่ายเต็ม ใน 12 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
22	NUM_PART_PMT_LAST_3	จำนวนครั้งที่จ่ายบางส่วน ใน 3 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
23	NUM_PART_PMT_LAST_6	จำนวนครั้งที่จ่ายบางส่วน ใน 6 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
24	NUM_PART_PMT_LAST_9	จำนวนครั้งที่จ่ายบางส่วน ใน 9 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
25	NUM_PART_PMT_LAST_12	จำนวนครั้งที่จ่ายบางส่วน ใน 12 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
26	MAX_DELQ_3	จำนวนงวดที่ค้างชำระมากที่สุด ใน 3 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
27	MAX_DELQ_6	จำนวนงวดที่ค้างชำระมากที่สุด ใน 6 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
28	MAX_DELQ_9	จำนวนงวดที่ค้างชำระมากที่สุด ใน 9 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
29	MAX_DELQ_12	จำนวนงวดที่ค้างชำระมากที่สุด ใน 12 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
30	ACT_STS3	จำนวนครั้งที่ติดต่อกับได้ ใน 3 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
31	ACT_STS6	จำนวนครั้งที่ติดต่อกับได้ ใน 6 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
32	ACT_STS9	จำนวนครั้งที่ติดต่อกับได้ ใน 9 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
33	ACT_STS12	จำนวนครั้งที่ติดต่อกับได้ ใน 12 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
34	N_CT3	จำนวนครั้งที่ติดต่อกับไม่ได้ ใน 3 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดลอกเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.2 (ต่อ) คุณลักษณะของตัวแปรที่ใช้ในงานวิจัย

ลำดับ	ชื่อตัวแปร	คำอธิบาย	ประเภทข้อมูลตามระดับมาตรวัด
35	N_CT6	จำนวนครั้งที่ติดต่อกันไม่ได้ ใน 6 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
36	N_CT9	จำนวนครั้งที่ติดต่อกันไม่ได้ ใน 9 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
37	N_CT12	จำนวนครั้งที่ติดต่อกันไม่ได้ ใน 12 เดือนล่าสุด (ครั้ง)	อัตราส่วน (Ratio)
38	TARGET	การผิมนัดชำระหนี้	นามบัญญัติ (Nominal)

3.3 การเตรียมข้อมูล (Data Preparation)

ขั้นตอนการเตรียมข้อมูลเป็นกระบวนการสำคัญโดยมีวัตถุประสงค์เพื่อปรับปรุงคุณภาพของข้อมูลให้มีความพร้อมในการป้อนเข้าสู่กระบวนการสร้างแบบจำลองโดยมีรายละเอียดดังนี้

3.3.1 การเลือกข้อมูล (Data Selection)

ในการดำเนินการเลือกข้อมูล ได้มีการคัดกรองเฉพาะข้อมูลและตัวแปรที่มีความเกี่ยวข้องโดยตรงกับวัตถุประสงค์ของการวิเคราะห์ เพื่อขจัดข้อมูลที่ไม่จำเป็น ลดจำนวนมิติของชุดข้อมูล และหลีกเลี่ยงปัญหาความซ้ำซ้อนของข้อมูล โดยกำหนดเงื่อนไขให้เลือกเฉพาะบัญชีสินเชื่อที่มีอายุการเปิดบัญชีไม่น้อยกว่า 1 ปี ณ เดือนธันวาคม พ.ศ. 2566 ซึ่งเป็นเงื่อนไขสำคัญในการประเมินพฤติกรรมการชำระหนี้ย้อนหลัง นอกจากนี้ ยังมีการลบคอลัมน์ที่ไม่เกี่ยวข้องกับการวิเคราะห์ ได้แก่ หมายเลขบัญชีสินเชื่อ และหมายเลขลูกหนี้ เนื่องจากเป็นข้อมูลเฉพาะรายที่ไม่มีความสัมพันธ์เชิงวิเคราะห์กับพฤติกรรมการณ์ผิมนัดชำระหนี้ และอาจส่งผลให้เกิดการเรียนรู้ที่ไม่เหมาะสมของแบบจำลอง

3.3.2 การทำความสะอาดข้อมูล (Data Cleansing)

ในการเตรียมข้อมูลเพื่อการวิเคราะห์และการสร้างแบบจำลอง Machine Learning จำเป็นต้องดำเนินการ "ทำความสะอาดข้อมูล" (Data Cleansing) อย่างเป็นระบบ เพื่อให้ได้ชุดข้อมูลที่มีคุณภาพสูง มีความสมบูรณ์ ถูกต้อง และเหมาะสมต่อการเรียนรู้ของแบบจำลอง โดยเฉพาะอย่างยิ่งในการประยุกต์ใช้กับโมเดลเชิงพยากรณ์ที่ต้องการความแม่นยำสูง

ในงานวิจัยนี้ ผู้วิจัยได้ดำเนินการทำความสะอาดข้อมูลตามขั้นตอนที่แสดงไว้ใน ตารางที่ 3.3 ซึ่งเริ่มต้นจากชุดข้อมูลตั้งต้นที่ประกอบด้วยจำนวน 44,251 แถว และ 38 คอลัมน์ โดยแต่ละขั้นตอนจะระบุอย่างชัดเจนว่าได้ดำเนินการลบแถวหรือคอลัมน์หรือไม่ พร้อมอธิบายหลักเกณฑ์ที่ใช้ในการตัดสินใจ ดังรายละเอียดต่อไปนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.3 รายละเอียดแต่ละขั้นตอนการประมวลผลข้อมูล

ลำดับ	ขั้นตอนการประมวลผล	จำนวนแถว (Rows)	จำนวนคอลัมน์ (Columns)
1	ลบข้อมูลซ้ำ	44,251	38
2	ลบคอลัมน์ที่มีข้อมูลสูญหายจำนวนมาก	44,251	38
3	ลบแถวที่มีข้อมูลสูญหายจำนวนมาก	43,870	38
4	ลบคอลัมน์ที่มีค่าเพียงค่าเดียว	43,870	38
5	ลบข้อมูลที่มีความหลากหลายของข้อมูลสูง	43,870	38

ขั้นตอนที่ 1 ตรวจสอบข้อมูลซ้ำซ้อน (Duplicate Records) โดยพิจารณาความซ้ำกันของข้อมูลในระดับแถว หากพบว่าแถวใดมีค่าทุกคอลัมน์เหมือนกันทั้งหมดจะถูกลบทิ้ง อย่างไรก็ตามผลการประมวลผลพบว่าไม่มีข้อมูลใดที่ซ้ำกันทุกค่าในระดับแถว (Row-Level Duplication) จึงไม่มีการลบแถวหรือคอลัมน์ในขั้นตอนนี้ ทั้งนี้การไม่มีข้อมูลซ้ำซ้อนให้เห็นถึงความสมบูรณ์และคุณภาพของข้อมูลเบื้องต้นที่ดี ซึ่งเป็นปัจจัยพื้นฐานที่ส่งผลต่อความน่าเชื่อถือของการวิเคราะห์ในขั้นตอนถัดไป

ขั้นตอนที่ 2 ตรวจสอบคอลัมน์ที่มีข้อมูลสูญหายจำนวนมาก โดยกำหนดเกณฑ์ว่าหากคอลัมน์ใดมีค่าข้อมูลสูญหายเกินร้อยละ 30 ของจำนวนแถวทั้งหมด จะถือว่าไม่สามารถใช้ในการวิเคราะห์ได้อย่างมีประสิทธิภาพ จึงควรถูกตัดออกจากชุดข้อมูล อย่างไรก็ตาม จากการประเมินไม่พบคอลัมน์ใดที่มีค่าข้อมูลสูญหายเกินเกณฑ์ดังกล่าว จึงไม่มีการลบคอลัมน์ในขั้นตอนนี้

ขั้นตอนที่ 3 ตรวจสอบแถวที่มีข้อมูลสูญหาย โดยผู้วิจัยได้กำหนดเงื่อนไขว่าหากแถวใดมีข้อมูลสูญหายตั้งแต่ 3 คอลัมน์ขึ้นไป จะพิจารณาว่าเป็นข้อมูลที่ไม่สมบูรณ์และควรถูกลบออกจากชุดข้อมูล ผลการประมวลผลในขั้นตอนนี้ส่งผลให้มีการลบแถวออกจำนวน 381 แถว ส่งผลให้จำนวนแถวลดลงจาก 44,251 เหลือ 43,870 แถว ขณะที่จำนวนคอลัมน์ยังคงเดิม

ขั้นตอนที่ 4 ตรวจสอบคอลัมน์ที่ไม่มีความหลากหลายของข้อมูล (Zero Variance Features) โดยพิจารณาว่าหากคอลัมน์ใดมีค่าซ้ำกันทั้งหมด หรือมีค่าหนึ่งค่าปรากฏมากกว่าร้อยละ 99 ของข้อมูลทั้งคอลัมน์ จะถือว่าไม่มีประโยชน์ในการวิเคราะห์และควรถูกลบทิ้ง อย่างไรก็ตามจากการตรวจสอบพบว่าไม่มีคอลัมน์ใดเข้าข่ายตามเกณฑ์ จึงไม่มีการลบคอลัมน์ในขั้นตอนนี้

ขั้นตอนที่ 5 ตรวจสอบคอลัมน์ที่มีความหลากหลายน้อย (Low Variance Features) โดยวิเคราะห์ค่าความแปรปรวน (Variance) และค่าเบี่ยงเบนมาตรฐาน (Standard Deviation) หากคอลัมน์ใดมีค่าเบี่ยงเบนมาตรฐานต่ำกว่า 0.01 หรือมีค่าที่พบมากที่สุดครอบคลุมข้อมูลในสัดส่วนสูงเกินไป จะพิจารณาลบทิ้ง ผลการวิเคราะห์พบว่าคอลัมน์ทั้งหมดมีค่าความแปรปรวนสูงกว่าเกณฑ์ที่กำหนด จึงไม่มีการลบคอลัมน์ใดเพิ่มเติมในขั้นตอนนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

โดยสรุปจากกระบวนการทำความสะอาดข้อมูลทั้ง 5 ขั้นตอน พบว่ามีการลบแถวเพียงขั้นตอนเดียว คือ การลบแถวที่มีข้อมูลสูญหายในขั้นตอนที่ 3 ส่งผลให้จำนวนแถวลดลงจาก 44,251 เหลือ 43,870 แถว ส่วนจำนวนคอลัมน์ยังคงอยู่ที่ 38 คอลัมน์ตลอดกระบวนการ การดำเนินการดังกล่าวทำให้ได้ชุดข้อมูลมีความสมบูรณ์ โปร่งใส และเหมาะสมต่อการนำไปใช้ในการวิเคราะห์และการสร้างแบบจำลองอย่างมีประสิทธิภาพในขั้นตอนต่อไป

3.3.3 การเตรียมข้อมูล (Data Preparation)

ในการเตรียมข้อมูลสำหรับการวิเคราะห์และสร้างแบบจำลอง ผู้วิจัยได้ดำเนินการแปลงข้อมูลให้เหมาะสมกับลักษณะของแต่ละตัวแปร โดยข้อมูลที่ใช้ในงานวิจัยนี้สามารถแบ่งออกเป็น 2 ประเภทหลัก ได้แก่

- 1) ข้อมูลเชิงตัวเลข (Numerical Data) ข้อมูลประเภทนี้ประกอบด้วยตัวแปรที่มีลักษณะเป็นตัวเลข หรือมีค่าต่อเนื่อง เช่น อายุ รายได้ จำนวนบุตร เป็นต้น สำหรับข้อมูลเชิงตัวเลขไม่จำเป็นต้องแปลงค่าก่อนนำเข้าแบบจำลองแต่ในบางกรณีเพื่อให้ข้อมูลอยู่ในช่วงที่เหมาะสมต่อการเรียนรู้ของแบบจำลอง จะมีการปรับขนาดข้อมูล เช่น การทำ Normalization หรือ Standardization
- 2) ข้อมูลเชิงกลุ่ม (Categorical Data) ข้อมูลประเภทนี้ประกอบด้วยตัวแปรที่แสดงถึงกลุ่มหรือประเภท ซึ่งไม่สามารถนำไปประมวลผลด้วยแบบจำลองเชิงตัวเลขได้โดยตรง จำเป็นต้องแปลงข้อมูลให้อยู่ในรูปแบบที่แบบจำลองสามารถนำไปใช้งานได้ โดยข้อมูลเชิงกลุ่มสามารถแบ่งย่อยได้อีก 2 ประเภท ได้แก่
 - ก) ข้อมูลนามบัญญัติ (Nominal Data) เป็นข้อมูลที่แสดงถึงประเภทหรือกลุ่มซึ่งไม่มีลำดับความสำคัญ ตัวอย่างเช่น เพศ (ชาย/หญิง), ภูมิภาค (เหนือ/กลาง/อีสาน/ใต้) เป็นต้น ตัวแปรประเภทนี้จะถูกแปลงโดยใช้เทคนิค One-hot Encoding หรือ Dummy Variable ตามที่แสดงในตารางที่ 3.4
 - ข) ข้อมูลเชิงอันดับ (Ordinal Data) เป็นข้อมูลที่แสดงถึงลำดับขั้นหรือระดับ ซึ่งลำดับของข้อมูลมีความหมาย ตัวอย่างเช่น ระดับความเสี่ยง (ต่ำ/ปานกลาง/สูง) เป็นต้น การแปลงข้อมูลประเภทนี้นิยมใช้เทคนิค Label Encoding หรือการแทนค่าด้วยตัวเลขที่สะท้อนลำดับขั้น ดังที่แสดงในตารางที่ 3.5

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.4 คุณลักษณะของตัวแปรต้นที่เป็นข้อมูลนามบัญญัติ

ชื่อตัวแปร	ค่าที่เป็นไปได้
GENDER	ชาย, หญิง
MARITAL	โสด, สมรส, หย่าร้าง, หม้าย
REGION	ภาคกลาง, ภาคเหนือ, ภาคตะวันออกเฉียงเหนือ, ภาคตะวันออก, ภาคใต้, กรุงเทพมหานคร
FLAG_PRG	เข้าร่วม, ไม่เข้าร่วม

ตารางที่ 3.5 คุณลักษณะของตัวแปรต้นที่เป็นข้อมูลเชิงอันดับ

ชื่อตัวแปร	ค่าที่เป็นไปได้
TBAND	ความเสี่ยงต่ำ, ความเสี่ยงปานกลาง, ความเสี่ยงสูง, ความเสี่ยงสูงมาก, ความเสี่ยงสูงสุด
CLASS	กลุ่มที่ไม่มีความเสี่ยงด้านเครดิต, กลุ่มที่มีความเสี่ยงเพิ่มขึ้นอย่างมีนัยยะสำคัญ, กลุ่มที่ไม่ก่อให้เกิดรายได้

3.3.4 การแปลงข้อมูล (Feature Encoding)

จากการสำรวจคุณลักษณะของข้อมูลตัวแปรอิสระ 37 ตัวแปร พบว่าในจำนวนนั้นมีข้อมูลเชิงกลุ่ม (Categorical Features) จำนวน 6 ตัวแปร ซึ่งไม่สามารถนำไปใช้ได้โดยตรงในการเรียนรู้ของเครื่อง เนื่องจากแบบจำลองเหล่านี้ต้องการข้อมูลในเชิงปริมาณหรือเชิงตัวเลข ดังนั้นในขั้นตอนต่อมาคือการนำข้อมูลเชิงกลุ่ม (Categorical Data) ไปแปลงให้อยู่ในรูปแบบตัวเลข โดยการเข้ารหัสข้อมูลเชิงกลุ่มนี้สามารถดำเนินการได้ดังนี้

การเข้ารหัสแบบวัน-ฮอต (One Hot Encoding) วิธีนี้นิยมใช้กับข้อมูลเชิงกลุ่มที่ไม่มีลำดับ (Nominal Data) โดยจะเปลี่ยนค่าของแต่ละกลุ่มให้เป็นคอลัมน์ใหม่ และแทนค่าด้วย 0 หรือ 1 ตามการปรากฏของข้อมูลในแต่ละประเภท ในตัวอย่างรูปที่ 3.1 ตัวแปร “Class” ที่ประกอบด้วยกลุ่ม “A”, “M”, และ “NPL” จะถูกแยกออกเป็นคอลัมน์ย่อยตามแต่ละกลุ่ม และระบุค่า 1 ในตำแหน่งของกลุ่มที่ตรงกัน ส่วนค่าที่ไม่ตรงกันจะเป็น 0 ทำให้แต่ละกลุ่มสามารถแสดงอยู่บนมิติข้อมูลใหม่ที่ไม่ทับซ้อนกัน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

OneHot Encoding

CLASS	A	M	NPL
A	1	0	0
M	0	1	0
NPL	0	0	1

รูปที่ 3.1 การแปลงข้อมูลเชิงกลุ่มแบบไม่มีลำดับด้วยวิธี One Hot Encoding

3.4 การสร้างแบบจำลอง (Modeling)

ในขั้นตอนนี้ ผู้วิจัยได้นำข้อมูลที่ผ่านการเตรียมมาใช้ในการพัฒนาแบบจำลองเพื่อทำนายการผิวน้ำขำระหนึ่ โดยเลือกใช้เทคนิคการเรียนรู้ของเครื่อง (Machine Learning) และเทคนิคการเรียนรู้เชิงลึก (Deep Learning) รวมทั้งสิ้น 5 วิธี ได้แก่ Logistic Regression, Decision Tree, Neural Network, Long Short-Term Memory (LSTM) และ Bidirectional LSTM พร้อมทั้งดำเนินการปรับแต่งไฮเปอร์พารามิเตอร์ (Hyperparameter Tuning) เพื่อค้นหาค่าที่เหมาะสมที่สุดสำหรับแต่ละแบบจำลอง นอกจากนี้ ยังได้มีการแบ่งชุดข้อมูลออกเป็นชุดฝึกฝนและชุดทดสอบ เพื่อประเมินประสิทธิภาพของแต่ละแบบจำลอง ก่อนนำผลลัพธ์ไปเปรียบเทียบและวิเคราะห์ในขั้นตอนถัดไป โดยมีขั้นตอนดังนี้

3.4.1 การเลือกวิธีแบ่งข้อมูลเป็นข้อมูลฝึกฝนและข้อมูลทดสอบ (Training/Testing)

ในการดำเนินการสร้างแบบจำลองเพื่อทำนายการผิวน้ำขำระหนึ่ ผู้วิจัยได้ทำการทดลองด้วยวิธีการแบ่งข้อมูลออกเป็นชุดฝึกฝน (Training Set) และชุดทดสอบ (Testing Set) โดยใช้สองวิธี ได้แก่ วิธีแบบแยกชุด (Hold-out) และวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) โดยมีรายละเอียดขั้นตอนดังนี้

1) วิธีการแบ่งแบบแยกชุด (Hold-out)

ในขั้นตอนนี้ ผู้วิจัยได้แบ่งข้อมูลออกเป็นสองส่วน คือ ชุดข้อมูลฝึกอบรมคิดเป็นร้อยละ 80 และชุดข้อมูลทดสอบคิดเป็นร้อยละ 20 โดยการสุ่มตัวอย่างแบบสุ่มอย่างง่าย (Simple Random Sampling) ชุดข้อมูลฝึกอบรมถูกนำไปฝึกฝนแบบจำลอง ขณะที่ชุดข้อมูลทดสอบถูกใช้ในการประเมินประสิทธิภาพของแบบจำลอง

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
print(f"Training: {X_train.shape[0]}")
print(f"Testing: {X_test.shape[0]}")
```

Training: 35096
Testing: 8774

รูปที่ 3.2 วิธีแบ่งข้อมูลแบบแยกชุด (Hold-Out)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2) วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation)

ในขั้นตอนนี้ ผู้วิจัยแบ่งข้อมูลออกเป็น 5 ส่วน (Folds) ที่มีขนาดเท่ากัน กระบวนการฝึกรูปแบบและทดสอบดำเนินการทั้งหมด 5 รอบ โดยแต่ละรอบใช้ข้อมูล 4 ส่วนในการฝึกรูปแบบ และ 1 ส่วนในการทดสอบ ซึ่งสลับส่วนข้อมูลทดสอบในแต่ละรอบ ประสิทธิภาพของแบบจำลองจะถูกประเมินจากค่าเฉลี่ยของผลลัพธ์ที่ได้จากทั้ง 5 รอบ

```
from sklearn.model_selection import StratifiedKFold

skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

for fold, (train_index, test_index) in enumerate(skf.split(X, y)):
    print(f"Fold {fold+1}")
    |
    X_train, X_test = X.iloc[train_index], X.iloc[test_index]
    y_train, y_test = y.iloc[train_index], y.iloc[test_index]

    print(f"X_train shape: {X_train.shape}")
    print(f"X_test shape: {X_test.shape}")
    print("-" * 30)
```

```
Fold 1
X_train shape: (35096, 47)
X_test shape: (8774, 47)
-----
```

```
Fold 2
X_train shape: (35096, 47)
X_test shape: (8774, 47)
-----
```

```
Fold 3
X_train shape: (35096, 47)
X_test shape: (8774, 47)
-----
```

```
Fold 4
X_train shape: (35096, 47)
X_test shape: (8774, 47)
-----
```

```
Fold 5
X_train shape: (35096, 47)
X_test shape: (8774, 47)
-----
```

รูปที่ 3.3 วิธีแบ่งข้อมูลแบบไขว้ (K-Fold Cross Validation)

3.4.2 การเลือกไฮเปอร์พารามิเตอร์ด้วยการค้นหาแบบกริด (Grid Search)

เพื่อเพิ่มประสิทธิภาพของแบบจำลองในการทำนายการผิנדัดชำระหนี้ ผู้วิจัยได้ดำเนินการปรับแต่งไฮเปอร์พารามิเตอร์ (Hyperparameter Tuning) ของแบบจำลองทั้ง 5 แบบ ได้แก่ Logistic Regression, Decision Tree, Neural Network (MLP), Long Short-Term Memory (LSTM) และ Bidirectional LSTM โดยใช้เทคนิค Grid Search ร่วมกับ K-Fold Cross Validation (k=5) และประเมินผลด้วยค่าประสิทธิภาพโดยรวม (F1-Score) เป็นเกณฑ์ในการเลือกค่าที่ดีที่สุด โดยผู้วิจัยได้ดำเนินการค้นหาไฮเปอร์พารามิเตอร์ที่เหมาะสมที่สุดสำหรับแบบจำลองทั้งหมดสองรอบ (Grid

Search) เพื่อเพิ่มประสิทธิภาพสูงสุดในการทำนาย โดยมีรายละเอียดขั้นตอนดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับใช้งานเพื่อการศึกษาค้นคว้าเท่านั้น เมื่อผู้ยู่ได้เห็นใบโฆษณาการดำเนินการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

- 1) การปรับจูนไฮเปอร์พารามิเตอร์ที่เหมาะสมด้วยการค้นหาแบบกริดของแบบจำลอง
Logistics Regression

ตารางที่ 3.6 การปรับแต่งพารามิเตอร์ของแบบจำลอง Logistic Regression

พารามิเตอร์ที่ปรับแต่ง	พารามิเตอร์ที่เหมาะสม	ค่าประสิทธิภาพโดยรวม
{'C': 0.01, 0.085, 0.1, 1, 10, 100, 'penalty': 'l2', 'solver': 'lbfgs'}	{'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}	0.7768

จากผลการปรับจูนไฮเปอร์พารามิเตอร์ของแบบจำลอง Logistic Regression ได้กำหนดช่วงค่าของพารามิเตอร์ C ตั้งแต่ 0.01 ไปจนถึง 100 เพื่อสำรวจผลกระทบของค่าดังกล่าวต่อประสิทธิภาพของแบบจำลอง โดยใช้ Penalty = 'l2' และ Solver = 'lbfgs' เป็นค่าคงที่ตลอดการทดลอง ผลการปรับจูนพบว่าค่าที่เหมาะสมที่สุดคือ C = 0.085 ซึ่งให้ค่าประสิทธิภาพโดยรวม (F1-Score) ที่ร้อยละ 77.68 ผลลัพธ์ดังกล่าวแสดงให้เห็นว่าการเลือกช่วงค่าของพารามิเตอร์ C อย่างเหมาะสมมีผลต่อความสามารถในการทำนายของแบบจำลอง โดยค่าที่ไม่สูงมากสามารถช่วยควบคุมความซับซ้อนของแบบจำลองและลดความเสี่ยงของการเกิด Overfitting ได้อย่างมีประสิทธิภาพ

- 2) การปรับจูนไฮเปอร์พารามิเตอร์ที่เหมาะสมด้วยการค้นหาแบบกริดของแบบจำลอง
Decision Tree

ตารางที่ 3.7 การปรับแต่งพารามิเตอร์ของแบบจำลอง Decision Tree

พารามิเตอร์ที่ปรับแต่ง	พารามิเตอร์ที่เหมาะสม	ค่าประสิทธิภาพโดยรวม
{'criterion': ['gini', 'entropy'], 'max_depth': [3, 5, 7, 10]}	{'criterion': 'gini', 'max_depth': 3}	0.7846

จากผลการปรับจูนไฮเปอร์พารามิเตอร์ของแบบจำลอง Decision Tree ในการทดลองครั้งนี้ ได้ทำการประเมินผลโดยใช้พารามิเตอร์หลักสองรายการ ได้แก่ Criterion ซึ่งใช้เพื่อกำหนดเกณฑ์การแยกข้อมูลภายในโหนด (Node Splitting) โดยมีตัวเลือกเป็น 'gini' และ 'entropy' และพารามิเตอร์ Max_Depth ซึ่งควบคุมความลึกสูงสุดของต้นไม้ โดยกำหนดช่วงค่าทดสอบเป็น [3, 5, 7, 10] ผลการปรับจูนพบว่าค่าพารามิเตอร์ที่ให้ประสิทธิภาพดีที่สุด คือ Criterion = 'gini' และ Max_Depth = 3 ซึ่งให้ค่าประสิทธิภาพโดยรวม (F1-Score) สูงสุดร้อยละ 78.46 ผลลัพธ์ดังกล่าวชี้ให้เห็นว่าการควบคุมความลึกของต้นไม้ให้อยู่ในระดับที่เหมาะสม เช่น ความลึกเพียง 3 ชั้น ช่วยลดความซับซ้อนของแบบจำลองและสามารถป้องกันการเกิด Overfitting ได้อย่างมีประสิทธิภาพ

- 3) การปรับจูนไฮเปอร์พารามิเตอร์ที่เหมาะสมด้วยการค้นหาแบบกริดของแบบจำลอง
Neural Network

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 3.8 การปรับแต่งพารามิเตอร์ของแบบจำลอง Neural Network

พารามิเตอร์ที่ปรับแต่ง	พารามิเตอร์ที่เหมาะสม	ค่าประสิทธิภาพโดยรวม
'activation': ['relu'], 'alpha': [0.0005, 0.001], 'hidden_layer_sizes': [(10,), (50,), (100,), (100,50)]	{'activation': 'relu', 'alpha': 0.0005, 'hidden_layer_sizes': (10,)}	0.7684

จากผลการปรับจูนไฮเปอร์พารามิเตอร์ของแบบจำลอง Neural Network ได้กำหนดช่วงค่าของพารามิเตอร์ต่างๆ ได้แก่ Activation ใช้ฟังก์ชันกระตุ้นแบบ 'relu' พารามิเตอร์ Alpha ที่ใช้ควบคุมค่าการเรียนรู้แบบ L2 Regularization โดยกำหนดไว้ที่ 0.0005 และ 0.001 และ Hidden_Layer_Sizes ซึ่งระบุโครงสร้างของจำนวนโน้ดในชั้นซ่อนโดยทดลองทั้งหมด 4 รูปแบบ ได้แก่ (10,), (50,), (100,) และ (100, 50) ผลพบว่าค่าพารามิเตอร์ที่เหมาะสมที่สุดคือ Activation = 'relu', Alpha = 0.0005 และ Hidden_Layer_Sizes = (10,) ซึ่งให้ค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 76.84 ผลลัพธ์ดังกล่าวสะท้อนให้เห็นว่าโครงข่ายประสาทเทียมมีโครงสร้างไม่ซับซ้อนมากนัก (มีชั้นซ่อนเพียง 1 ชั้นและมีเพียง 10 โน้ด) การใช้ค่าการเรียนรู้ที่เหมาะสม โดยไม่จำเป็นต้องเพิ่มความซับซ้อนของแบบจำลอง ซึ่งอาจนำไปสู่การเกิด Overfitting ได้

4) การปรับจูนไฮเปอร์พารามิเตอร์ที่เหมาะสมด้วยการค้นหาแบบกริดของแบบจำลอง LSTM

ตารางที่ 3.9 ผลการปรับแต่งพารามิเตอร์ของแบบจำลอง LSTM

พารามิเตอร์ที่ปรับแต่ง	พารามิเตอร์ที่เหมาะสม	ค่าประสิทธิภาพโดยรวม
'model_units': [64, 128, 256], 'model_activation': ['relu', 'tanh'], 'model_learning_rate': [0.01, 0.005, 0.001], 'epochs': [30], 'batch_size': [16, 32]	{'batch_size': 16, 'epochs': 30, 'model_activation': 'tanh', 'model_learning_rate': 0.001, 'model_units': 256}	0.7776

จากผลการปรับจูนไฮเปอร์พารามิเตอร์ของแบบจำลอง Long Short-Term Memory (LSTM) ได้กำหนดช่วงค่าของพารามิเตอร์สำคัญ ได้แก่ Model_Units ซึ่งระบุจำนวนหน่วยความจำใน LSTM Layer ที่ทดลองตั้งแต่ 64, 128 และ 256 หน่วย พารามิเตอร์ Model_Activation ใช้ฟังก์ชันกระตุ้นแบบ 'relu' และ 'tanh' อัตราการเรียนรู้ (Model_Learning_Rate) ที่กำหนดไว้ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ในช่วง 0.01, 0.005 และ 0.001 จำนวนรอบการฝึก (Epochs) ที่ 30 รอบ และขนาดของแบตช์ (Batch_Size) กำหนดไว้ที่ 16 และ 32 จากการปรับจนพบว่าชุดพารามิเตอร์ที่เหมาะสมที่สุดคือ Model_Units = 256, Model_Activation = 'tanh', Model_Learning_Rate = 0.001 ในส่วน Epochs = 30 และ Batch_Size = 16 ซึ่งให้ค่าประสิทธิภาพโดยรวม (F1-Score) ที่ร้อยละ 77.76 ผลดังกล่าวชี้ให้เห็นว่าการเลือกจำนวนหน่วยความจำที่มากขึ้นร่วมกับการใช้ฟังก์ชันกระตุ้นแบบ Tanh และอัตราการเรียนรู้ต่ำ มีส่วนช่วยให้แบบจำลองสามารถเรียนรู้ลำดับข้อมูลได้อย่างมีประสิทธิภาพ โดยเฉพาะข้อมูลที่มีลักษณะเชิงเวลา ซึ่งเหมาะสมกับโครงสร้างของ LSTM

5) การปรับจูนไฮเปอร์พารามิเตอร์ที่เหมาะสมด้วยการค้นหาแบบกริดของแบบจำลอง Bidirectional LSTM

ตารางที่ 3.10 ผลการปรับแต่งพารามิเตอร์ของแบบจำลอง Bidirectional LSTM

พารามิเตอร์ที่ปรับแต่ง	พารามิเตอร์ที่เหมาะสม	ค่าประสิทธิภาพโดยรวม
'model_units': [64, 128], 'model_activation': ['tanh', 'relu'], 'model_learning_rate': [0.001, 0.0005], 'model_dropout_rate': [0.2, 0.3, 0.4], 'epochs': [30], 'batch_size': [32]	{'batch_size': 32, 'epochs': 30, 'model_activation': 'tanh', 'model_dropout_rate': 0.4, 'model_learning_rate': 0.0005, 'model_units': 128}	0.7788

จากผลการปรับจูนไฮเปอร์พารามิเตอร์ของแบบจำลอง Bidirectional Long Short-Term Memory (BiLSTM) ได้ทำการกำหนดช่วงค่าของพารามิเตอร์หลัก ได้แก่ Model_Units จำนวนหน่วยความจำในแต่ละ LSTM Layer โดยพิจารณา 2 ค่า คือ 64 และ 128, พารามิเตอร์ Model_Activation ที่ใช้ฟังก์ชันกระตุ้นแบบ 'tanh' และ 'relu', Model_Learning_Rate ซึ่งกำหนดค่าไว้ที่ 0.001 และ 0.0005, ค่าการดรอปเอาต์ (Model Dropout Rate) ที่ทดสอบในช่วง 0.2, 0.3 และ 0.4, รวมถึงจำนวนรอบการฝึก (Epochs) กำหนดไว้ที่ 30 และขนาดของ Batch_Size อยู่ที่ 32 ผลการทดลองแสดงให้เห็นว่าพารามิเตอร์ที่ให้ประสิทธิภาพที่ดีที่สุดได้แก่ Model_Units = 128, Model_Activation = 'tanh', Model_Learning_Rate = 0.0005, Model_Dropout_Rate = 0.4, Epochs = 30 และ Batch_Size = 32 ซึ่งให้ค่าประสิทธิภาพโดยรวม (F1-Score) สูงสุดที่ร้อยละ 77.88 จากผลดังกล่าวสามารถสรุปได้ว่าการใช้จำนวนหน่วยความจำที่เหมาะสมร่วมกับอัตราการเรียนรู้ต่ำ และอัตราการดรอปเอาต์ในระดับสูง สามารถช่วยเพิ่มความสามารถในการเรียนรู้ลำดับข้อมูลในสองทิศทางได้อย่างมีประสิทธิภาพ และช่วยลดปัญหา Overfitting ของแบบจำลอง BiLSTM ได้อย่างเหมาะสม

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3.4.3 ขั้นตอนการสร้างแบบจำลอง

ในขั้นตอนการสร้างแบบจำลอง ผู้วิจัยได้ดำเนินการตามกระบวนการที่มีความเป็นระบบ โดยเริ่มจากการแบ่งข้อมูลออกเป็น 2 ชุดหลัก ได้แก่ ชุดข้อมูลฝึกฝน (Training Data) และชุดข้อมูลทดสอบ (Testing Data) โดยได้เลือกใช้วิธีการแบ่งข้อมูลแบบไขว้ (K-Fold Cross Validation) ซึ่งเป็นผลจากการวิเคราะห์ในข้อ 3.4.1 ในการประเมินประสิทธิภาพของแบบจำลอง หลังจากได้พารามิเตอร์ที่เหมาะสมที่สุดจากกระบวนการปรับจูนไฮเปอร์พารามิเตอร์โดยใช้วิธี Grid Search ในข้อ 3.4.2 แล้ว ผู้วิจัยจึงนำค่าพารามิเตอร์เหล่านั้นมาใช้ในการสร้างแบบจำลองจริง โดยมีขั้นตอนดำเนินการที่สามารถสรุปไว้ในตาราง 3.11 ดังนี้

ตารางที่ 3.11 ขั้นตอนการสร้างแบบจำลองของแต่ละเทคนิค

ขั้นตอน	Logistic Regression	Decision Tree	Neural Network	LSTM/BiLSTM
การเลือก Feature	ใช้ Feature ทั้งหมดที่ผ่านการเตรียมข้อมูล			แยกข้อมูลออกเป็น Temporal และ Non-Temporal
ตั้งค่าพารามิเตอร์จาก Best Hyperparameter	C, penalty, solver	max_depth, criterion, min_sample_s_split	hidden_layer_sizes, activation, alpha	units, activation, lr, dropout, batch_size, epochs
การแบ่งข้อมูล	ใช้ StratifiedKFold แบ่งข้อมูล 5 Fold			
การฝึกฝนและทดสอบ (Train-Test)	ฝึกด้วย 4 Folds และทดสอบกับ 1 Fold			Reshape Temporal เป็น [samples, 4, features] แล้วฝึกด้วย Temporal + Non-Temporal Inputs
การบันทึกผลลัพธ์	Accuracy, Precision, Recall, F1-score			
การหาค่าเฉลี่ยของผล	คำนวณค่าเฉลี่ยจากทั้ง 5 Fold โดยเน้น F1-Score และ Recall			
การสรุปผล	รายงานประสิทธิภาพโดยรวมของแบบจำลอง			

จากตาราง 3.11 สามารถอธิบายขั้นตอนการสร้างแบบจำลองโดยสรุปได้ดังนี้ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1) การเลือกใช้คุณลักษณะ (Feature Selection) สำหรับแบบจำลอง Logistic Regression, Decision Tree และ Neural Network ได้ใช้ฟีเจอร์ทั้งหมดที่ผ่านการเตรียมข้อมูลไว้ล่วงหน้า ในขณะที่แบบจำลอง LSTM และ BiLSTM ซึ่งเป็นแบบจำลองที่ออกแบบมาสำหรับข้อมูลอนุกรมเวลา (Time Series) ได้แยกฟีเจอร์ออกเป็น 2 กลุ่ม ได้แก่ (1) ข้อมูลเชิงเวลา (Temporal Features) และ (2) ข้อมูลไม่เชิงเวลา (Non-Temporal Features) เพื่อให้เหมาะสมกับโครงสร้างของแบบจำลอง

2) การตั้งค่าพารามิเตอร์ (Hyperparameter Tuning) ได้ใช้เทคนิค Grid Search ร่วมกับ K-Fold Cross-Validation เพื่อค้นหาค่าพารามิเตอร์ที่เหมาะสมที่สุด โดยพารามิเตอร์ที่ได้รับการปรับแต่งในแต่ละแบบจำลองประกอบด้วย

Logistic Regression: ค่าคงที่ (C) ฟังก์ชันลงโทษ (Penalty) และอัลกอริทึม Solver

Decision Tree: เกณฑ์การแบ่งข้อมูล (Criterion) ความลึกของต้นไม้ (Max_Depth) และจำนวนตัวอย่างขั้นต่ำในการแบ่ง (Min_Samples_Split)

Neural Network: ขนาดของชั้นซ่อน (Hidden_Layer_Sizes) ฟังก์ชันกระตุ้น (Activation) และค่า Regularization Alpha

LSTM / BiLSTM: จำนวนหน่วยความจำ (Units) ฟังก์ชันกระตุ้น (Activation) อัตราการเรียนรู้ (Learning_Rate) ค่าการดรอปเอาต์ (Dropout) ขนาดแบตช์ (Batch_Size) และจำนวนรอบการฝึก (Epochs)

3) การแบ่งข้อมูล (Data Splitting: K-Fold Cross-Validation) แบบจำลองทั้งหมดใช้เทคนิค Stratified K-Fold Cross-Validation เพื่อให้แต่ละชุดย่อย (Fold) มีการกระจายตัวของข้อมูลในแต่ละคลาสอย่างสมดุลกัน ซึ่งวิธีนี้ช่วยลดความลำเอียง (Bias) และเพิ่มความน่าเชื่อถือในการประเมินประสิทธิภาพของแบบจำลอง

4) การฝึกฝนและทดสอบแบบจำลอง (Train-Test Procedure) ในแต่ละรอบของ K-Fold ได้ใช้ข้อมูล 4 Folds สำหรับการฝึก (Training) และ 1 Fold สำหรับการทดสอบ (Testing) ในกรณีของ LSTM และ BiLSTM ได้มีการจัดโครงสร้างข้อมูลเชิงเวลาให้อยู่ในรูป [Samples, 4, Features] เพื่อให้สอดคล้องกับการประมวลผลแบบลำดับ

5) การประเมินผลในแต่ละรอบ (Per-Fold Evaluation) ในแต่ละรอบของกระบวนการฝึกฝนและทดสอบ มีการประเมินประสิทธิภาพของแบบจำลองด้วยตัวชี้วัดที่สำคัญ 4 ด้าน ได้แก่ Accuracy, Precision, Recall และ F1-score เพื่อสะท้อนประสิทธิภาพของแบบจำลองในด้านต่างๆ

6) การคำนวณค่าเฉลี่ยผลลัพธ์ (Average Evaluation) หลังจากสิ้นสุดการประเมินครบทั้ง 5 Folds แล้ว จะทำการคำนวณค่าเฉลี่ยของผลลัพธ์ทั้งหมดโดยเฉพาะให้ความสำคัญกับค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall)

7) การสรุปผล (Final Model Evaluation) ผลลัพธ์จากการประเมินทั้งหมดจะถูก

รวบรวมและสรุปเป็นค่าประสิทธิภาพโดยรวมของแต่ละแบบจำลอง ซึ่งจะถูกใช้เปรียบเทียบเพื่อเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปเผยแพร่โดยไม่ได้รับอนุญาต
ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

คัดเลือกแบบจำลองที่มีความเหมาะสมที่สุดในการนำไปใช้กับการทำนายการผิคนัดชำระหนี้ในบริษัทของการวิจัยครั้งนี้

3.5 การประเมินผล (Evaluation)

การประเมินประสิทธิภาพของแบบจำลองเป็นขั้นตอนสำคัญที่ช่วยให้ผู้วิจัยสามารถวัดและเปรียบเทียบศักยภาพของแต่ละแบบจำลองได้อย่างเป็นระบบ ในขั้นตอนนี้จะใช้ตัวชี้วัดมาตรฐานสำหรับงานจำแนกประเภท (Classification) ได้แก่ เมทริกซ์ความสับสน (Confusion Matrix) และจำแนกประเภท (Classification Report) เพื่อวิเคราะห์ผลการทำนายทั้งในชุดข้อมูลฝึกฝน (Train Set) และชุดข้อมูลทดสอบ (Test Set)

3.5.1 เมทริกซ์ความสับสน (Confusion Matrix)

Confusion Matrix เป็นเครื่องมือพื้นฐานที่ใช้ประเมินผลการจำแนกประเภทของแบบจำลอง โดยแสดงให้เห็นว่าการทำนายแต่ละกลุ่มเป้าหมายมีความถูกต้องหรือผิดพลาดมากน้อยเพียงใด ช่วยให้สามารถวิเคราะห์ข้อดีข้อเสียของแบบจำลองในเชิงลึก เช่น แบบจำลองมีแนวโน้มทำนายกลุ่มเสี่ยงได้แม่นยำหรือไม่ ผู้วิจัยสามารถสร้าง Confusion Matrix ได้ด้วยฟังก์ชันจากไลบรารี scikit-learn และแสดงผลด้วย heatmap จากไลบรารี seaborn เพื่อให้อ่านค่าได้ง่าย ดังตัวอย่างโค้ด

```
from sklearn.metrics import confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

cm_train = confusion_matrix(y_train, y_pred_train)
sns.heatmap(cm_train, annot=True, fmt='d', cmap='Oranges')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()
```

รูปที่ 3.4 การประเมินผลด้วยเมทริกซ์ความสับสน (Confusion Matrix) ของแบบจำลอง

3.5.2 การจำแนกประเภท (Classification Report)

Classification Report เป็นการสรุปผลการทำนายของแบบจำลองอย่างละเอียด โดยนำเสนอค่าตัวชี้วัดหลัก ได้แก่ Precision, Recall, F1-Score และ Accuracy สำหรับแต่ละกลุ่มเป้าหมายรายงานนี้ช่วยให้เห็นประสิทธิภาพเชิงลึกของแบบจำลอง โดยเฉพาะในกรณีข้อมูลไม่สมดุล (Imbalanced Data) สามารถสร้าง Classification Report ได้โดยใช้ฟังก์ชันจากไลบรารี Scikit-Learn ดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

```

from sklearn.metrics import classification_report

# Logistic Regression
print("Model Logistic Regression")
print(classification_report(y_train, y_pred_train_lr))

# Decision Tree
print("Model Decision Tree")
print(classification_report(y_train, y_pred_train_dt))

# Neural Network
print("Model Neural Network")
print(classification_report(y_train, y_pred_train_nn))

# Long Short-Term Memory (LSTM)
print("Model LSTM")
print(classification_report(y_train, y_pred_train_lstm))

# Bidirectional LSTM
print("Model Bidirectional LSTM")
print(classification_report(y_train, y_pred_train_bilstm))

```

รูปที่ 3.5 การประเมินผลข้อมูลโดยการวิเคราะห์ความแม่นยำของแบบจำลอง

3.6 การนำแบบจำลองไปใช้งาน (Deployment)

หลังจากที่ได้แบบจำลองที่มีประสิทธิภาพสูงสุดจากการประเมินผลในขั้นตอนก่อนหน้านี้ ขั้นตอนสุดท้ายคือการนำแบบจำลองดังกล่าวไปประยุกต์ใช้กับข้อมูลจริงในทางปฏิบัติ เพื่อสนับสนุนกระบวนการตัดสินใจขององค์กรด้านการบริหารความเสี่ยงสินเชื่อรายย่อย (Nano Finance) โดยแบบจำลองที่พัฒนาขึ้นนี้ สามารถนำไปใช้สำหรับคาดการณ์โอกาสการผิดนัดชำระหนี้ของลูกค้าในเดือนถัดไป ซึ่งช่วยให้สถาบันการเงินสามารถติดตามกลุ่มลูกค้าเสี่ยง วางแผนกลยุทธ์ในการบริหารจัดการหนี้ หรือออกแบบมาตรการช่วยเหลือลูกหนี้ได้อย่างทันที่

นอกจากนี้ ยังได้เสนอแนะแนวทางการพัฒนาและปรับปรุงแบบจำลองอย่างต่อเนื่อง เพื่อให้สอดคล้องกับข้อมูลและพฤติกรรมของลูกค้าที่เปลี่ยนแปลงไปในอนาคต ทั้งนี้ การนำแบบจำลองไปใช้จริงจะต้องคำนึงถึงความเหมาะสมของเทคโนโลยี การจัดเก็บข้อมูล การรักษาความปลอดภัยของข้อมูลส่วนบุคคล และการสร้างความเข้าใจให้กับผู้ปฏิบัติงานที่เกี่ยวข้อง

สำหรับงานวิจัยนี้ยังไม่ครอบคลุมถึงขั้นตอนการนำแบบจำลองไปใช้งานจริง แต่อยู่ในระหว่างการวางแผนสำหรับการประยุกต์ใช้แบบจำลองในกระบวนการปฏิบัติงานต่อไป

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 4

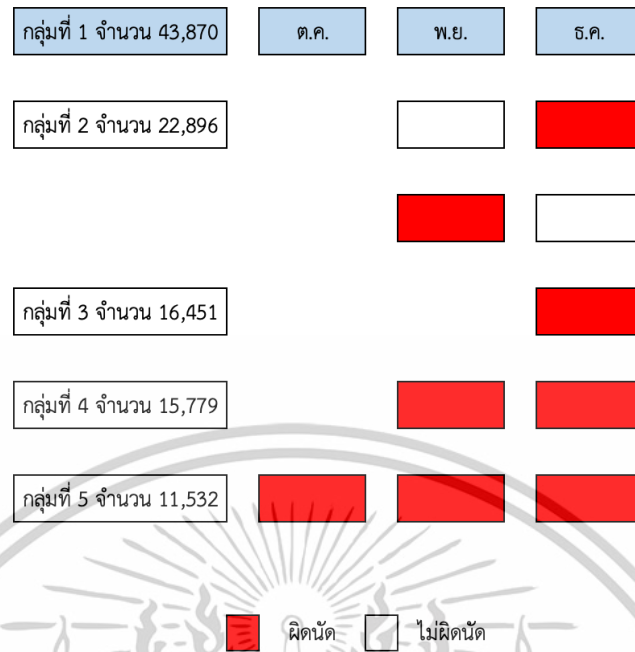
ผลการวิจัยและอภิปรายผล

บทนี้นำเสนอผลการวิเคราะห์การทำนายการผัดขันธ์ชำระหนี้ของลูกค้าสินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพ โดยใช้เทคนิคการเรียนรู้ของเครื่องและการเรียนรู้เชิงลึก 5 แบบ ได้แก่ Logistic Regression, Decision Tree, Neural Network, Long Short-Term Memory (LSTM) และ Bidirectional LSTM เพื่อพัฒนาแบบจำลองสำหรับประเมินความเสี่ยงการผัดขันธ์ชำระหนี้ในเดือนถัดไปอย่างแม่นยำ โดยเฉพาะการเลือกใช้ Bidirectional LSTM อ้างอิงจากงานของ Ala'raj, Abbod, & Majdalawieh. (2021) ซึ่งชี้ให้เห็นศักยภาพของแบบจำลองนี้ในการวิเคราะห์ข้อมูลเชิงลำดับ (Sequential Data) ที่สามารถเรียนรู้ความสัมพันธ์ของข้อมูลทั้งอดีตและอนาคต ส่งผลให้สามารถระบุรูปแบบพฤติกรรมทางการเงินที่ซับซ้อนและต่อเนื่องได้ดีกว่าแบบจำลองแบบดั้งเดิมสำหรับเกณฑ์การประเมินประสิทธิภาพของแบบจำลอง ผู้วิจัยใช้ค่าประสิทธิภาพโดยรวม (F1-Score) ร่วมกับค่าความระลึก (Recall) เพื่อสะท้อนถึงความแม่นยำและความสามารถในการระบุลูกค้าที่มีความเสี่ยงสูงต่อการผัดขันธ์ชำระหนี้

ข้อมูลที่ใช้ประกอบด้วยบัญชีสินเชื่อรายย่อยที่มีสถานะเคลื่อนไหว ณ เดือนมกราคม พ.ศ. 2567 รวม 43,870 บัญชี โดยแบ่งเป็นชุดฝึก (Train Set) และชุดทดสอบ (Test Set) เพื่อการพัฒนาและประเมินแบบจำลอง โดยได้ทดลองเปรียบเทียบการแบ่งข้อมูลสองวิธี คือ วิธีการแบ่งแบบแยกชุด (Hold-Out: 80% Training, 20% Testing) และวิธีการแบ่งแบบไขว้ (5-Fold Cross Validation) เพื่อเลือกวิธีที่เหมาะสมที่สุด

เพื่อประเมินศักยภาพของแบบจำลองภายใต้บริบทความเสี่ยงที่แตกต่างกัน งานวิจัยนี้ได้จำแนกกลุ่มลูกค้าออกเป็น 5 กลุ่มตามลักษณะประวัติการผัดขันธ์ชำระหนี้ย้อนหลัง ได้แก่ กลุ่มที่ 1 (กลุ่มลูกค้าทั่วไป) กลุ่มที่ 2 (กลุ่มความเสี่ยงต่ำ) กลุ่มที่ 3 (กลุ่มความเสี่ยงปานกลาง) กลุ่มที่ 4 (กลุ่มความเสี่ยงสูง) และกลุ่มที่ 5 (กลุ่มความเสี่ยงวิกฤต) ดังแสดงในรูปที่ 4.1 ซึ่งแสดงวิธีการจำแนกกลุ่มลูกค้าโดยใช้พฤติกรรมผัดขันธ์ชำระหนี้ในลำดับเวลาเป็นเกณฑ์ โดยสัญลักษณ์สีแดงในแต่ละเดือนแทนการผัดขันธ์ชำระหนี้ และสีขาวแทนการชำระปกติ การแบ่งกลุ่มดังกล่าวช่วยให้สามารถวิเคราะห์และเปรียบเทียบประสิทธิภาพของแบบจำลองในการรับมือกับความเสี่ยงที่หลากหลายได้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.1 วิธีการแบ่งกลุ่มลูกค้า

นอกจากนี้ ได้ดำเนินการวิเคราะห์ความสำคัญของตัวแปร (Feature Importance) ในกลุ่มลูกค้าหลัก เพื่อระบุปัจจัยสำคัญที่ส่งผลต่อความเสี่ยงผิดนัดชำระหนี้โดยโครงสร้างเนื้อหาของบทรนี้ประกอบด้วย 8 หัวข้อหลัก ดังนี้

1. สถิติเชิงพรรณนา (Descriptive Statistics)
2. ผลการคัดเลือกการแบ่งข้อมูลชุดทดสอบด้วยวิธีแบบแยกชุด (Hold-out) และวิธีแบบไขว้ (K-Fold Cross Validation)
3. ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 1 ลูกค้าทั่วไป (General Customer)
4. ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 2 ความเสี่ยงต่ำ (Low Risk)
5. ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 3 ความเสี่ยงปานกลาง (Moderate Risk)
6. ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 4 ความเสี่ยงสูง (High Risk)
7. ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 5 ความเสี่ยงวิกฤต (Critical Risk)
8. อภิปรายผล

การนำเสนอผลการวิจัยในหัวข้อต่างๆ ตามโครงสร้างข้างต้น ผู้วิจัยได้วางแนวทางการนำเสนอผลการวิเคราะห์โดยแบ่งออกเป็น 2 ลักษณะหลัก ดังนี้ 1) กลุ่มที่ 1 ซึ่งเป็นกลุ่มลูกค้าทั่วไป (General Customer) จะมีการนำเสนอรายละเอียดเชิงลึกในทุกมิติของการประเมินผลแบบจำลอง ได้แก่ ประสิทธิภาพการทำนายของแต่ละเทคนิค สมมติฐานทางสถิติ ตลอดจนการวิเคราะห์

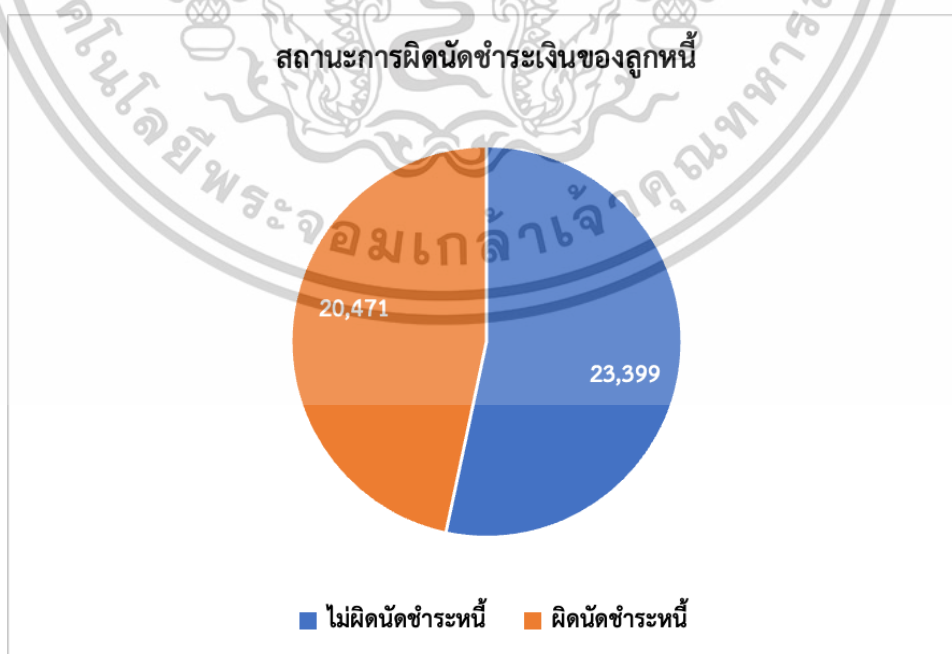
ความสำคัญของตัวแปร (Feature Importance) เพื่อสะท้อนศักยภาพของแต่ละแบบจำลองในบริบทเอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ของลูกค้ำทั้งหมด ซึ่งจะช่วยให้อาจมองเห็นจุดเด่น จุดด้อย และความแตกต่างระหว่างเทคนิคได้อย่างครอบคลุม และ 2) กลุ่มที่ 2 ถึงกลุ่มที่ 5 (กลุ่มความเสี่ยงต่ำ กลุ่มความเสี่ยงปานกลาง กลุ่มความเสี่ยงสูง และกลุ่มความเสี่ยงวิกฤต) ผู้วิจัยจะมุ่งเน้นการเปรียบเทียบประสิทธิภาพของแต่ละเทคนิค (Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM) โดยใช้ค่า F1-Score และ Recall เป็นหลักสำหรับแต่ละกลุ่มความเสี่ยงเท่านั้น ทั้งนี้เพื่อเป็นการนำเสนอมีความกระชับและเน้นการวิเคราะห์เชิงเปรียบเทียบข้ามกลุ่มความเสี่ยง เพื่อแสดงศักยภาพของแบบจำลองในแต่ละกลุ่มและสนับสนุนการเลือกใช้เทคนิคที่เหมาะสมในเชิงนโยบายการกำหนดแนวทางดังกล่าวมีวัตถุประสงค์เพื่อให้การนำเสนอผลการวิเคราะห์มีความเหมาะสมสอดคล้องกับบริบทของแต่ละกลุ่มลูกค้ำ และเอื้อต่อการนำผลลัพธ์ไปใช้ประโยชน์ในการบริหารความเสี่ยงและกำหนดกลยุทธ์การบริหารสินเชื่อในทางปฏิบัติ

4.1 สถิติเชิงพรรณนา (Descriptive Statistics)

งานวิจัยนี้ดำเนินการโดยใช้ชุดข้อมูลจากสถาบันการเงินแห่งหนึ่งในประเทศไทย ซึ่งผู้วิจัยได้ทำการคัดเลือกมาจากรฐานข้อมูล โดยอิงจากเงื่อนไขที่กำหนด คือ ต้องเป็นข้อมูลเฉพาะบัญชีสินเชื่อที่ยังมีการเคลื่อนไหวในเดือนมกราคม พ.ศ. 2567 เท่านั้น โดยมีข้อมูลทั้งหมด 43,870 บัญชี

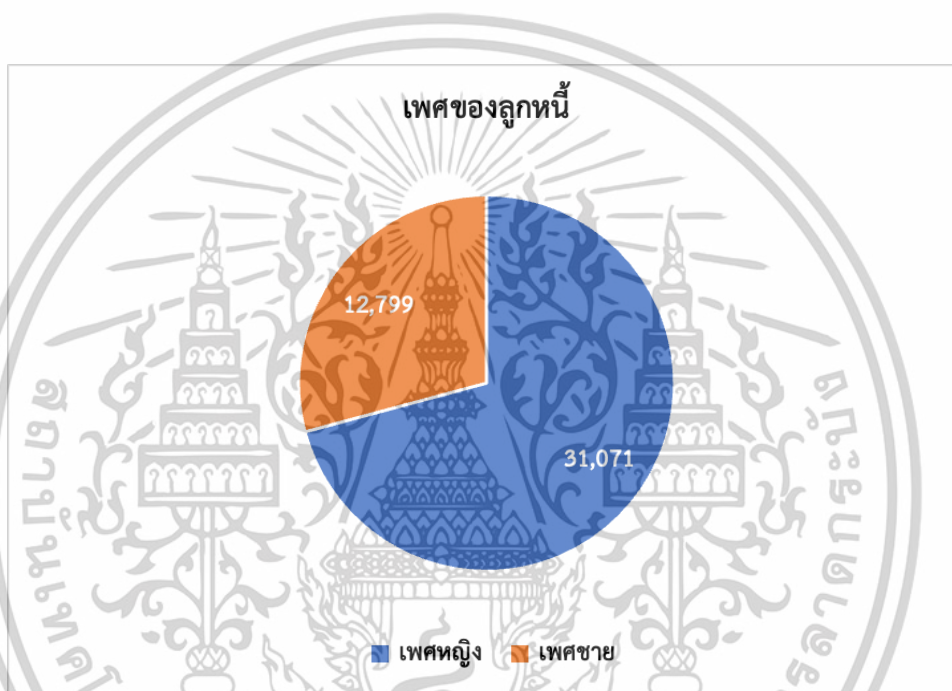
ตัวแปรอิสระ ได้มาจากสองส่วนหลัก ได้แก่ ข้อมูลด้านลักษณะสินเชื่อของลูกค้าหนี้ที่บัญชียังมีการเคลื่อนไหวในเดือนมกราคม พ.ศ. 2567 และ ข้อมูลด้านพฤติกรรมของลูกค้าหนี้ ตั้งแต่เดือนมกราคม พ.ศ. 2566 ถึงเดือนธันวาคม พ.ศ. 2566 และในส่วนของตัวแปรตาม คือ สถานะการผิดนัดชำระของลูกค้าหนี้ในเดือน มกราคม พ.ศ. 2567



รูปที่ 4.2 สถานะผิดนัดชำระเงินของลูกค้าหนี้ ณ เดือน มกราคม พ.ศ. 2567

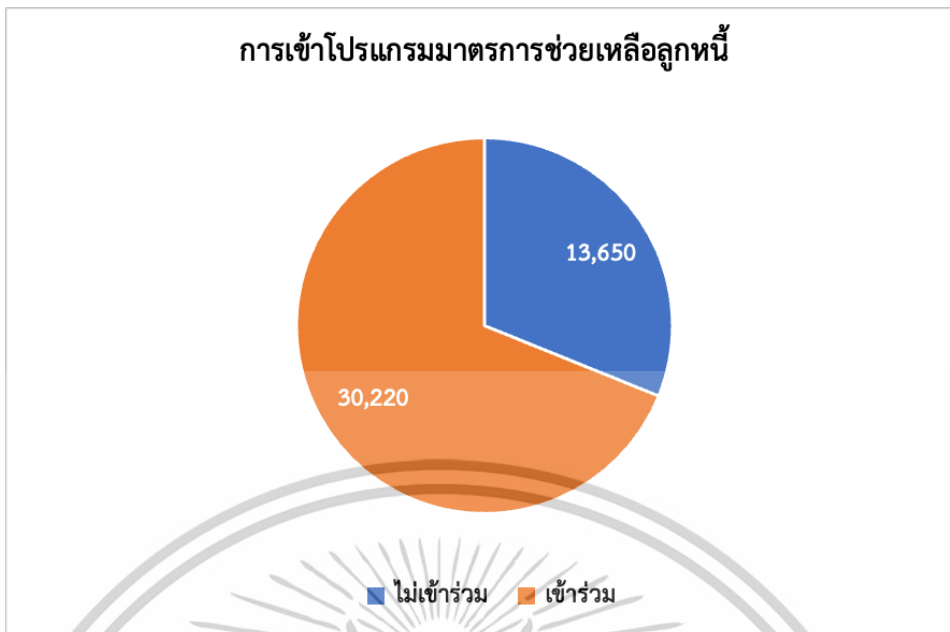
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากรูปที่ 4.2 แสดงสถานการณ์ผิคนัดชำระเงินของลูกหนี้ ณ เดือน มกราคม พ.ศ. 2567 โดยแบ่งออกเป็น 2 กลุ่มได้แก่ กลุ่มบัญชีลูกหนี้ที่ไม่ผิคนัดชำระ จำนวน 23,399 บัญชี คิดเป็นร้อยละ 53.33 ของบัญชีลูกหนี้ทั้งหมด และกลุ่มบัญชีลูกหนี้ที่ผิคนัดชำระ จำนวน 20,471 บัญชี คิดเป็นร้อยละ 46.67 ของบัญชีลูกหนี้ทั้งหมด ข้อมูลดังกล่าวแสดงให้เห็นว่าลูกหนี้ส่วนใหญ่ยังสามารถชำระหนี้ได้ แต่สัดส่วนของกลุ่มลูกหนี้ที่ผิคนัดชำระยังคงอยู่ในระดับที่ค่อนข้างสูงเกือบครึ่งหนึ่งของบัญชีทั้งหมด ซึ่งอาจเป็นสัญญาณถึงความเปราะบางด้านความสามารถในการชำระหนี้ โดยเฉพาะอย่างยิ่งในสินเชื่อรายย่อยเพื่อผู้ประกอบการอาชีพ ที่มักปล่อยกู้ให้กับกลุ่มผู้มีรายได้น้อยและไม่แน่นอนและไม่มีหลักประกัน



รูปที่ 4.3 เพศของลูกหนี้รายบัญชี

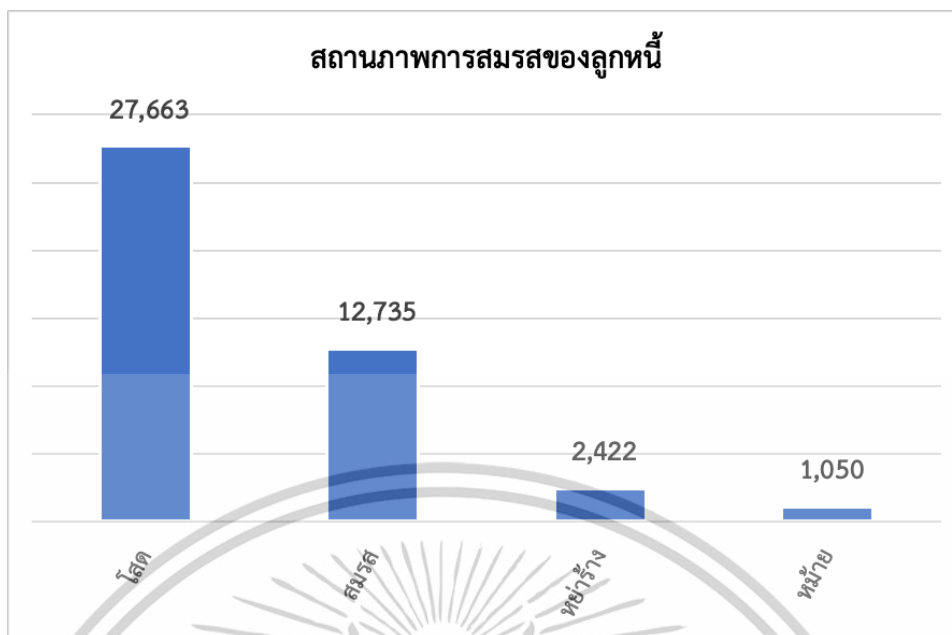
จากรูปที่ 4.3 แสดงข้อมูลจำแนกตามเพศของลูกหนี้รายบัญชี โดยแบ่งออกเป็น 2 กลุ่ม ได้แก่ กลุ่มเพศหญิง จำนวน 31,071 บัญชี คิดเป็นร้อยละ 70.83 ของบัญชีลูกหนี้ทั้งหมด และกลุ่มเพศชาย จำนวน 12,799 บัญชี คิดเป็นร้อยละ 29.17 ของบัญชีลูกหนี้ทั้งหมด จากข้อมูลดังกล่าวแสดงให้เห็นว่า กลุ่มบัญชีลูกหนี้ส่วนใหญ่เป็นเพศหญิง ซึ่งมีจำนวนมากกว่ากลุ่มเพศชายอย่างชัดเจน ซึ่งอาจสะท้อนถึงโครงสร้างประชากรของประเทศไทยที่มีสัดส่วนเพศหญิงมากกว่าเพศชาย จึงเป็นเหตุให้บัญชีสินเชื่อในกลุ่มตัวอย่างส่วนใหญ่เป็นเพศหญิง



รูปที่ 4.4 การเข้าโปรแกรมมาตรการช่วยเหลือลูกหนี้รายบัญชี

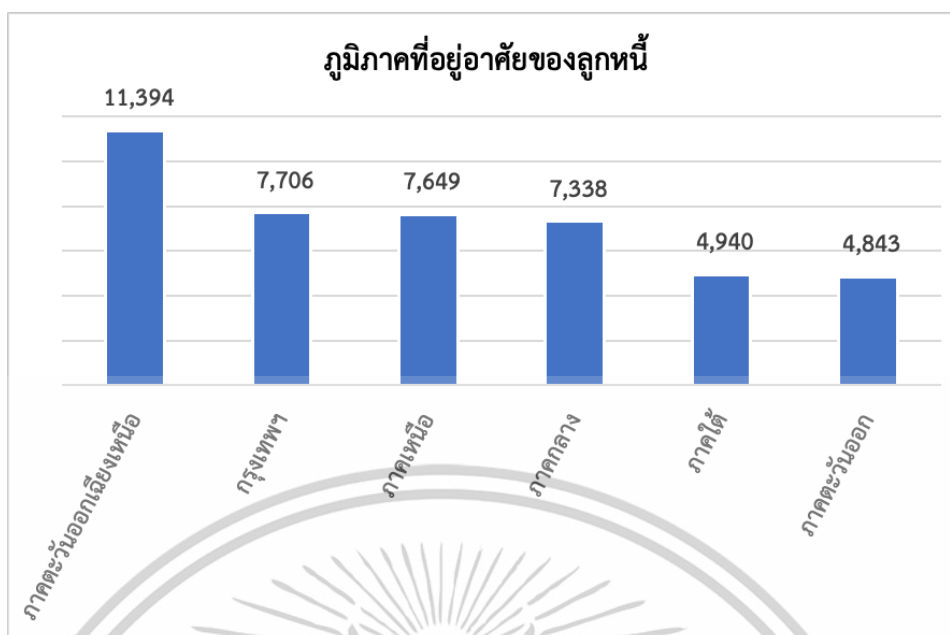
จากรูปที่ 4.4 แสดงการเข้าโปรแกรมมาตรการช่วยเหลือลูกหนี้รายบัญชี โดยแบ่งออกเป็น 2 กลุ่ม ได้แก่ กลุ่มที่ เข้าร่วมโปรแกรม และกลุ่มที่ ไม่ได้เข้าร่วมโปรแกรม พบว่ากลุ่มที่เข้าร่วมโปรแกรม มีจำนวน 30,220 บัญชี คิดเป็น ร้อยละ 68.89 ของบัญชีลูกหนี้ทั้งหมด ขณะที่กลุ่มที่ไม่ได้เข้าร่วมมีจำนวน 13,650 บัญชี คิดเป็น ร้อยละ 31.11 จากข้อมูลดังกล่าวแสดงให้เห็นว่าลูกหนี้ส่วนใหญ่มีความจำเป็นต้องพึ่งพามาตรการช่วยเหลือจากสถาบันการเงิน ซึ่งอาจเป็นผลมาจากสถานการณ์ทางการเงินที่เปราะบาง เช่น รายได้ไม่แน่นอน ภาระหนี้สูง หรือได้รับผลกระทบจากภาวะเศรษฐกิจ โดยการเข้าร่วมโปรแกรมช่วยเหลือเหล่านี้มีบทบาทสำคัญในการบรรเทาภาระหนี้ และอาจมีผลต่อพฤติกรรมการผัดนัดชำระในระยะยาว

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.5 สถานภาพการสมรสของลูกหนี้รายบัญชี

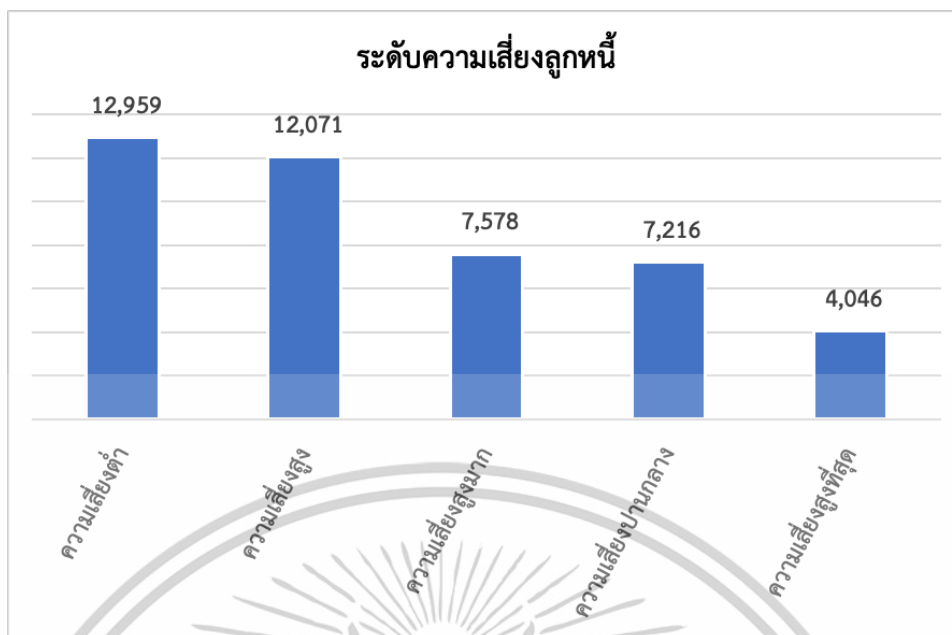
จากรูปที่ 4.5 แสดงสถานภาพการสมรสของลูกหนี้รายบัญชี โดยจำแนกออกเป็น 4 กลุ่ม ได้แก่ กลุ่มโสด กลุ่มสมรส กลุ่มหย่าร้าง และกลุ่มหม้าย โดยพบว่ากลุ่มที่มีจำนวนมากที่สุดคือกลุ่มลูกหนี้ที่มีสถานภาพโสด จำนวน 27,663 บัญชี คิดเป็นร้อยละ 63.06 ของบัญชีลูกหนี้ทั้งหมด รองลงมาคือกลุ่มที่สมรส จำนวน 12,735 บัญชี คิดเป็นร้อยละ 29.03 ของบัญชีลูกหนี้ทั้งหมด ตามด้วยกลุ่มหย่าร้าง จำนวน 2,422 บัญชี คิดเป็นร้อยละ 5.52 ของบัญชีลูกหนี้ทั้งหมด และกลุ่มหม้าย จำนวน 1,050 บัญชี คิดเป็นร้อยละ 2.39 ของบัญชีลูกหนี้ทั้งหมด ตามลำดับ จากข้อมูลดังกล่าวแสดงให้เห็นว่าลูกหนี้ส่วนใหญ่อยู่ในสถานภาพโสด กลุ่มที่มีสถานภาพหย่าร้างหรือหม้าย แม้จะมีสัดส่วนน้อย แต่ก็อาจมีความเสี่ยงด้านรายได้หรือความสามารถในการชำระหนี้ที่แตกต่างกัน



รูปที่ 4.6 ภูมิภาคที่อยู่อาศัยของบัญชีลูกหนี้รายบัญชี

จากรูปที่ 4.6 แสดงภูมิภาคที่อยู่อาศัยของลูกหนี้รายบัญชี แบ่งออกเป็น 6 ภูมิภาค ได้แก่ ภาคตะวันออกเฉียงเหนือ จำนวน 11,394 บัญชี คิดเป็นร้อยละ 25.97 ของบัญชีลูกหนี้ทั้งหมด รองลงมาคือกรุงเทพมหานคร จำนวน 7,706 บัญชี คิดเป็นร้อยละ 17.57 ของบัญชีลูกหนี้ทั้งหมด อันดับสามคือภาคเหนือ จำนวน 7,649 บัญชี คิดเป็นร้อยละ 17.44 ของบัญชีลูกหนี้ทั้งหมด อันดับถัดมาคือภาคกลาง จำนวน 7,338 บัญชี คิดเป็นร้อยละ 16.73 ของบัญชีลูกหนี้ทั้งหมด ตามด้วยภาคใต้ จำนวน 4,940 บัญชี คิดเป็นร้อยละ 11.26 และภาคตะวันออก จำนวน 4,843 บัญชี คิดเป็นร้อยละ 11.04 ของบัญชีลูกหนี้ทั้งหมด ตามลำดับ จากข้อมูลดังกล่าวสะท้อนให้เห็นว่า ลูกหนี้ส่วนใหญ่อาศัยอยู่ในภูมิภาคภาคเหนือและภาคตะวันออกเฉียงเหนือ ซึ่งมีสัดส่วนมากที่สุดในกลุ่มตัวอย่าง ขณะที่ภาคตะวันออกและภาคใต้มีสัดส่วนน้อยที่สุด อย่างไรก็ตาม ความแตกต่างของจำนวนบัญชีในแต่ละภูมิภาคอาจสะท้อนถึงปัจจัยด้านประชากร ความต้องการใช้สินเชื่อ และการเข้าถึงแหล่งเงินทุนในแต่ละพื้นที่

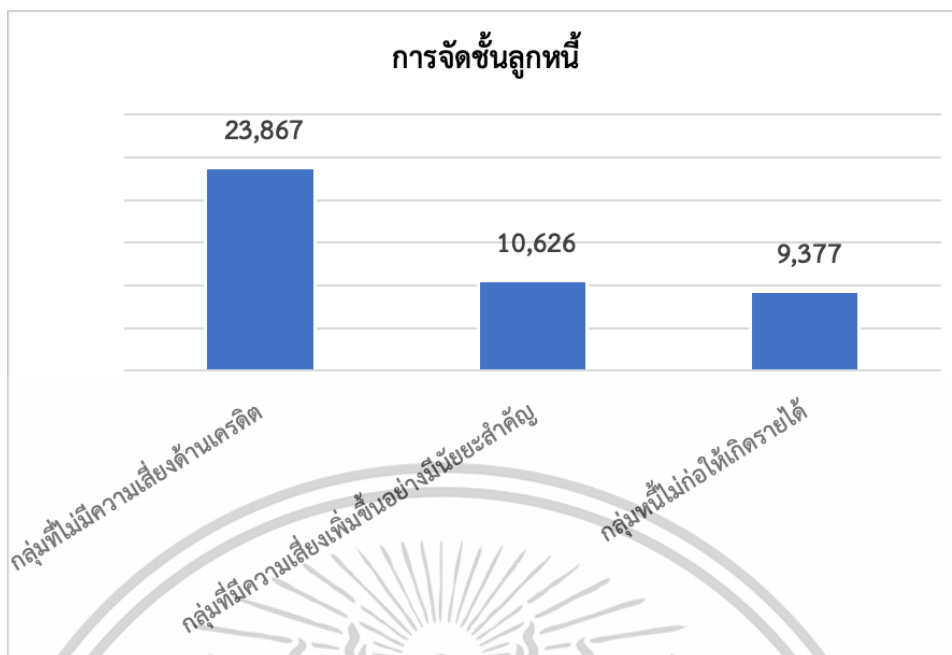
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.7 ระดับความเสี่ยงของลูกหนี้รายบัญชี

จากรูปที่ 4.7 แสดงระดับความเสี่ยงของลูกหนี้รายบัญชี ที่จัดระดับโดยสถาบันการเงิน ซึ่งแบ่งออกเป็น 5 ระดับ ได้แก่ ความเสี่ยงต่ำ ความเสี่ยงปานกลาง ความเสี่ยงสูง ความเสี่ยงสูงมาก และ ความเสี่ยงสูงที่สุด โดยพบว่า กลุ่มลูกหนี้ที่มีความเสี่ยงอยู่ในระดับต่ำมีจำนวนมากที่สุด คือ 12,959 บัญชี คิดเป็นร้อยละ 29.54 ของบัญชีทั้งหมด รองลงมาคือกลุ่มที่มีความเสี่ยงสูง จำนวน 12,071 บัญชี คิดเป็นร้อยละ 27.52 และกลุ่มความเสี่ยงสูงมาก จำนวน 7,578 บัญชี คิดเป็นร้อยละ 17.27 ตามด้วยกลุ่มความเสี่ยงปานกลาง จำนวน 7,216 บัญชี คิดเป็นร้อยละ 16.45 ขณะที่กลุ่มที่มีความเสี่ยงสูงที่สุดมีจำนวนบัญชีต่ำที่สุด คือ 4,046 บัญชี คิดเป็นร้อยละ 9.23 ของบัญชีทั้งหมด จากข้อมูลดังกล่าวสะท้อนให้เห็นถึงการกระจายระดับความเสี่ยงของลูกหนี้ในภาพรวม โดยแม้ว่ากลุ่มความเสี่ยงต่ำจะมีจำนวนมากที่สุด แต่กลุ่มที่มีระดับความเสี่ยงปานกลางถึงสูงก็ยังมีสัดส่วนรวมกันเกินครึ่งหนึ่งของข้อมูลทั้งหมด ซึ่งอาจส่งผลต่อการวางแผนบริหารความเสี่ยงด้านสินเชื่อของสถาบันการเงินในเชิงกลยุทธ์ ทั้งในด้านการพิจารณาอนุมัติสินเชื่อ การกำหนดเงื่อนไขในการปล่อยกู้ รวมถึงการวิเคราะห์แนวโน้มการผิดนัดชำระหนี้ในกลุ่มความเสี่ยงที่แตกต่างกัน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.8 ชั้นหนี้ของลูกหนี้รายบัญชี

จากรูปที่ 4.8 แสดงหนี้ของลูกหนี้รายบัญชี โดยแบ่งออกเป็น 3 กลุ่ม ได้แก่ กลุ่มหนี้ที่ไม่ก่อให้เกิดรายได้ กลุ่มที่มีความเสี่ยงเพิ่มขึ้นอย่างมีนัยยะสำคัญ และกลุ่มที่ไม่มีความเสี่ยงด้านเครดิต โดยพบว่า กลุ่มที่ไม่มีความเสี่ยงด้านเครดิตมีจำนวนบัญชีมากที่สุด คือ 23,867 บัญชี คิดเป็นร้อยละ 54.52 ของบัญชีทั้งหมด รองลงมาคือกลุ่มที่มีความเสี่ยงเพิ่มขึ้นอย่างมีนัยยะสำคัญ จำนวน 10,626 บัญชี คิดเป็นร้อยละ 24.27 ของบัญชีทั้งหมด และกลุ่มหนี้ที่ไม่ก่อให้เกิดรายได้มีจำนวน 9,377 บัญชี คิดเป็นร้อยละ 21.41 ของบัญชีทั้งหมด จากข้อมูลดังกล่าวสะท้อนให้เห็นว่ากลุ่มลูกหนี้ส่วนใหญ่ยังคงอยู่ในสถานะที่ไม่มีความเสี่ยงด้านเครดิต ซึ่งเป็นสัญญาณเชิงบวกต่อเสถียรภาพของพอร์ตสินเชื่อโดยรวม อย่างไรก็ตาม ยังมีกลุ่มลูกหนี้ที่มีแนวโน้มความเสี่ยงเพิ่มขึ้นและกลุ่มหนี้เสียในระดับที่ไม่อาจมองข้ามได้ เนื่องจากอาจส่งผลกระทบต่อคุณภาพของสินทรัพย์ในอนาคต

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.1 สถิติพรรณนาของตัวแปรเชิงปริมาณ

ชื่อตัวแปร	ค่าเฉลี่ย	ค่าต่ำสุด	ค่าสูงสุด
อายุ (ปี)	44.3299	21.0000	67.0000
อัตราส่วนความสามารถในการชำระหนี้	12.7871	0.0000	624.2447
จำนวนบุตร (คน)	1.3822	0.0000	10.0000
จำนวนเดือนนับจากวันทำสินเชื่อ (เดือน)	30.1656	12.0000	133.0000
ระยะเวลากู้ยืม (เดือน)	12.4577	12.0000	156.0000
อัตราดอกเบี้ยของเงินกู้	0.3183	0.2400	0.3300
จำนวนวันค้างชำระ (วัน)	61.9211	0.0000	1076.0000
อัตราส่วนยอดเงินกู้คงเหลือต่อจำนวนเงินกู้	0.8631	-0.0386	1.5162
อัตราส่วนค้างงวดต่อยอดหนี้	0.0699	0.0000	0.9474
ค่าเฉลี่ยของอัตราส่วนการเบิกใช้วงเงิน ใน 3 เดือน	0.8666	-0.0385	1.5162
จำนวนครั้งที่จ่ายเต็ม ใน 12 เดือนล่าสุด (ครั้ง)	4.7833	0.0000	12.0000
จำนวนครั้งที่จ่ายเต็ม ใน 9 เดือนล่าสุด (ครั้ง)	3.5594	0.0000	9.0000
จำนวนครั้งที่จ่ายเต็ม ใน 6 เดือนล่าสุด (ครั้ง)	2.3163	0.0000	6.0000
จำนวนครั้งที่จ่ายเต็ม ใน 3 เดือนล่าสุด (ครั้ง)	1.1509	0.0000	3.0000
จำนวนครั้งที่จ่ายบางส่วน ใน 12 เดือนล่าสุด (ครั้ง)	5.0616	0.0000	12.0000
จำนวนครั้งที่จ่ายบางส่วน ใน 9 เดือนล่าสุด (ครั้ง)	4.2968	0.0000	9.0000
จำนวนครั้งที่จ่ายบางส่วน ใน 6 เดือนล่าสุด (ครั้ง)	3.2612	0.0000	6.0000
จำนวนครั้งที่จ่ายบางส่วน ใน 3 เดือนล่าสุด (ครั้ง)	1.7783	0.0000	3.0000
จำนวนงวดที่ค้างชำระมากที่สุด ใน 3 เดือนล่าสุด (ครั้ง)	2.0038	0.0000	5.0000
จำนวนงวดที่ค้างชำระมากที่สุด ใน 6 เดือนล่าสุด (ครั้ง)	2.1009	0.0000	5.0000
จำนวนงวดที่ค้างชำระมากที่สุด ใน 9 เดือนล่าสุด (ครั้ง)	2.1451	0.0000	5.0000
จำนวนงวดที่ค้างชำระมากที่สุด ใน 12 เดือนล่าสุด (ครั้ง)	2.1574	0.0000	5.0000
จำนวนครั้งที่ติดต่อดี ใน 12 เดือนล่าสุด (ครั้ง)	1.9099	0.0000	2.0000
จำนวนครั้งที่ติดต่อดี ใน 9 เดือนล่าสุด (ครั้ง)	1.8857	0.0000	2.0000
จำนวนครั้งที่ติดต่อดี ใน 6 เดือนล่าสุด (ครั้ง)	1.8127	0.0000	2.0000
จำนวนครั้งที่ติดต่อดี ใน 3 เดือนล่าสุด (ครั้ง)	1.6449	0.0000	2.0000
จำนวนครั้งที่ติดต่อดำเนินไม่ได้ ใน 12 เดือนล่าสุด (ครั้ง)	20.3137	0.0000	304.0000
จำนวนครั้งที่ติดต่อดำเนินไม่ได้ ใน 9 เดือนล่าสุด (ครั้ง)	16.6517	0.0000	282.0000
จำนวนครั้งที่ติดต่อดำเนินไม่ได้ ใน 6 เดือนล่าสุด (ครั้ง)	12.1035	0.0000	214.0000
จำนวนครั้งที่ติดต่อดำเนินไม่ได้ ใน 3 เดือนล่าสุด (ครั้ง)	7.0465	0.0000	124.0000

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า

ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากตารางที่ 4.1 แสดงให้เห็นถึงรายละเอียดเชิงสถิติของตัวแปรที่เป็นอัตราส่วน (Ratio) โดยเริ่มจาก อายุของลูกหนี้ มีค่าเฉลี่ยอยู่ที่ 44.3299 ปี มีค่าต่ำสุด 21 ปี และค่าสูงสุด 67 ปี สะท้อนว่า กลุ่มลูกหนี้ส่วนใหญ่อยู่ในช่วงวัยทำงานตอนกลางถึงปลาย อัตราส่วนความสามารถในการชำระหนี้ มีค่าเฉลี่ย 12.7871 หน่วย ค่าต่ำสุดคือ 0 และค่าสูงสุดคือ 624.2447 แสดงให้เห็นถึงความหลากหลาย ในความสามารถในการชำระหนี้ของลูกหนี้ จำนวนบุตร เฉลี่ย 1.3822 คน มีค่าต่ำสุด 0 คน และค่าสูงสุด 10 คน จำนวนเดือนนับจากวันทำสินเชื่อ เฉลี่ย 30.1656 เดือน ต่ำสุดคือ 12 เดือน และสูงสุด 133 เดือน แสดงถึงระยะเวลาการถือครองสินเชื่อที่หลากหลาย ในส่วนของ ระยะเวลากู้ยืม เฉลี่ยอยู่ที่ 12.4577 เดือน โดยมีค่าต่ำสุด 12 เดือน และค่าสูงสุด 156 เดือน อัตราดอกเบี้ยของเงินกู้ มีค่าเฉลี่ย 0.3183 โดยมีค่าต่ำสุดเป็น 0.2400 ส่วนจำนวนวันค้างชำระ เฉลี่ย 61.9211 วัน มีลูกหนี้บางบัญชีที่ไม่มีค้างชำระเลย (0 วัน) และบางบัญชีมีค้างชำระสูงสุดถึง 1,076 วันสำหรับ อัตราส่วน ยอดเงินกู้คงเหลือต่อจำนวนเงินกู้ มีค่าเฉลี่ย 0.8631 โดยมีค่าต่ำสุดเป็นลบที่ -0.0386 และค่าสูงสุด 1.5162 ซึ่งอาจเกิดจากการเบิกเกินวงเงินในบางกรณีทำให้ค่าต่ำสุดติดลบ อัตราส่วนค้างงวดต่อยอดหนี้ เฉลี่ย 0.0699 ต่ำสุด 0 และสูงสุด 0.9474 ค่าเฉลี่ยของอัตราส่วนการเบิกใช้วงเงินในช่วง 3 เดือน เท่ากับ 0.8666 อยู่ในช่วง -0.0385 ถึง 1.5162 ในส่วนของตัวแปร ส่วนของพฤติกรรมชำระหนี้ พบว่า จำนวนครั้งที่ชำระเต็มจำนวนในช่วง 12 เดือนล่าสุด มีค่าเฉลี่ย 4.7833 ครั้ง ต่ำสุด 0 ครั้ง และสูงสุด 12 ครั้ง และลดหลั่นลงมาตามช่วงเวลา (9 เดือน 6 เดือน และ 3 เดือน) ขณะที่ จำนวนครั้งที่ชำระบางส่วน ในแต่ละช่วงเวลาก็ตกหล่นกันเช่นกัน โดยช่วง 12 เดือนล่าสุดเฉลี่ย 5.0616 ครั้ง ส่วนช่วง 3 เดือนล่าสุดเฉลี่ย 1.7783 ครั้ง ในด้านพฤติกรรมการผิดนัด พบว่า จำนวนงวดที่ค้างชำระมากที่สุด ในช่วง 3-12 เดือน มีค่าเฉลี่ยอยู่ที่ประมาณ 2.00-2.15 งวด และสูงสุด 5 งวด สำหรับจำนวนครั้งที่สามารถติดต่อได้ เฉลี่ยประมาณ 1.6-1.9 ครั้ง ในแต่ละช่วง และ จำนวนครั้งที่ติดต่อไม่ได้ เฉลี่ยสูงถึง 20.3137 ครั้งในช่วง 12 เดือนล่าสุด โดยมีสูงสุดถึง 304 ครั้ง ในบางบัญชี ค่าดังกล่าวบ่งชี้ถึงความเสี่ยงในการติดตามหนี้ที่เพิ่มขึ้น ซึ่งอาจสะท้อนปัญหาด้านความร่วมมือของลูกหนี้หรือปัจจัยด้านการเข้าถึงช่องทางสื่อสาร

4.2 ผลการคัดเลือกการแบ่งข้อมูลชุดทดสอบด้วยวิธีแบบแยกชุด (Hold-out) และวิธีแบบไขว้ (K-Fold Cross Validation)

ในการพัฒนาแบบจำลองการทำนาย จำเป็นต้องแบ่งชุดข้อมูลออกเป็นชุดฝึกฝน (Train Set) และชุดทดสอบ (Testing Set) เพื่อประเมินความสามารถของแบบจำลองในการทำนายข้อมูลใหม่ที่ไม่เคยถูกใช้ในการฝึกมาก่อน วิธีการแบ่งข้อมูลมีผลโดยตรงต่อประสิทธิภาพและความน่าเชื่อถือของแบบจำลอง ในการศึกษาครั้งนี้ ผู้วิจัยได้เปรียบเทียบ 2 วิธีที่ได้รับความนิยม ได้แก่ 1) วิธีการแบ่งแบบแยกชุด (Hold-out) เป็นการสุ่มแบ่งข้อมูลทั้งหมดเป็น 2 ส่วน โดยกำหนดให้ 80% เป็นชุดฝึกฝน และ 20% เป็นชุดทดสอบและ 2) วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) เป็นการแบ่ง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ข้อมูลออกเป็น 5 ส่วนเท่าๆ กัน จากนั้นฝึกและทดสอบแบบจำลองสลับกัน 5 ครั้ง โดยแต่ละส่วนจะถูกใช้เป็นชุดทดสอบอย่างน้อยหนึ่งครั้ง การเปรียบเทียบวิธีการแบ่งชุดข้อมูลทั้งสองวิธีนี้ใช้ชุดข้อมูลของกลุ่มที่ 1 (กลุ่มลูกค้าทั่วไป) โดยมีวัตถุประสงค์เพื่อประเมินว่าเทคนิคใดให้ผลลัพธ์ที่มีประสิทธิภาพและเสถียรมากกว่า โดยใช้ค่าประสิทธิภาพโดยรวม (F1-Score) ร่วมกับค่าความระลึก (Recall) ของกลุ่มลูกค้าที่ผิดนัดชำระหนี้ (Class 1) เป็นเกณฑ์หลักในการประเมิน เพราะกลุ่มนี้เป็นเป้าหมายสำคัญในการวิเคราะห์ความเสี่ยงทางการเงิน ผลการทดลองจะเน้นการวิเคราะห์ค่าประสิทธิภาพทั้งในภาพรวมและความสามารถในการตรวจจับกลุ่มผิดนัดชำระ เพื่อประกอบการตัดสินใจเลือกวิธีการแบ่งข้อมูลที่เหมาะสมที่สุดสำหรับการประเมินแบบจำลองทั้ง 5 เทคนิค

ในกระบวนการปรับจูนไฮเปอร์พารามิเตอร์ของแต่ละแบบจำลอง ได้มีการกำหนดช่วงค่าพารามิเตอร์ที่เกี่ยวข้อง และดำเนินการค้นหาค่าที่เหมาะสมที่สุดโดยใช้เทคนิค Grid Search เพื่อให้ได้ชุดพารามิเตอร์ที่ส่งผลให้แบบจำลองมีประสิทธิภาพสูงสุด ผลลัพธ์จากการปรับจูนไฮเปอร์พารามิเตอร์ในแต่ละกรณีจะถูกนำเสนอในรูปแบบตาราง เพื่อเปรียบเทียบผลลัพธ์ของแต่ละแบบจำลองอย่างเป็นระบบและชัดเจนในรูปแบบตารางดังนี้

ตารางที่ 4.2 สรุปไฮเปอร์พารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองสำหรับวิธีการแบ่งแบบแยกชุด (Hold-out)

แบบจำลอง	ค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters)
Logistic Regression	$C = 0.1$, penalty = 'l1', solver = 'liblinear', class_weight = None
Decision Tree	criterion = 'gini', max_depth = 5
Neural Network	activation = 'relu', alpha = 0.01, hidden_layer_sizes = (50,)
LSTM	batch_size = 8, epochs = 30, activation = 'elu', learning_rate = 0.003, units = 256
Bidirectional LSTM	batch_size = 32, epochs = 30, activation = 'elu', learning_rate = 0.0005, units = 64

จากตารางที่ 4.2 ซึ่งสรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลอง สำหรับวิธีการแบ่งข้อมูลแบบแยกชุด (Hold-out) สามารถสังเกตแนวโน้มที่น่าสนใจของค่าพารามิเตอร์ที่ได้ ดังนี้

Logistic Regression ค่าพารามิเตอร์ที่เหมาะสมสะท้อนถึงการเลือกแบบจำลองที่ให้ความสำคัญต่อระหว่างความยืดหยุ่นและความสามารถในการตีความ ส่งผลให้แบบจำลองสามารถรับมือกับข้อมูลที่มีความไม่สมดุลได้อย่างเหมาะสม ไม่ซับซ้อนเกินไป

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Decision Tree แบบจำลองมีแนวโน้มเลือกโครงสร้างที่มีความลึกไม่มาก ซึ่งให้เห็นว่าข้อมูลชุดนี้ไม่ต้องการความซับซ้อนสูงในการแบ่งแยกกลุ่มเป้าหมาย อีกทั้งยังช่วยลดโอกาสเกิด Overfitting และทำให้แบบจำลองคงความเรียบง่าย

Neural Network ค่าพารามิเตอร์ที่ได้มักแสดงให้เห็นถึงการเลือกขนาดชั้นซ่อนและอัตราการเรียนรู้ที่เหมาะสม เพื่อป้องกัน Overfitting และรองรับข้อมูลที่มีความซับซ้อนระดับหนึ่ง ทำให้แบบจำลองมีความยืดหยุ่นในการจับความสัมพันธ์ของข้อมูลมากขึ้น

LSTM และ Bidirectional LSTM มีแนวโน้มเลือกใช้ขนาด Batch Size ที่ไม่ใหญ่จนเกินไป และหน่วยความจำ (Units) ที่เหมาะสม ประกอบกับการใช้งานฟังก์ชันกระตุ้นและอัตราการเรียนรู้ที่ตอบสนองกับลักษณะ Sequence ของข้อมูล ช่วยให้แบบจำลองทั้งสองสามารถเรียนรู้ข้อมูลเชิงลำดับและความสัมพันธ์ในข้อมูลได้อย่างมีประสิทธิภาพ

สรุปจากไฮเปอร์พารามิเตอร์ที่ได้ในแต่ละแบบจำลอง จะเห็นได้ว่าแบบจำลองที่ดีจะเลือกความซับซ้อนที่เหมาะสมกับข้อมูลและลดโอกาสเกิด Overfitting ขณะเดียวกันก็รักษาความสามารถในการเรียนรู้รูปแบบสำคัญของข้อมูลไว้ได้อย่างครบถ้วน ผลการปรับจูนจึงช่วยยืนยันว่าการเลือกค่าพารามิเตอร์ที่เหมาะสมในแต่ละเทคนิค มีผลโดยตรงต่อประสิทธิภาพและความน่าเชื่อถือของแบบจำลองในการทำนายข้อมูลใหม่

ตารางที่ 4.3 สรุปไฮเปอร์พารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองสำหรับวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation)

แบบจำลอง	ค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters)
Logistic Regression	C = 0.01, penalty = 'l2', solver = 'lbfgs', class_weight = None
Decision Tree	criterion = 'gini', max_depth = 3
Neural Network	activation = 'relu', alpha = 0.005, hidden_layer_sizes = (10,)
LSTM	batch_size = 8, epochs = 30, activation = 'tanh', learning_rate = 0.003, units = 256
Bidirectional LSTM	batch_size = 32, epochs = 30, activation = 'tanh', learning_rate = 0.0005, units = 128

จากตารางที่ 4.3 ซึ่งสรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลอง สำหรับวิธีการแบ่งข้อมูลแบบไขว้ (K-Fold Cross Validation) พบว่าการปรับจูนพารามิเตอร์ด้วยเทคนิค Grid Search ส่งผลให้แบบจำลองแต่ละแบบสามารถเลือกชุดค่าพารามิเตอร์ที่เหมาะสมกับลักษณะข้อมูลและโครงสร้างของแบบจำลองได้อย่างมีประสิทธิภาพ โดยสามารถสรุปแนวโน้มของแต่ละแบบจำลองได้ดังนี้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Logistic Regression การเลือกพารามิเตอร์ที่เหมาะสมช่วยให้แบบจำลองมีความสามารถในการจัดการกับข้อมูลที่มีความไม่สมดุล รวมถึงยังคงจุดเด่นด้านความเรียบง่ายและความโปร่งใสในการตีความผลลัพธ์

Decision Tree พบว่าโครงสร้างของต้นไม้ที่ไม่ซับซ้อนและความลึกที่เหมาะสม สามารถลดโอกาสการเกิด Overfitting และให้ประสิทธิภาพการทำนายที่น่าเชื่อถือในบริบทของข้อมูลจริง

Neural Network การปรับจูนช่วยให้แบบจำลองสามารถเรียนรู้ข้อมูลที่ซับซ้อนและหลากหลายมากขึ้น ขณะเดียวกันก็ป้องกันปัญหา Overfitting ผ่านการกำหนดขนาดและโครงสร้างของชั้นซ่อนและอัตราการเรียนรู้ที่เหมาะสม

LSTM และ Bidirectional LSTM ค่าพารามิเตอร์ที่ได้สะท้อนถึงการเลือกขนาด Batch Size และจำนวนหน่วยความจำ (Units) ที่สมดุลกับความซับซ้อนของข้อมูล Sequence ช่วยให้แบบจำลองสามารถเรียนรู้ความสัมพันธ์เชิงลำดับในข้อมูลได้อย่างมีประสิทธิภาพ อีกทั้งยังสามารถป้องกัน Overfitting ได้ดีเมื่อมีการกำหนดอัตราการเรียนรู้และฟังก์ชันกระตุ้นที่เหมาะสม

สรุปการปรับจูนไฮเปอร์พารามิเตอร์โดยใช้วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ช่วยให้การเลือกค่าพารามิเตอร์มีความแม่นยำและเหมาะสมมากขึ้นสำหรับแต่ละแบบจำลอง ช่วยเพิ่มศักยภาพในการทำนายและลดความเสี่ยงของการ Overfitting ผลลัพธ์ที่ได้จึงสะท้อนถึงความสำคัญของกระบวนการปรับจูนไฮเปอร์พารามิเตอร์ในแต่ละเทคนิค เพื่อให้ได้แบบจำลองที่มีประสิทธิภาพสูงสุดและเหมาะสมกับข้อมูลในแต่ละสถานการณ์ เมื่อได้ค่าไฮเปอร์พารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) แต่ละแบบจำลองของวิธีการแบ่งแบบแยกชุด (Hold-out) และวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) แล้ว ในหัวข้อต่อไปจะเป็นการนำเสนอผลการทดลองเปรียบเทียบประสิทธิภาพของแต่ละแบบจำลองภายใต้เงื่อนไขการแบ่งข้อมูลทั้งสองวิธี

4.2.1 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Logistic Regression

จากการแบ่งข้อมูลชุดทดสอบโดยการใช้เทคนิคการแบ่งชุดข้อมูลด้วย 1) วิธีแบบแยกชุด (Hold-out) และ 2) วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) สำหรับแบบจำลอง Logistic Regression พบว่าทั้งสองเทคนิคมีค่าความเที่ยงตรง (Accuracy), ค่าความระลึก (Recall), ค่าความแม่นยำ (Precision) และค่าประสิทธิภาพโดยรวม (F1-Score) ที่ใกล้เคียงกัน โดยเทคนิคการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าความเที่ยงตรง (Accuracy) ที่ร้อยละ 79.88 และค่าประสิทธิภาพโดยรวม (F1-Score) ที่ร้อยละ 79.86 ซึ่งสูงกว่าวิธีแบบแยกชุด (Hold-out) เล็กน้อย โดยมีค่าความเที่ยงตรง (Accuracy) ที่ร้อยละ 79.48 และค่าประสิทธิภาพโดยรวม (F1-Score) ที่ร้อยละ 79.44 เมื่อพิจารณาค่าความระลึก (Recall) ของกลุ่มผิคนัดชำระหนี้ (Class 1) พบว่า วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าสูงกว่า คือร้อยละ 76.70 ขณะที่วิธีแบบแยกชุด (Hold-out) ให้ค่าร้อยละ 75.41 ส่วนค่าความแม่นยำ (Precision) ของกลุ่มผิคนัดชำระหนี้ (Class 1)

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

พบว่า วิธีแบบแยกชุด (Hold-out) มีค่าสูงกว่าเล็กน้อย คือร้อยละ 79.71 ขณะที่วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าร้อยละ 79.48 ดังแสดงในตารางที่ 4.4

ตารางที่ 4.4 ผลการจำแนกประเภทของข้อมูลชุดทดสอบในการเลือกวิธีการแบ่งชุดข้อมูลของแบบจำลอง Logistic Regression

เทคนิคการแบ่งชุดข้อมูล	วิธีแบบแยกชุด	วิธีการแบ่งแบบไขว้
ค่าความเที่ยงตรง	0.7948	0.7988
ค่าประสิทธิภาพโดยรวม	0.7944	0.7986
ค่าความระลึก	0	0.8308
	1	0.7541
ค่าความแม่นยำ	0	0.7930
	1	0.7971

ดังนั้นสรุปได้ว่า ในการทำนายการผิมนัดชำระหนี้โดยใช้แบบจำลอง Logistic Regression ควรเลือกใช้เทคนิคการแบ่งข้อมูลด้วยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) เนื่องจากให้ค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall) ของกลุ่มผิมนัดชำระหนี้ (Class 1) สูงกว่า แม้ความแตกต่างจะมีเพียงเล็กน้อย แต่ก็แสดงให้เห็นถึงความเสถียรและความน่าเชื่อถือของผลลัพธ์ที่ดีกว่าในการประเมินแบบจำลอง

4.2.2 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Decision Tree

จากการแบ่งข้อมูลชุดทดสอบโดยใช้แบบจำลอง Decision Tree เปรียบเทียบระหว่างเทคนิคการแบ่งข้อมูลด้วยวิธีแบบแยกชุด (Hold-out) และวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) พบว่า ทั้งสองเทคนิคให้ผลลัพธ์ที่ใกล้เคียงกันในด้านค่าความเที่ยงตรง (Accuracy), ค่าประสิทธิภาพโดยรวม (F1-Score), ค่าความระลึก (Recall) และค่าความแม่นยำ (Precision) โดยเทคนิควิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าประสิทธิภาพโดยรวม (F1-Score) สูงสุดที่ร้อยละ 79.59 ในขณะที่วิธีแบบแยกชุด (Hold-out) ได้ค่าประสิทธิภาพโดยรวม (F1-Score) อยู่ที่ร้อยละ 79.29 ซึ่งแตกต่างกันเล็กน้อย ส่วนค่าความเที่ยงตรง (Accuracy) ของวิธีแบบแยกชุด (Hold-out) อยู่ที่ร้อยละ 79.44 ขณะที่วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) มีค่าเท่ากับ 79.58 เมื่อพิจารณาค่าความระลึก (Recall) สำหรับกลุ่มผิมนัดชำระหนี้ (Class 1) พบว่าวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าสูงกว่า คือร้อยละ 79.63 เทียบกับวิธีแบบแยกชุด (Hold-out) ที่ร้อยละ 71.76 ส่วนค่าความแม่นยำ (Precision) ของกลุ่มผิมนัดชำระหนี้ (Class 1) พบว่าวิธีแบบแยกชุด (Hold-out) ให้ค่าสูงกว่า คือร้อยละ 82.10 เทียบกับวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ที่ร้อยละ 77.29 ดังแสดงในตารางที่ 4.5

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.5 ผลการจำแนกประเภทของข้อมูลชุดทดสอบในการเลือกวิธีการแบ่งชุดข้อมูลของแบบจำลอง Decision Tree

เทคนิคการแบ่งชุดข้อมูล		วิธีแบบแยกชุด	วิธีการแบ่งแบบไขว้
ค่าความเที่ยงตรง		0.7944	0.7958
ค่าประสิทธิภาพโดยรวม		0.7929	0.7959
ค่าความระลึก	0	0.8621	0.7953
	1	0.7176	0.7963
ค่าความแม่นยำ	0	0.7759	0.8169
	1	0.8210	0.7729

ดังนั้นสรุปได้ว่า ในการทำนายการผิมนัดชำระหนี้โดยใช้แบบจำลอง Decision Tree ควรเลือกใช้เทคนิคการแบ่งข้อมูลด้วยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) เนื่องจากให้ค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall) ของกลุ่มผิมนัดชำระหนี้ (Class 1) สูงกว่า แม้ค่าความแม่นยำจะต่ำกว่าวิธีแบบแยกชุดเล็กน้อย แต่ความสามารถในการตรวจจับกลุ่มเป้าหมายได้ดีขึ้นแสดงให้เห็นถึงความเสถียรและความน่าเชื่อถือของผลลัพธ์ที่ดีกว่าในการประเมินแบบจำลอง

4.2.3 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Neural Network

จากการแบ่งข้อมูลชุดทดสอบโดยใช้แบบจำลอง Neural Network เปรียบเทียบระหว่างเทคนิคการแบ่งข้อมูลด้วยวิธีแบบแยกชุด (Hold-out) และวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) พบว่าทั้งสองเทคนิคให้ผลลัพธ์ที่ใกล้เคียงกันในหลายด้าน โดยเทคนิควิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าความเที่ยงตรง (Accuracy) สูงกว่าที่ร้อยละ 80.59 เทียบกับวิธีแบบแยกชุด (Hold-out) ที่ร้อยละ 79.44 เมื่อพิจารณาค่าประสิทธิภาพโดยรวม (F1-Score) ซึ่งเป็นเกณฑ์หลักของการศึกษานี้ พบว่าวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าประสิทธิภาพโดยรวมเท่ากับร้อยละ 80.45 สูงกว่าวิธีแบบแยกชุด (Hold-out) ที่ให้ค่าร้อยละ 79.29 ด้านค่าความระลึก (Recall) พบว่าวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าความระลึก (Recall) ของกลุ่มผิมนัดชำระหนี้ (Class 1) ที่ร้อยละ 73.14 สูงกว่าวิธีแบบแยกชุด (Hold-out) ที่ได้ร้อยละ 71.76 เช่นเดียวกับค่าความแม่นยำ (Precision) ของกลุ่มผิมนัดชำระหนี้ (Class 1) ที่วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าสูงกว่า คือร้อยละ 83.24 เทียบกับวิธีแบบแยกชุด (Hold-out) ที่ร้อยละ 82.10 ดังแสดงในตารางที่ 4.6

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.6 ผลการจำแนกประเภทของข้อมูลชุดทดสอบในการเลือกวิธีการแบ่งชุดข้อมูลของแบบจำลอง Neural Network

เทคนิคการแบ่งชุดข้อมูล	วิธีแบบแยกชุด	วิธีการแบ่งแบบไขว้
ค่าความเที่ยงตรง	0.7944	0.8059
ค่าประสิทธิภาพโดยรวม	0.7929	0.8045
ค่าความระลึก	0	0.8621
	1	0.7176
ค่าความแม่นยำ	0	0.7759
	1	0.8210

ดังนั้นสรุปได้ว่า ในการทำนายการผิมนัดชำระหนี้โดยใช้แบบจำลอง Neural Network ควรเลือกใช้เทคนิคการแบ่งชุดข้อมูลด้วยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) เนื่องจากให้ค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall) ของกลุ่มผิมนัดชำระหนี้ (Class 1) สูงกว่า ซึ่งแสดงให้เห็นถึงความสามารถของแบบจำลองในการตรวจจับกลุ่มเป้าหมายได้ดีกว่า และยังมี ความเสถียรในการประเมินผลแบบจำลองที่เหมาะสมกว่า

4.2.4 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง LSTM

จากการแบ่งข้อมูลชุดทดสอบโดยใช้แบบจำลอง LSTM โดยเปรียบเทียบระหว่างเทคนิคการแบ่งข้อมูลด้วยวิธีแบบแยกชุด (Hold-out) และวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) พบว่า ทั้งสองเทคนิคให้ผลลัพธ์ที่ใกล้เคียงกันในด้านค่าความเที่ยงตรง (Accuracy) โดยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าสูงกว่าเล็กน้อยที่ร้อยละ 80.51 ขณะที่วิธีแบบแยกชุด (Hold-out) ได้ร้อยละ 80.18 เมื่อพิจารณาค่าประสิทธิภาพโดยรวม (F1-Score) ซึ่งเป็นตัวชี้วัดหลักของการศึกษา พบว่าวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าประสิทธิภาพโดยรวม (F1-Score) สูงสุดที่ร้อยละ 80.41 ขณะที่วิธีแบบแยกชุด (Hold-out) ได้ค่าเท่ากับร้อยละ 80.07 ในด้านค่าความระลึก (Recall) เทคนิคการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าความระลึก (Recall) ของกลุ่มผิมนัดชำระหนี้ (Class 1) สูงกว่าคือร้อยละ 74.21 เทียบกับวิธีแบบแยกชุด (Hold-out) ที่ร้อยละ 73.75 ขณะที่ค่าความระลึก (Recall) ของกลุ่มไม่ผิมนัดชำระหนี้ (Class 0) วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าเท่ากับร้อยละ 86.02 สูงกว่าวิธีแบบแยกชุด (Hold-out) ที่ร้อยละ 85.85 ส่วนค่าความแม่นยำ (Precision) ของกลุ่มผิมนัดชำระหนี้ (Class 1) วิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าสูงกว่าคือร้อยละ 82.29 เทียบกับวิธีแบบแยกชุด (Hold-out) ที่ร้อยละ 82.12 ดังตารางที่ 4.7

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.7 ผลการจำแนกประเภทของข้อมูลชุดทดสอบในการเลือกวิธีการแบ่งชุดข้อมูลของแบบจำลอง LSTM

เทคนิคการแบ่งชุดข้อมูล	วิธีแบบแยกชุด	วิธีการแบ่งแบบไขว้
ค่าความเที่ยงตรง	0.8018	0.8051
ค่าประสิทธิภาพโดยรวม	0.8007	0.8041
ค่าความระลึก	0	0.8585
	1	0.7375
ค่าความแม่นยำ	0	0.7877
	1	0.8212

ดังนั้นสรุปได้ว่า ในการทำนายการผิנדชนิดชำระหนี้โดยใช้แบบจำลอง LSTM ควรเลือกใช้เทคนิคการแบ่งข้อมูลด้วยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) เนื่องจากให้ค่าประสิทธิภาพโดยรวม (F1-Score) สูงที่สุด อีกทั้งยังแสดงถึงความเสถียรของแบบจำลองในภาพรวม และยังให้ค่าความระลึก (Recall) และค่าความแม่นยำ (Precision) ของกลุ่มผิנדชำระหนี้ (Class 1) สูงกว่า แม้ความแตกต่างจะมีเพียงเล็กน้อยก็ตาม จึงเหมาะสมมากกว่าในการนำไปใช้ประเมินผลการทำนาย

4.2.5 ผลการแบ่งข้อมูลชุดฝึกฝนและชุดทดสอบของแบบจำลอง Bidirectional LSTM

จากการแบ่งข้อมูลชุดทดสอบโดยใช้แบบจำลอง Bidirectional LSTM โดยเปรียบเทียบระหว่างเทคนิคการแบ่งข้อมูลด้วยวิธีแบบแยกชุด (Hold-out) และวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) พบว่า ทั้งสองเทคนิคให้ผลลัพธ์ที่ใกล้เคียงกันในด้านค่าความเที่ยงตรง (Accuracy) โดยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าสูงกว่าเล็กน้อยที่ร้อยละ 80.68 ขณะที่วิธีแบบแยกชุด (Hold-out) ได้ร้อยละ 80.39 เมื่อพิจารณาค่าประสิทธิภาพโดยรวม (F1-Score) พบว่าวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าประสิทธิภาพโดยรวม (F1-Score) สูงสุดที่ร้อยละ 80.56 ขณะที่วิธีแบบแยกชุด (Hold-out) ได้ค่าเท่ากับร้อยละ 80.20 ในด้านค่าความระลึก (Recall) พบว่าวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ให้ค่าความระลึกของกลุ่มผิנדชำระหนี้ (Class 1) สูงกว่าคือร้อยละ 73.65 เทียบกับวิธีแบบแยกชุด (Hold-out) ที่ร้อยละ 71.56 ขณะที่ค่าความระลึก (Recall) ของกลุ่มไม่ผิנדชำระหนี้ (Class 0) วิธีแบบแยกชุด (Hold-out) ให้ค่าสูงกว่าคือร้อยละ 88.16 เทียบกับวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ที่ได้ร้อยละ 86.83 สำหรับค่าความแม่นยำ (Precision) ของกลุ่มผิנדชำระหนี้ (Class 1) วิธีแบบแยกชุด (Hold-out) ให้ค่าสูงกว่าคือร้อยละ 84.20 เทียบกับวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) ที่ได้ร้อยละ 83.04 ดังแสดงในตารางที่ 4.8

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.8 ผลการจำแนกประเภทของข้อมูลชุดทดสอบในการเลือกวิธีการแบ่งชุดข้อมูลของแบบจำลอง Bidirectional LSTM

เทคนิคการแบ่งชุดข้อมูล	วิธีแบบแยกชุด	วิธีการแบ่งแบบไขว้
ค่าความเที่ยงตรง	0.8039	0.8129
ค่าประสิทธิภาพโดยรวม	0.8020	0.8121
ค่าความระลึก	0	0.8592
	1	0.7599
ค่าความแม่นยำ	0	0.8036
	1	0.8252

ดังนั้นสรุปได้ว่า ในการทำนายการผัดหน้าชำระหนี้โดยใช้แบบจำลอง Bidirectional LSTM ควรเลือกใช้เทคนิคการแบ่งข้อมูลด้วยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) เนื่องจากให้ค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall) ของกลุ่มผัดหน้าชำระหนี้ (Class 1) สูงกว่า แม้ว่าวิธีแบบแยกชุด (Hold-out) จะให้ค่าความแม่นยำสูงกว่าเล็กน้อย แต่ K-Fold Cross Validation แสดงถึงความเสถียรและความสามารถในการตรวจจับกลุ่มเป้าหมายได้ดีกว่า จึงเหมาะสมมากกว่าสำหรับการนำไปใช้ประเมินผลการทำนาย

สรุปผู้วิจัยเลือกใช้วิธีการแบ่งชุดข้อมูลด้วยวิธีการแบ่งแบบไขว้ (K-Fold Cross Validation) สำหรับแบบจำลองทั้ง 5 แบบ ได้แก่ Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM เนื่องจากให้ค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall: Class 1) สูงกว่าวิธีแบบแยกชุด (Hold-out) ในทุกแบบจำลอง ซึ่งสอดคล้องกับหลักเกณฑ์ที่ใช้ในงานวิจัยนี้

4.3 ผลการทำนายการผัดหน้าชำระหนี้ในกลุ่มที่ 1 ลูกค้าทั่วไป (General Customer)

ในหัวข้อนี้ ผู้วิจัยได้ดำเนินการวิเคราะห์ผลการทำนายการผัดหน้าชำระหนี้ในกลุ่มลูกค้าทั้งหมด ซึ่งถือเป็นกลุ่มลูกค้าทั่วไป (General Customer) ประกอบด้วยลูกค้าทุกบัญชีในชุดข้อมูลที่อยู่ในขอบเขตการศึกษา โดยไม่แบ่งแยกตามประวัติการชำระหนี้ การวิเคราะห์ในกลุ่มนี้มีวัตถุประสงค์เพื่อประเมินความสามารถของแต่ละแบบจำลองในการคัดแยกกลุ่มลูกค้าที่มีแนวโน้มผัดหน้าชำระหนี้ที่เป็นภาพรวมของลูกค้าทั้งหมด

จากข้อสรุปในหัวข้อ 4.2 ผู้วิจัยได้เลือกใช้วิธีการแบ่งชุดข้อมูลแบบไขว้ (K-Fold Cross Validation) สำหรับการประเมินประสิทธิภาพของแบบจำลองทั้ง 5 เทคนิค ได้แก่ Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM ผู้วิจัยใช้เกณฑ์ประสิทธิภาพโดยรวม (F1-Score) ควบคู่กับค่าความระลึก (Recall) เพื่อสะท้อนศักยภาพในการ

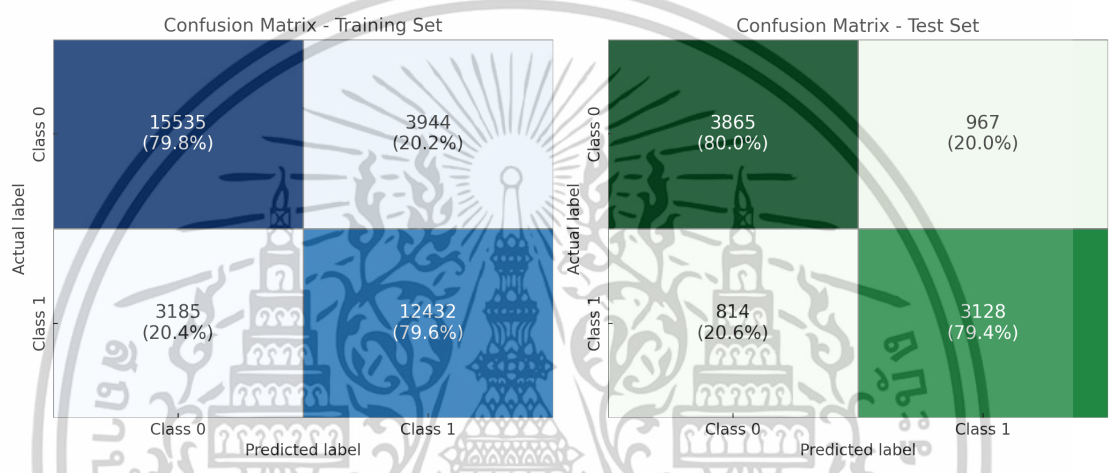
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ทำนายลูกค้าที่มีความเสี่ยงผิดนัดชำระหนี้ได้อย่างถูกต้องและครอบคลุม เพื่อสนับสนุนการบริหารความเสี่ยงและการกำหนดกลยุทธ์เชิงป้องกันในเชิงนโยบายอย่างมีประสิทธิภาพ

4.3.1 ผลการทำนายการผิดนัดชำระหนี้ของแบบจำลอง Logistic Regression

1) เมทริกซ์ความสับสน (Confusion Matrix)

เพื่อประเมินประสิทธิภาพของแบบจำลอง Logistic Regression ในการทำนายสถานะของลูกค้า ได้จัดทำเมทริกซ์ความสับสน (Confusion Matrix) เปรียบเทียบระหว่างข้อมูลชุดฝึกฝน (Train Set) และชุดทดสอบ (Test Set) ดังแสดงในรูปที่ 4.9



รูปที่ 4.9 เมทริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Logistic Regression

จากรูปที่ 4.9 พบว่าแบบจำลอง Logistic Regression สามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ในชุดฝึกฝนได้ถูกต้อง 15,535 บัญชี หรือคิดเป็นร้อยละ 79.8 และทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 12,432 บัญชี หรือคิดเป็นร้อยละ 79.6 ในขณะที่ข้อมูลชุดทดสอบพบว่าแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 3,865 บัญชี คิดเป็นร้อยละ 80.0 และทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 3,128 บัญชี คิดเป็นร้อยละ 79.4 สำหรับในส่วนของการทำนายผิดพลาด พบว่าทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ว่าเป็นผิดนัดชำระหนี้ (Default) ในชุดฝึกฝนจำนวน 3,944 บัญชี คิดเป็นร้อยละ 20.2 และในชุดทดสอบจำนวน 967 บัญชี คิดเป็นร้อยละ 20.0 ในขณะที่แบบจำลองทำนายกลุ่มผิดนัดชำระหนี้ (Default) เป็นไม่ผิดนัดชำระหนี้ (Non-Default) ในชุดฝึกฝนจำนวน 3,185 บัญชี คิดเป็นร้อยละ 20.4 และในชุดทดสอบจำนวน 814 บัญชี คิดเป็นร้อยละ 20.6 เมื่อพิจารณาผลการเปรียบเทียบระหว่างชุดฝึกฝนและชุดทดสอบ พบว่าสัดส่วนของการทำนายผิดพลาดของแบบจำลองในทั้งสองชุดข้อมูลมีความใกล้เคียงกัน ซึ่งสะท้อนให้เห็นว่าแบบจำลอง Logistic Regression มีความสามารถในการ

การเรียนรู้ที่เหมาะสม เนื่องจากไม่พบปัญหา Overfitting. เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2) การจำแนกประเภท (Classification Report)

ตารางที่ 4.9 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Logistic Regression

ตัวชี้วัดประสิทธิภาพ		Logistic Regression	จำนวนข้อมูล
ค่าความเที่ยงตรง		0.7972	8,774
ค่าประสิทธิภาพโดยรวม		0.7969	8,774
ค่าความระลึก	0	0.8262	4,679
	1	0.7641	4,095
ค่าความแม่นยำ	0	0.8001	4,679
	1	0.7938	4,095

จากตารางที่ 4.9 แสดงผลการจำแนกประเภทของการทำนายข้อมูลชุดทดสอบของแบบจำลอง Logistic Regression มีค่าความเที่ยงตรง (Accuracy) เท่ากับ 79.72 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับ 79.69 เมื่อพิจารณาค่าความแม่นยำ (Precision) และความระลึก (Recall) ในแต่ละกลุ่มพบว่า

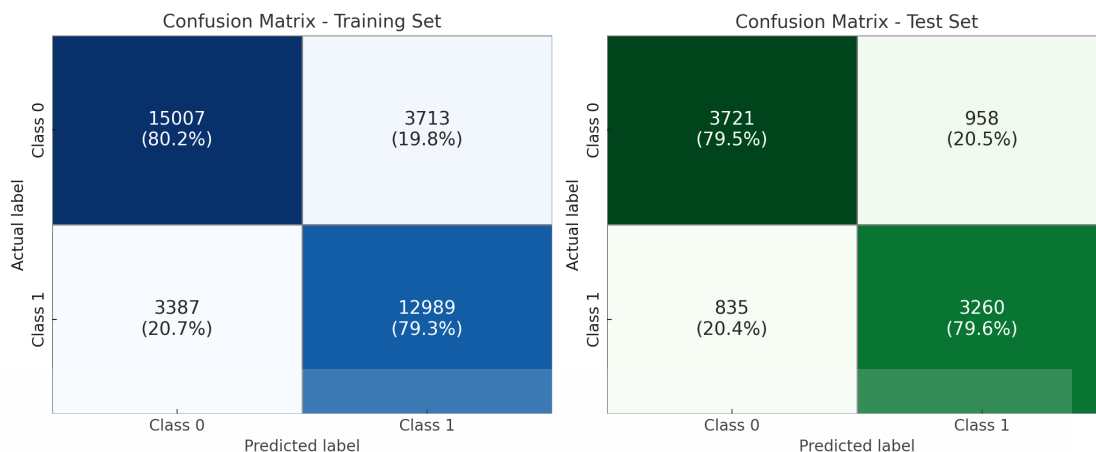
2) ประสิทธิภาพการจำแนกบัญชีที่ไม่ผิดนัดชำระหนี้ (Non-Default) แบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 82.62 ซึ่งหมายความว่าสามารถตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราสูง นอกจากนี้ ค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 80.01

ในการศึกษาครั้งนี้ใช้เกณฑ์เปรียบเทียบจากค่าประสิทธิภาพโดยรวม (F1-Score) ที่มีค่าเท่ากับร้อยละ 79.69 และค่าความระลึก (Recall) สำหรับกลุ่มที่ผิดนัดชำระหนี้ (Default : Class 1) ที่มีค่าเท่ากับร้อยละ 76.41

4.3.2 ผลการทำนายการผิดนัดชำระหนี้ของแบบจำลอง Decision Tree

1) เมทริกซ์ความสับสน (Confusion Matrix)

เพื่อประเมินประสิทธิภาพของแบบจำลอง Decision Tree ในการทำนายสถานะของลูกค้า ได้จัดทำเมทริกซ์ความสับสน (Confusion Matrix) เปรียบเทียบระหว่างข้อมูลชุดฝึกฝน (Train Set) และชุดทดสอบ (Test Set) ดังแสดงในรูปที่ 4.10



รูปที่ 4.10 เมตริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Decision Tree

จากรูปที่ 4.10 พบว่าในชุดฝึกฝนแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 15,007 บัญชี หรือคิดเป็นร้อยละ 80.2 และสามารถทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 12,989 บัญชี หรือคิดเป็นร้อยละ 79.3 ขณะที่ในชุดทดสอบแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 3,721 บัญชี หรือคิดเป็นร้อยละ 79.5 และสามารถทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 3,260 บัญชี หรือคิดเป็นร้อยละ 79.6 สำหรับการทำนายผิดพลาด พบว่าแบบจำลองทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ว่าเป็นผิดนัดชำระหนี้ (Default) ในชุดฝึกฝนจำนวน 3,713 บัญชี คิดเป็นร้อยละ 19.8 และในชุดทดสอบจำนวน 958 บัญชี คิดเป็นร้อยละ 20.5 ขณะที่แบบจำลองทำนายกลุ่มผิดนัดชำระหนี้ (Default) ว่าเป็นไม่ผิดนัดชำระหนี้ (Non-Default) ในชุดฝึกฝนจำนวน 3,387 บัญชี คิดเป็นร้อยละ 20.7 และในชุดทดสอบจำนวน 835 บัญชี คิดเป็นร้อยละ 20.4 เมื่อพิจารณาผลการเปรียบเทียบระหว่างชุดฝึกฝนและชุดทดสอบ พบว่าสัดส่วนของการทำนายผิดพลาดของแบบจำลองในทั้งสองชุดข้อมูลมีความใกล้เคียงกัน ซึ่งสะท้อนให้เห็นว่าแบบจำลอง Decision Tree มีความสามารถในการเรียนรู้ที่เหมาะสม และสามารถนำไปประยุกต์ใช้กับข้อมูลใหม่ได้จริง เนื่องจากไม่พบปัญหา Overfitting

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2) การจำแนกประเภท (Classification Report)

ตารางที่ 4.10 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Decision Tree

ตัวชี้วัดประสิทธิภาพ		Decision Tree	จำนวนข้อมูล
ค่าความเที่ยงตรง		0.7958	8,774
ค่าประสิทธิภาพโดยรวม		0.7959	8,774
ค่าความระลึก	0	0.7953	4,679
	1	0.7963	4,095
ค่าความแม่นยำ	0	0.8169	4,679
	1	0.7729	4,095

จากตารางที่ 4.10 แสดงผลการจำแนกประเภทของการทำนายข้อมูลชุดทดสอบของแบบจำลอง Decision Tree มีค่าความเที่ยงตรง (Accuracy) เท่ากับ 79.58 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับ 79.59 เมื่อพิจารณาค่าความแม่นยำ (Precision) และความระลึก (Recall) ในแต่ละกลุ่มพบว่า

1) ประสิทธิภาพการจำแนกบัญชีที่ผิดนัดชำระหนี้ (Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 79.63 ซึ่งหมายความว่าสามารถตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราสูงเช่นกัน โดยค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 77.29

2) ประสิทธิภาพการจำแนกบัญชีที่ไม่ผิดนัดชำระหนี้ (Non-Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 79.53 ซึ่งหมายความว่าสามารถตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราสูง นอกจากนี้ ค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 81.69

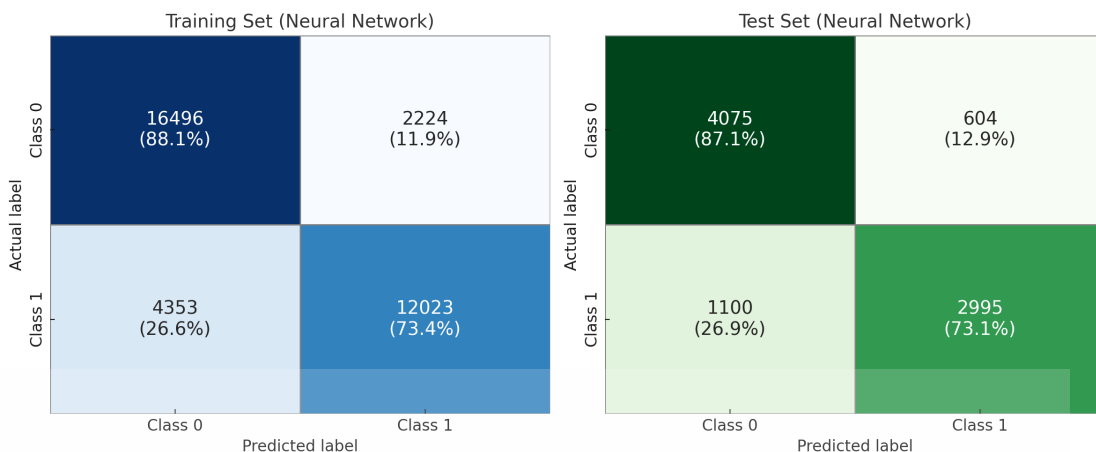
ในการศึกษาครั้งนี้ใช้เกณฑ์เปรียบเทียบจากค่าประสิทธิภาพโดยรวม (F1-Score) ที่มีค่าเท่ากับร้อยละ 79.59 และค่าความระลึก (Recall) สำหรับกลุ่มที่ผิดนัดชำระหนี้ (Default : Class 1) ที่มีค่าเท่ากับร้อยละ 79.63

4.3.3 ผลการทำนายการผิดนัดชำระหนี้ของแบบจำลอง Neural Network

1) เมทริกซ์ความสับสน (Confusion Matrix)

เพื่อประเมินประสิทธิภาพของแบบจำลอง Neural Network ในการทำนายสถานะของลูกค้า ได้จัดทำเมทริกซ์ความสับสน (Confusion Matrix) เปรียบเทียบระหว่างข้อมูลชุดฝึกฝน (Train Set) และชุดทดสอบ (Test Set) ดังแสดงในรูปที่ 4.11

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.11 เมตริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Neural Network

จากรูปที่ 4.11 พบว่าในชุดฝึกฝนแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 16,496 บัญชี หรือคิดเป็นร้อยละ 88.1 และสามารถทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 12,023 บัญชี หรือคิดเป็นร้อยละ 73.4 ขณะที่ในชุดทดสอบแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 4,075 บัญชี หรือคิดเป็นร้อยละ 87.1 และสามารถทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 2,995 บัญชี หรือคิดเป็นร้อยละ 73.1 สำหรับการทำนายผิดพลาด พบว่าแบบจำลองทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ว่าเป็นผิดนัดชำระหนี้ (Default) ในชุดฝึกฝนจำนวน 2,224 บัญชี คิดเป็นร้อยละ 11.9 และในชุดทดสอบจำนวน 604 บัญชี คิดเป็นร้อยละ 12.9 ขณะที่แบบจำลองทำนายกลุ่มผิดนัดชำระหนี้ (Default) ว่าเป็นไม่ผิดนัดชำระหนี้ (Non-Default) ในชุดฝึกฝนจำนวน 4,353 บัญชี คิดเป็นร้อยละ 26.6 และในชุดทดสอบจำนวน 1,100 บัญชี คิดเป็นร้อยละ 26.9 เมื่อพิจารณาผลการเปรียบเทียบระหว่างชุดฝึกฝนและชุดทดสอบ พบว่าสัดส่วนของการทำนายผิดพลาดของแบบจำลองในทั้งสองชุดข้อมูลมีความใกล้เคียงกัน สะท้อนให้เห็นว่าแบบจำลอง Neural Network ไม่พบปัญหา Overfitting

2) การจำแนกประเภท (Classification Report)

ตารางที่ 4.11 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Neural Network

ตัวชี้วัดประสิทธิภาพ	Neural Network	จำนวนข้อมูล
ค่าความเที่ยงตรง	0.8059	8,774
ค่าประสิทธิภาพโดยรวม	0.8045	8,774
ค่าความระลึก	0	4,679
	1	0.7314
ค่าความแม่นยำ	0	4,679
	1	0.8324

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จากตารางที่ 4.11 แสดงผลการจำแนกประเภทของการทำนายข้อมูลชุดทดสอบของแบบจำลอง Neural Network ค่าความเที่ยงตรง (Accuracy) เท่ากับ 80.59 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับ 80.45 เมื่อพิจารณาค่าความแม่นยำ (Precision) และความระลึก (Recall) ในแต่ละกลุ่มพบว่า

1) ประสิทธิภาพการจำแนกบัญชีที่ผิดนัดชำระหนี้ (Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 73.14 ซึ่งหมายความว่าสามารถตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราสูงเช่นกัน โดยค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 83.24

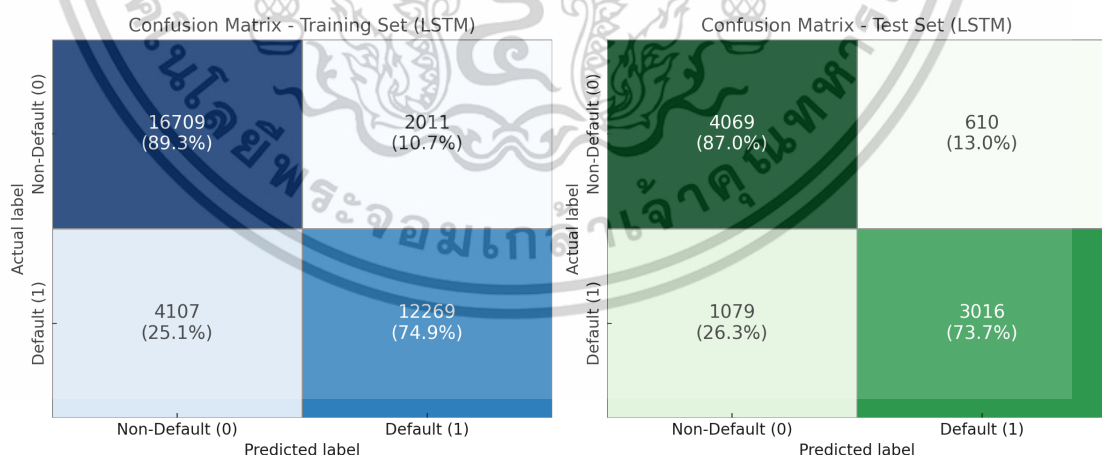
2) ประสิทธิภาพการจำแนกบัญชีที่ไม่ผิดนัดชำระหนี้ (Non-Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 87.11 แสดงถึงความสามารถในการตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราที่สูงมาก นอกจากนี้ ค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 78.75

ในการศึกษาครั้งนี้ใช้เกณฑ์เปรียบเทียบจากค่าประสิทธิภาพโดยรวม (F1-Score) ที่มีค่าเท่ากับร้อยละ 80.45 และค่าความระลึก (Recall) สำหรับกลุ่มที่ผิดนัดชำระหนี้ (Default : Class 1) ที่มีค่าเท่ากับร้อยละ 73.14

4.3.4 ผลการทำนายการผิดนัดชำระหนี้ของแบบจำลอง LSTM

1) เมทริกซ์ความสับสน (Confusion Matrix)

เพื่อประเมินประสิทธิภาพของแบบจำลอง Long Short-Term Memory (LSTM) ในการทำนายสถานะของลูกค้า ได้จัดทำเมทริกซ์ความสับสน (Confusion Matrix) เปรียบเทียบระหว่างข้อมูลชุดฝึกฝน (Train Set) และชุดทดสอบ (Test Set) ดังแสดงในรูปที่ 4.12



รูปที่ 4.12 เมทริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง LSTM

จากรูปที่ 4.12 พบว่าในชุดฝึกฝนแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 16,709 บัญชี หรือคิดเป็นร้อยละ 89.3 และสามารถทำนายกลุ่มผิดนัดชำระหนี้เอกสารนี้เป็นเอกสารที่ส่งวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

(Default) ได้ถูกต้อง 12,269 บัญชี หรือคิดเป็นร้อยละ 74.9 ขณะที่ในชุดทดสอบแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 4,069 บัญชี หรือคิดเป็นร้อยละ 87.0 และสามารถทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 3,016 บัญชี หรือคิดเป็นร้อยละ 73.7 สำหรับการทำนายผิดพลาด พบว่าแบบจำลองทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ว่าเป็นผิดนัดชำระหนี้ (Default) ในชุดฝึกฝนจำนวน 2,011 บัญชี คิดเป็นร้อยละ 10.7 และในชุดทดสอบจำนวน 610 บัญชี คิดเป็นร้อยละ 13.0 ขณะที่แบบจำลองทำนายกลุ่มผิดนัดชำระหนี้ (Default) เป็นไม่ผิดนัดชำระหนี้ (Non-Default) ในชุดฝึกฝนจำนวน 4,107 บัญชี คิดเป็นร้อยละ 25.1 และในชุดทดสอบจำนวน 1,079 บัญชี คิดเป็นร้อยละ 26.3 เมื่อพิจารณาผลการเปรียบเทียบระหว่างชุดฝึกฝนและชุดทดสอบ พบว่าสัดส่วนของการทำนายผิดพลาดของแบบจำลองในทั้งสองชุดข้อมูลมีความใกล้เคียงกัน ซึ่งสะท้อนให้เห็นว่าแบบจำลอง LSTM มีความสามารถในการเรียนรู้ที่เหมาะสมและสามารถนำไปประยุกต์ใช้กับข้อมูลใหม่ได้จริง เนื่องจากไม่พบปัญหา Overfitting

2) การจำแนกประเภท (Classification Report)

ตารางที่ 4.12 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง LSTM

ตัวชี้วัดประสิทธิภาพ	LSTM	จำนวนข้อมูล
ค่าความเที่ยงตรง	0.8052	8,774
ค่าประสิทธิภาพโดยรวม	0.8039	8,774
ค่าความระลึก	0	4,679
	1	0.7314
ค่าความแม่นยำ	0	4,679
	1	0.8310

จากตารางที่ 4.12 แสดงผลการจำแนกประเภทของการทำนายข้อมูลชุดทดสอบของแบบจำลอง LSTM มีค่าความเที่ยงตรง (Accuracy) เท่ากับ 80.52 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับ 80.39 เมื่อพิจารณาค่าความแม่นยำ (Precision) และความระลึก (Recall) ในแต่ละกลุ่มพบว่า

1) ประสิทธิภาพการจำแนกบัญชีที่ผิดนัดชำระหนี้ (Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 73.14 ซึ่งหมายความว่าสามารถตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราที่น่าพอใจ โดยค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 83.10

2) ประสิทธิภาพการจำแนกบัญชีที่ไม่ผิดนัดชำระหนี้ (Non-Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 86.98 ซึ่งแสดงถึงความสามารถในการตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราสูง นอกจากนี้ ค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 78.72

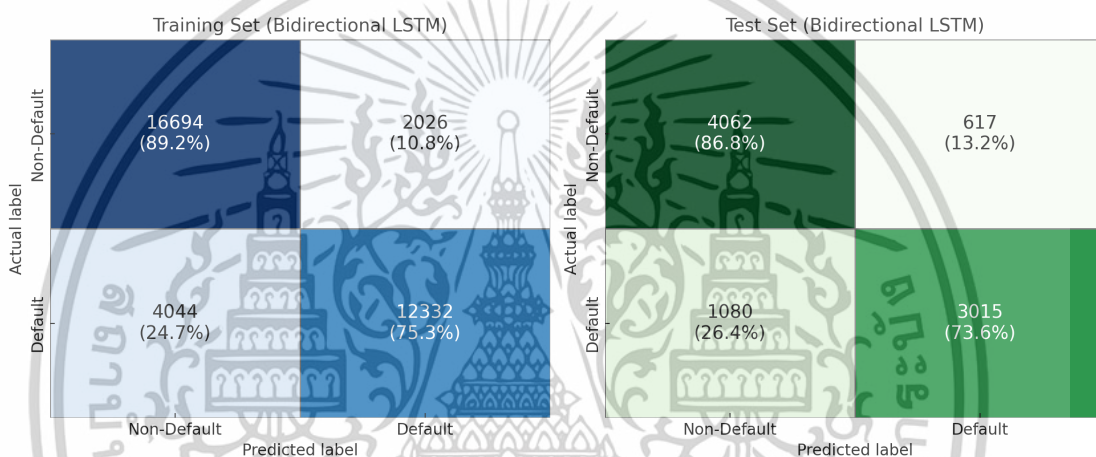
เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ในการศึกษาครั้งนี้ใช้เกณฑ์เปรียบเทียบจากค่าประสิทธิภาพโดยรวม (F1-Score) ที่มีค่าเท่ากับร้อยละ 80.39 และค่าความระลึก (Recall) สำหรับกลุ่มที่ผิดนัดชำระหนี้ (Default : Class 1) ที่มีค่าเท่ากับร้อยละ 73.14

4.3.5 ผลการทำนายการผิดนัดชำระหนี้ของแบบจำลอง Bidirectional LSTM

1) เมทริกซ์ความสับสน (Confusion Matrix)

เพื่อประเมินประสิทธิภาพของแบบจำลอง Bidirectional LSTM ในการทำนายสถานะของลูกค้า ได้จัดทำเมทริกซ์ความสับสน (Confusion Matrix) เปรียบเทียบระหว่างข้อมูลชุดฝึกฝน (Train Set) และชุดทดสอบ (Test Set) ดังแสดงในรูปที่ 4.13



รูปที่ 4.13 เมทริกซ์ความสับสนของการทำนายข้อมูลชุดฝึกฝนกับข้อมูลชุดทดสอบของแบบจำลอง Bidirectional LSTM

จากรูปที่ 4.13 พบว่าในชุดฝึกฝนแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 16,694 บัญชี หรือคิดเป็นร้อยละ 89.2 และสามารถทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 12,332 บัญชี หรือคิดเป็นร้อยละ 75.3 ขณะที่ในชุดทดสอบแบบจำลองสามารถทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ได้ถูกต้อง 4,062 บัญชี หรือคิดเป็นร้อยละ 86.8 และสามารถทำนายกลุ่มผิดนัดชำระหนี้ (Default) ได้ถูกต้อง 3,015 บัญชี หรือคิดเป็นร้อยละ 73.6 สำหรับการทำนายผิดพลาด พบว่าแบบจำลองทำนายกลุ่มไม่ผิดนัดชำระหนี้ (Non-Default) ว่าเป็นผิดนัดชำระหนี้ (Default) ในชุดฝึกฝนจำนวน 2,026 บัญชี คิดเป็นร้อยละ 10.8 และในชุดทดสอบจำนวน 617 บัญชี คิดเป็นร้อยละ 13.2 ขณะที่แบบจำลองทำนายกลุ่มผิดนัดชำระหนี้ (Default) ว่าเป็นไม่ผิดนัดชำระหนี้ (Non-Default) ในชุดฝึกฝนจำนวน 4,044 บัญชี คิดเป็นร้อยละ 24.7 และในชุดทดสอบจำนวน 1,080 บัญชี คิดเป็นร้อยละ 26.4 เมื่อพิจารณาผลการเปรียบเทียบระหว่างชุดฝึกฝนและชุดทดสอบ พบว่าสัดส่วนของการทำนายผิดพลาดของแบบจำลองในทั้งสองชุดข้อมูลมีความ

ใกล้เคียงกัน สามารถนำไปประยุกต์ใช้กับข้อมูลใหม่ได้จริง เนื่องจากไม่พบปัญหา Overfitting เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่นิยมนำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2) การจำแนกประเภท (Classification Report)

ตารางที่ 4.13 ผลการจำแนกประเภทการทำนายข้อมูลชุดทดสอบของแบบจำลอง Bidirectional LSTM

ตัวชี้วัดประสิทธิภาพ		Bidirectional LSTM	จำนวนข้อมูล
ค่าความเที่ยงตรง		0.8068	8,774
ค่าประสิทธิภาพโดยรวม		0.8056	8,774
ค่าความระลึก	0	0.8683	4,679
	1	0.7365	4,095
ค่าความแม่นยำ	0	0.7902	4,679
	1	0.8304	4,095

จากตารางที่ 4.13 แสดงผลการจำแนกประเภทของการทำนายข้อมูลชุดทดสอบของแบบจำลอง Bidirectional LSTM มีค่าความเที่ยงตรง (Accuracy) เท่ากับ 80.68 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับ 80.56 เมื่อพิจารณาค่าความแม่นยำ (Precision) และความระลึก (Recall) ในแต่ละกลุ่มพบว่า

1) ประสิทธิภาพการจำแนกบัญชีที่ผิดนัดชำระหนี้ (Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 73.65 ซึ่งหมายความว่าสามารถตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราที่น่าพึงพอใจ โดยค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 83.04

2) ประสิทธิภาพการจำแนกบัญชีที่ไม่ผิดนัดชำระหนี้ (Non-Default) พบว่าแบบจำลองมีค่าความระลึก (Recall) เท่ากับร้อยละ 86.83 แสดงถึงความสามารถในการตรวจจับกลุ่มนี้ได้ถูกต้องในอัตราสูง นอกจากนี้ ค่าความแม่นยำ (Precision) อยู่ที่ร้อยละ 79.02

ในการศึกษาครั้งนี้ใช้เกณฑ์เปรียบเทียบจากค่าประสิทธิภาพโดยรวม (F1-Score) ที่มีค่าเท่ากับร้อยละ 80.56 และค่าความระลึก (Recall) สำหรับกลุ่มที่ผิดนัดชำระหนี้ (Default : Class 1) ที่มีค่าเท่ากับร้อยละ 73.65

4.3.6 ผลการเปรียบเทียบประสิทธิภาพของ 5 เทคนิคในการทำนายการผิดนัดชำระหนี้

จากการวิเคราะห์ผลการจำแนกการทำนายการผิดนัดชำระหนี้โดยใช้แบบจำลองทั้ง 5 แบบ ได้แก่ Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM พบว่าแบบจำลองแต่ละประเภทมีระดับประสิทธิภาพที่แตกต่างกัน โดยใช้ค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall) เป็นเกณฑ์หลักในการประเมิน เนื่องจากเป็นดัชนีที่สะท้อนทั้งความสมดุลในการจำแนกข้อมูล และความสามารถในการตรวจจับกลุ่มลูกค้าที่ผิดนัดชำระหนี้ได้อย่างมีประสิทธิภาพ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.14 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 1 มลูกค้าทั่วไป (General Customer)

แบบจำลอง	ค่าความเที่ยงตรง	ค่าความแม่นยำ	ค่าความระลึก	ค่าประสิทธิภาพโดยรวม
Logistic Regression	0.7972	0.7971	0.7972	0.7969
Decision Tree	0.7958	0.7964	0.7958	0.7959
Neural Network	0.8059	0.8084	0.8053	0.8045
LSTM	0.8052	0.8077	0.8052	0.8039
Bidirectional LSTM	0.8068	0.8089	0.8068	0.8056

จากตารางที่ 4.14 แสดงให้เห็นว่าแบบจำลอง Bidirectional LSTM มีประสิทธิภาพในการทำนายการผิดนัดชำระหนี้สูงที่สุด โดยมีค่าประสิทธิภาพโดยรวม (F1-Score) อยู่ที่ร้อยละ 80.56 และค่าความระลึก (Recall) ที่ร้อยละ 80.68 ซึ่งสะท้อนให้เห็นถึงประสิทธิภาพที่สมดุลระหว่างความถูกต้องและความสามารถในการระบุลูกค้าที่มีแนวโน้มผิดนัดชำระหนี้ได้ดีที่สุด รองลงมาคือแบบจำลอง Neural Network และ LSTM ที่มีค่าประสิทธิภาพโดยรวม (F1-Score) อยู่ที่ร้อยละ 80.45 และร้อยละ 80.39 ตามลำดับ โดยมีค่าความระลึก (Recall) ใกล้เคียงกัน คือร้อยละ 80.53 และร้อยละ 80.52 ซึ่งสะท้อนถึงศักยภาพในการทำนายที่มีประสิทธิภาพในระดับสูง แต่ยังคงต่ำกว่า Bidirectional LSTM เล็กน้อย สำหรับแบบจำลองที่ใช้เทคนิค Logistic Regression และ Decision Tree แม้ว่าจะมีค่าประสิทธิภาพโดยรวม (F1-Score) ต่ำกว่าแบบจำลองอื่นๆ แต่ยังสามารถรักษาระดับประสิทธิภาพในการจำแนกทั้งสองกลุ่ม (ผิดนัดชำระหนี้และไม่ผิดนัดชำระหนี้) ได้ในระดับที่ดีพอสมควร โดยเฉพาะ Decision Tree ที่สามารถระบุลูกค้าที่ผิดนัดชำระได้แม่นยำดีในระดับหนึ่ง

ดังนั้นสรุปได้ว่าแบบจำลอง Bidirectional LSTM เป็นเทคนิคที่เหมาะสมที่สุดสำหรับการนำไปประยุกต์ใช้ในการบริหารจัดการความเสี่ยงด้านสินเชื่อของสถาบันการเงิน เนื่องจากมีค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall) สูงที่สุดเมื่อเทียบกับเทคนิคอื่นๆ และสอดคล้องกับวัตถุประสงค์ของการวิจัยที่มุ่งเน้นการตรวจจับกลุ่มลูกค้าที่มีแนวโน้มผิดนัดชำระได้อย่างแม่นยำ ครอบคลุม และนำไปใช้ได้จริงในเชิงนโยบาย

4.3.7 การทดสอบสมมติฐานด้วย McNemar's Test

ภายหลังจากการประเมินและเปรียบเทียบประสิทธิภาพของแบบจำลองทั้ง 5 เทคนิค ได้แก่ Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM ตามที่นำเสนอในหัวข้อ 4.3 โดยใช้ค่าประสิทธิภาพโดยรวม (F1-Score) และค่าความระลึก (Recall) เป็นเกณฑ์ ผลการทดลองพบว่าแบบจำลอง Bidirectional LSTM มีค่าประสิทธิภาพสูงสุด อย่างไรก็ตาม เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เพื่อเสริมความน่าเชื่อถือของข้อค้นพบในเชิงสถิติ ผู้วิจัยจึงดำเนินการทดสอบสมมติฐานโดยใช้ McNemar's Test เพื่อวิเคราะห์ว่า Bidirectional LSTM มีความแตกต่างจากแบบจำลองอื่นอย่างมีนัยสำคัญหรือไม่ การทดสอบ McNemar's Test เป็นเทคนิคที่เหมาะสมสำหรับการเปรียบเทียบผลการจำแนกประเภทของแบบจำลองสองแบบที่ทดสอบกับกลุ่มตัวอย่างเดียวกัน (Paired Comparison) โดยให้ความสำคัญกับกรณีที่แบบจำลองหนึ่งทำนายถูก ในขณะที่อีกแบบจำลองหนึ่งทำนายผิด ซึ่งจะวิเคราะห์จากตัวแปร b และ c ในตาราง 2x2 ตามนิยามและสมการ (2.7) ที่นำเสนอในบทที่ 2

การทดสอบนี้ กำหนดสมมติฐานดังนี้

H_0 : ผลการทำนายของแบบจำลองทั้งสองไม่มีความแตกต่างกัน

H_1 : ผลการทำนายของแบบจำลองทั้งสองมีความแตกต่างกัน

ตารางที่ 4.15 ผลการเปรียบเทียบระหว่าง Bidirectional LSTM กับแบบจำลองอื่นๆ

Model Pair	แบบจำลอง A ผิด - B ถูก	แบบจำลอง A ถูก - B ผิด	χ^2	p-value
BiLSTM vs Logistic Regression	480	665	29.88	1.05×10^{-7}
BiLSTM vs Decision Tree	472	672	34.96	3.92×10^{-9}
BiLSTM vs Neural Network	432	585	23.03	1.52×10^{-6}
BiLSTM vs LSTM	447	602	22.89	1.56×10^{-6}

จากตาราง 4.15 พบว่าทุกคู่เปรียบเทียบมีค่า p-value ต่ำกว่าระดับนัยสำคัญ 0.05 อย่างชัดเจน จึงปฏิเสธ H_0 กล่าวคือ ผลการทำนายของ Bidirectional LSTM มีความแตกต่างจากแบบจำลองอื่นๆ อย่างมีนัยสำคัญทางสถิติ

การทดสอบ McNemar's Test ในหัวข้อนี้ช่วยส่งเสริมหลักฐานทางสถิติที่ชัดเจนว่า Bidirectional LSTM ให้ผลลัพธ์ที่แตกต่างและเหนือกว่าแบบจำลองอื่นๆ เมื่อใช้ข้อมูลเดียวกัน

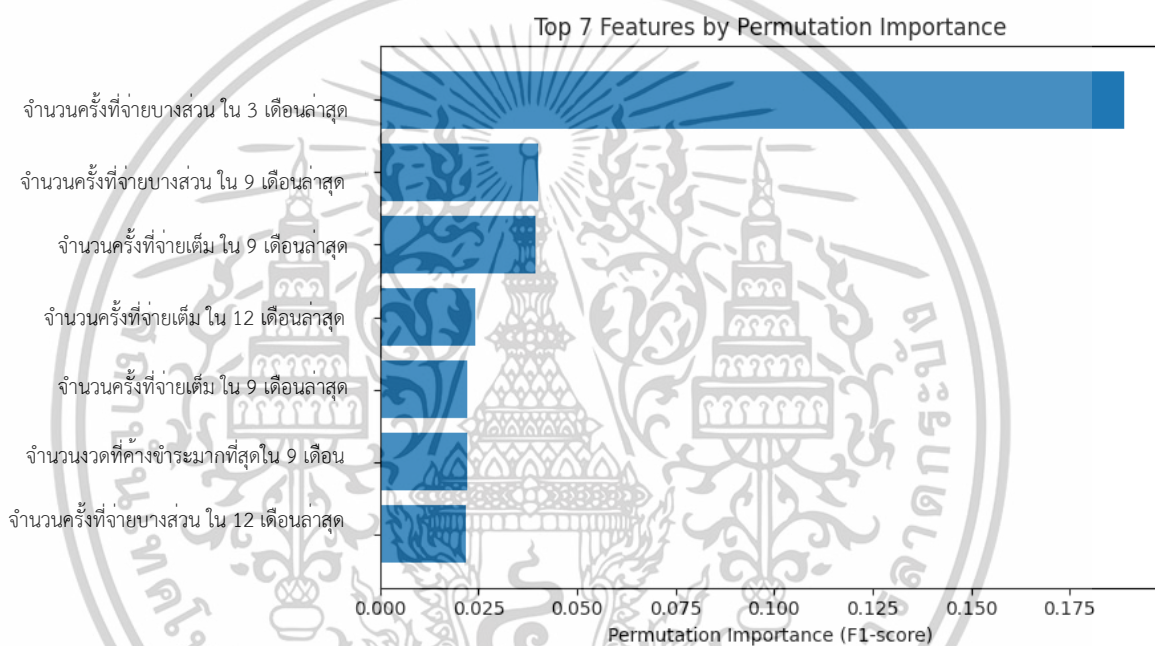
4.3.8 ผลการวิเคราะห์ความสำคัญของตัวแปร (Feature Importance)

การวิเคราะห์ความสำคัญของตัวแปร (Feature Importance) ถือเป็นขั้นตอนที่มีบทบาทอย่างยิ่งในการตีความการตัดสินใจของแบบจำลองการทำนาย ซึ่งในงานวิจัยนี้เน้นการศึกษาจากกลุ่มลูกค้าทั่วไปทั้งหมด (กลุ่มที่ 1) เพื่อให้สามารถนำผลวิเคราะห์ไปสรุปเป็นข้อเสนอแนะเชิงนโยบายและปรับใช้กับลูกค้าทั่วไปได้อย่างเหมาะสม เป้าหมายหลักของการวิเคราะห์นี้คือเพื่อคัดเลือกฟีเจอร์ที่มีอิทธิพลมากที่สุดต่อผลลัพธ์ของแบบจำลอง โดยใช้สองเทคนิคหลัก ได้แก่ 1) Permutation Importance ซึ่งวัดความสำคัญของฟีเจอร์โดยดูผลกระทบต่อค่าความแม่นยำของแบบจำลองเมื่อมีการสลับค่าของฟีเจอร์นั้นแบบสุ่ม และ 2) SHAP (SHapley Additive exPlanations) ซึ่งอธิบายการมีส่วนร่วมของฟีเจอร์แต่ละตัวต่อการทำนายของแบบจำลอง การวิเคราะห์ครอบคลุมแบบจำลองทั้ง 5 เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ประเภท ได้แก่ Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM ซึ่งแต่ละแบบจำลองอาจตีความความสำคัญของตัวแปรแตกต่างกันไปตามโครงสร้าง โดยสามารถอธิบายได้ดังนี้

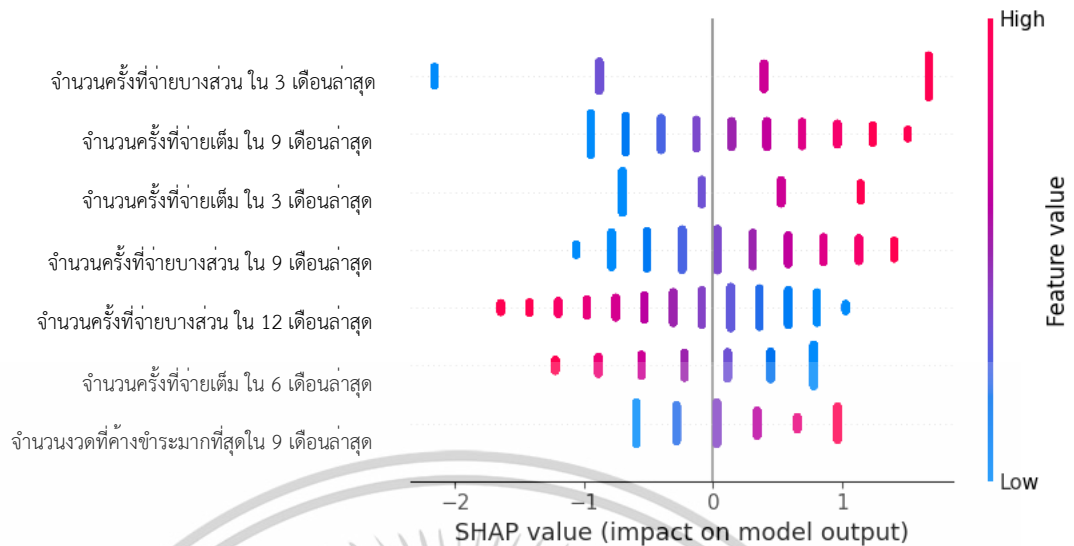
1) แบบจำลอง Logistic Regression

ในการศึกษาปัจจัยที่มีผลต่อการทำนายการผัดขี้หน้าของแบบจำลอง Logistic Regression ได้มีการวิเคราะห์ความสำคัญของตัวแปรโดยใช้ทั้งวิธี Permutation Importance และ SHAP ซึ่งทั้งสองเทคนิคนี้ช่วยให้สามารถตีความเชิงลึกถึงปัจจัยที่มีอิทธิพลมากที่สุดต่อการตัดสินใจของแบบจำลอง



รูปที่ 4.14 ผลจากการทำ Permutation Importance ของแบบจำลอง Logistic Regression

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



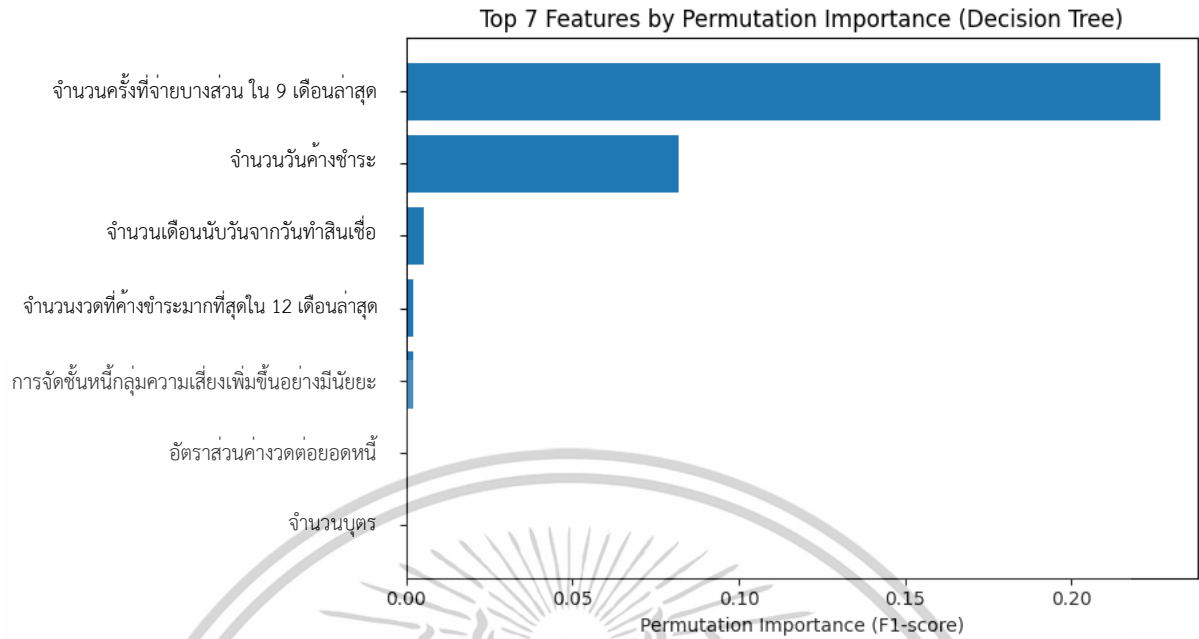
รูปที่ 4.15 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Logistic Regression

จากรูปที่ 4.14 ซึ่งแสดงผลการวิเคราะห์ Permutation Importance ของแบบจำลอง Logistic Regression ซึ่งระบุว่าตัวแปรที่เกี่ยวข้องกับพฤติกรรมการชำระเงินบางส่วนและการชำระเต็มจำนวนในช่วง 3–12 เดือนล่าสุดมีความสำคัญสูงต่อประสิทธิภาพของแบบจำลอง ตัวอย่างเช่น ตัวแปร "จำนวนครั้งที่จ่ายบางส่วนใน 3 เดือนล่าสุด" มีค่า Permutation Importance สูงที่สุด แสดงให้เห็นว่าพฤติกรรมการชำระบางส่วนในช่วงเวลาล่าสุดเป็นปัจจัยที่สำคัญอย่างยิ่งในการแยกแยะกลุ่มลูกค้าที่มีแนวโน้มผิดนัดชำระหนี้ เมื่อพิจารณาพร้อมกับผลการวิเคราะห์ค่า SHAP Value ในรูปที่ 4.15 แสดงให้เห็นผลกระทบของค่าต่างๆ ของตัวแปรที่สำคัญแต่ละตัวต่อผลลัพธ์ของแบบจำลอง โดยแกน X แสดงค่า SHAP value ซึ่งสะท้อนถึงอิทธิพลของตัวแปรที่มีต่อการคาดการณ์ของแบบจำลอง ส่วนสีของจุดข้อมูลแสดงระดับค่าของตัวแปรนั้นๆ โดยสีแดงหมายถึงค่าตัวแปรสูง ในขณะที่สีน้ำเงินหมายถึงค่าตัวแปรต่ำ ตัวอย่างเช่น ในกรณีของตัวแปร “จำนวนครั้งที่จ่ายบางส่วนใน 3 เดือนล่าสุด” หากค่าของฟีเจอร์นี้อยู่ในระดับสูง (สีแดง) จะส่งผลให้ค่าทำนายของแบบจำลองมีแนวโน้มไปในทิศทางที่บ่งชี้ถึงความเสี่ยงในการผิดนัดชำระที่เพิ่มขึ้น ในขณะที่ค่าฟีเจอร์ที่ต่ำ (สีน้ำเงิน) มีแนวโน้มผลักดันค่าทำนายไปในทางตรงข้าม

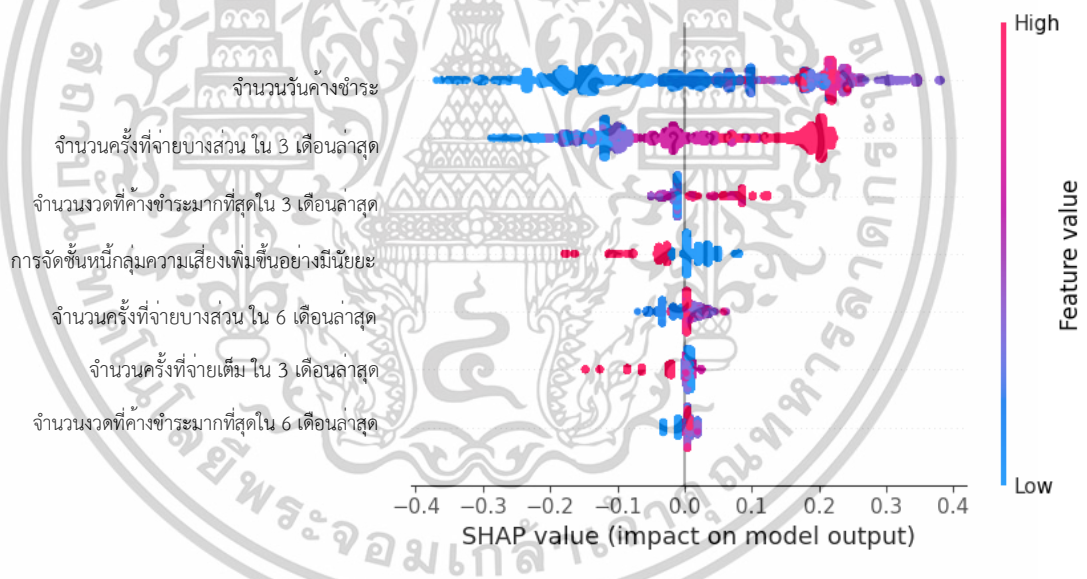
2) แบบจำลอง Decision Tree

ในการศึกษาปัจจัยที่มีผลต่อประสิทธิภาพการทำนายการผิดนัดชำระหนี้ของแบบจำลอง Decision Tree ได้มีการวิเคราะห์ความสำคัญของตัวแปรโดยใช้ทั้งวิธี Permutation Importance และ SHAP ซึ่งทั้งสองเทคนิคนี้ช่วยให้สามารถตีความเชิงลึกถึงปัจจัยที่มีอิทธิพลมากที่สุดต่อการตัดสินใจของแบบจำลอง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.16 ผลจากการทำ Permutation Importance ของแบบจำลอง Decision Tree



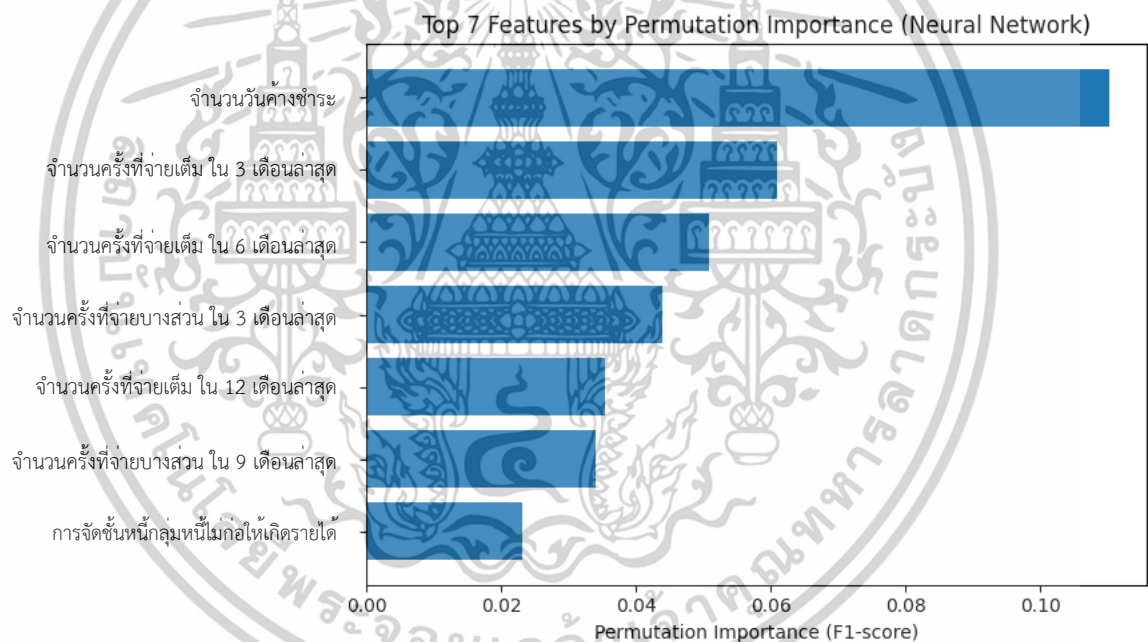
รูปที่ 4.17 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Decision Tree

จากรูปที่ 4.16 ซึ่งแสดงผลการวิเคราะห์ Permutation Importance ของแบบจำลอง Decision Tree พบว่า ตัวแปรที่เกี่ยวข้องกับจำนวนครั้งที่ลูกค้างชำระบางส่วนและจำนวนวันที่ค้างชำระมีอิทธิพลอย่างชัดเจนต่อผลลัพธ์ของแบบจำลอง โดยเฉพาะตัวแปร “จำนวนครั้งที่จ่ายบางส่วน ใน 9 เดือนล่าสุด” ซึ่งมีค่า Permutation Importance สูงที่สุด บ่งชี้ว่าพฤติกรรมการชำระหนี้แบบไม่เต็มจำนวนในช่วง 9 เดือนล่าสุดมีบทบาทสำคัญในการจำแนกลูกค้าที่มีแนวโน้มผิดนัดชำระหนี้ นอกจากนี้ ตัวแปร “จำนวนวันค้างชำระ” ก็ปรากฏเป็นอีกหนึ่งตัวแปรหลักที่มีอิทธิพลสูงต่อผลการทำนายของแบบจำลองเช่นกัน เมื่อพิจารณาร่วมกับผลการวิเคราะห์ค่า SHAP Value ในรูปที่ 4.17 เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ซึ่งแสดงให้เห็นผลกระทบของค่าต่างๆ ของตัวแปรแต่ละตัวต่อผลลัพธ์ของแบบจำลอง โดยแกน X แสดงค่า SHAP value ซึ่งสะท้อนถึงอิทธิพลของตัวแปรที่มีต่อการคาดการณ์ของแบบจำลอง ส่วนสีของจุดข้อมูลแสดงระดับค่าของตัวแปรนั้นๆ โดยสีแดงหมายถึงค่าพีเจอร์สูง ในขณะที่สีน้ำเงินหมายถึงค่าพีเจอร์ต่ำ ตัวอย่างเช่น ในกรณีของตัวแปร “จำนวนวันค้างชำระ” หากมีค่าพีเจอร์สูง (สีแดง) จะมีแนวโน้มผลักดันค่าทำนายไปในทิศทางที่แสดงถึงความเสี่ยงต่อการผิดนัดชำระหนี้ที่สูงขึ้น ขณะที่ค่าพีเจอร์ที่ต่ำ (สีน้ำเงิน) ส่งผลลดความเสี่ยงในการทำนายลง

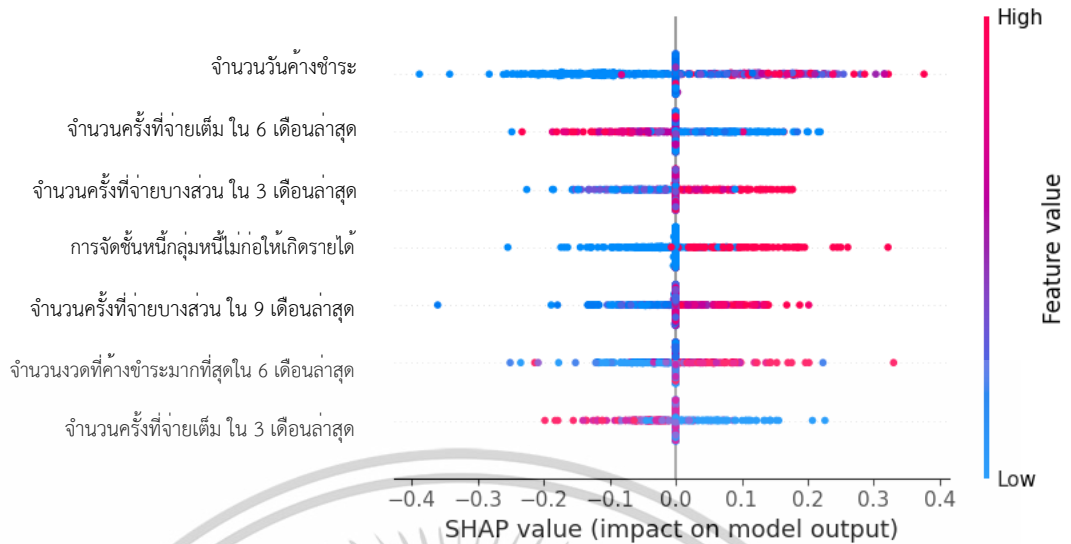
3) แบบจำลอง Neural Network

ในการศึกษาปัจจัยที่มีผลต่อประสิทธิภาพการทำนายการผิดนัดชำระหนี้ของแบบจำลอง Neural Network ได้มีการวิเคราะห์ความสำคัญของตัวแปรโดยใช้ทั้งวิธี Permutation Importance และ SHAP ซึ่งทั้งสองเทคนิคนี้ช่วยให้สามารถตีความเชิงลึกถึงปัจจัยที่มีอิทธิพลมากที่สุดต่อการตัดสินใจของแบบจำลอง



รูปที่ 4.18 ผลจากการทำ Permutation Importance ของแบบจำลอง Neural Network

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



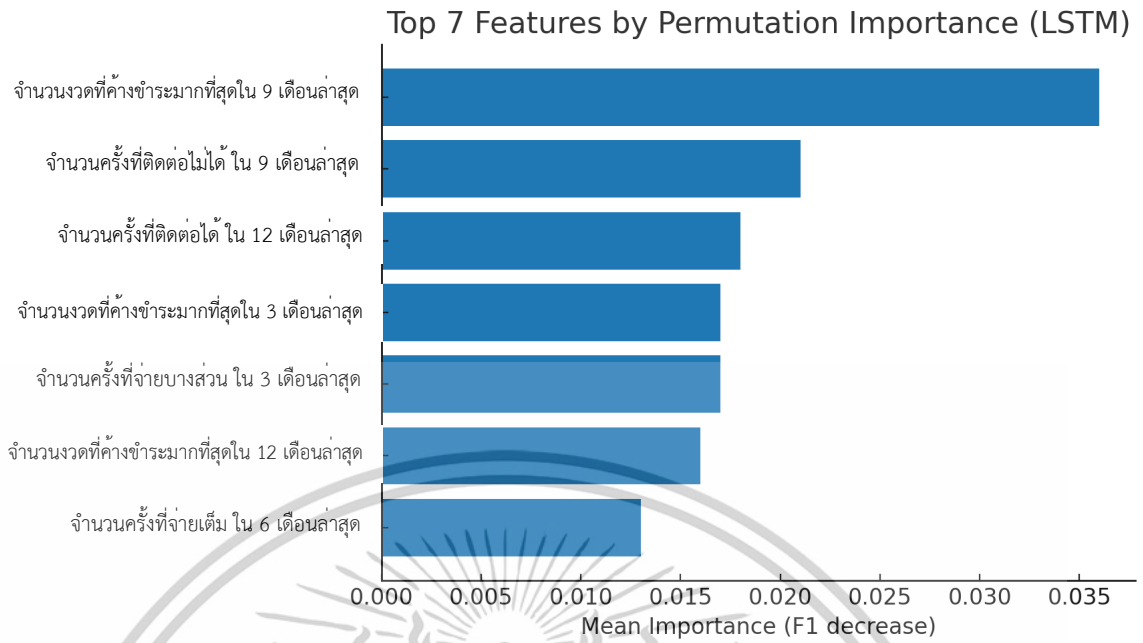
รูปที่ 4.19 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Neural Network

จากรูปที่ 4.18 ซึ่งแสดงผลการวิเคราะห์ Permutation Importance ของแบบจำลอง Neural Network พบว่า ตัวแปรที่สะท้อนถึงพฤติกรรมการค้างชำระและการชำระหนี้บางส่วนหรือเต็มจำนวนในช่วงเวลา 3-12 เดือนล่าสุด มีผลกระทบต่อประสิทธิภาพของแบบจำลอง โดยเฉพาะตัวแปร “จำนวนวันค้างชำระ” ซึ่งมีค่า Permutation Importance สูงที่สุด แสดงให้เห็นว่าจำนวนวันที่ลูกค้าค้างชำระสะสมในอดีตมีบทบาทสำคัญในการพยากรณ์ความเสี่ยงในการผิดนัดชำระหนี้ในอนาคต นอกจากนี้ ตัวแปร “จำนวนครั้งที่จ่ายเต็มใน 3 เดือนล่าสุด” และ “จำนวนครั้งที่จ่ายเต็มใน 6 เดือนล่าสุด” ก็เป็นอีกสองฟีเจอร์ที่มีอิทธิพลสูงต่อแบบจำลอง เมื่อพิจารณาพร้อมกับผลการวิเคราะห์ SHAP Value ในรูปที่ 4.19 ซึ่งแสดงผลกระทบของค่าตัวแปรแต่ละรายการต่อผลลัพธ์ของแบบจำลอง พบว่าแกน X ซึ่งแสดงค่า SHAP Value สะท้อนถึงทิศทางและความแรงของผลกระทบของตัวแปรนั้นๆ ต่อการทำนาย ส่วนสีของจุดข้อมูลแสดงระดับค่าของฟีเจอร์ โดยสีแดงแสดงถึงค่าฟีเจอร์ที่สูง และสีน้ำเงินแสดงค่าฟีเจอร์ที่ต่ำ ตัวอย่างเช่น ในกรณีของตัวแปร “จำนวนวันค้างชำระ” หากมีค่าสูง (สีแดง) จะมีแนวโน้มส่งผลให้ค่าทำนายของแบบจำลองแสดงถึงความเสี่ยงผิดนัดชำระหนี้ที่เพิ่มขึ้น ในขณะที่ค่าต่ำ (สีน้ำเงิน) จะผลักดันค่าทำนายไปในทิศทางตรงกันข้าม นอกจากนี้ยังพบว่าตัวแปร “จำนวนครั้งที่จ่ายบางส่วนใน 3 เดือนล่าสุด” ก็มีลักษณะการกระจายของ SHAP Value ในทิศทางเดียวกัน โดยค่าที่สูงส่งผลเชิงบวกต่อโอกาสผิดนัด

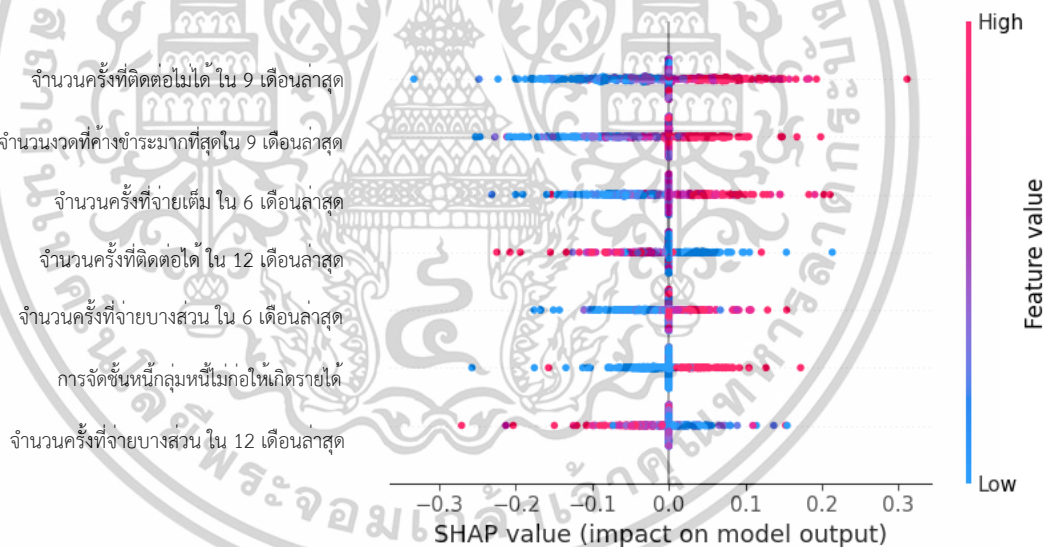
4) แบบจำลอง LSTM

ในการศึกษาปัจจัยที่มีผลต่อประสิทธิภาพการทำนายการผิดนัดชำระหนี้ของแบบจำลอง LSTM ได้มีการวิเคราะห์ความสำคัญของตัวแปรโดยใช้ทั้งวิธี Permutation Importance และ SHAP ซึ่งทั้งสองเทคนิคนี้ช่วยให้สามารถตีความเชิงลึกถึงปัจจัยที่มีอิทธิพลมากที่สุดต่อการตัดสินใจของแบบจำลอง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.20 ผลจากการทำ Permutation Importance ของแบบจำลอง LSTM



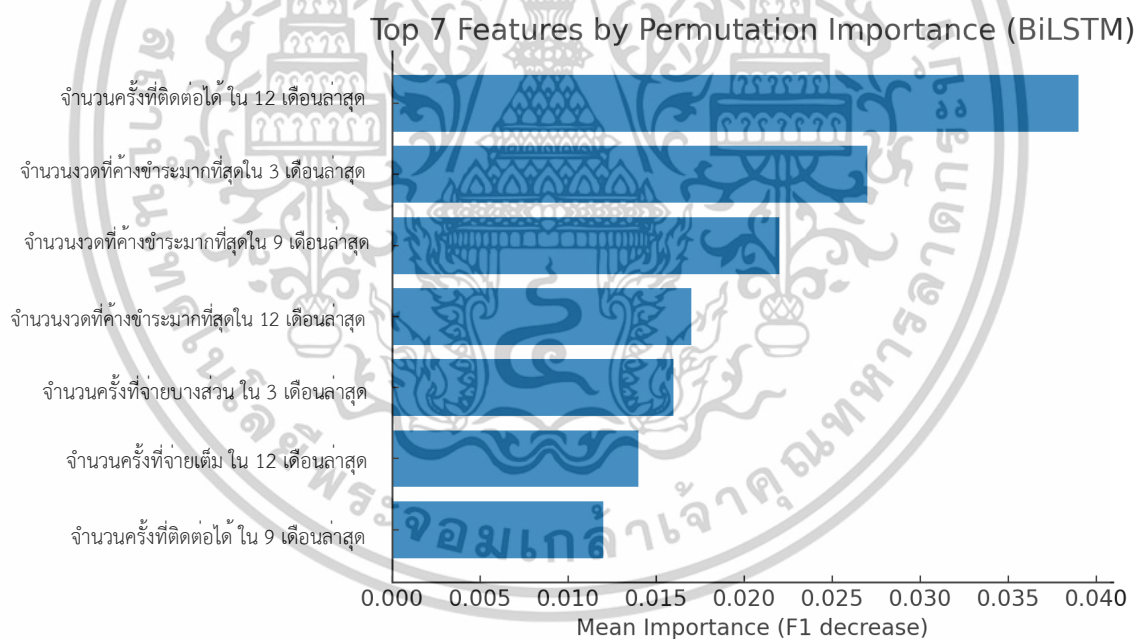
รูปที่ 4.21 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง LSTM

จากรูปที่ 4.20 ซึ่งแสดงผลการวิเคราะห์ Permutation Importance ของแบบจำลอง LSTM พบว่า ตัวแปรที่เกี่ยวข้องกับพฤติกรรมชำระเงินย้อนหลังในลักษณะเชิงลำดับเวลา โดยเฉพาะ “จำนวนงวดที่ทำการชำระมากที่สุดใน 9 เดือนล่าสุด” และ “จำนวนครั้งที่ติดต่อไม่ได้ใน 9 เดือนล่าสุด” เป็นฟีเจอร์ที่มีอิทธิพลต่อประสิทธิภาพของแบบจำลองอย่างมีนัยสำคัญ ตัวแปรเหล่านี้ บ่งชี้ถึงสถานการณ์ความไม่สม่ำเสมอในการชำระหนี้และการขาดการติดต่อกับลูกค้าซึ่งมักสะท้อนถึง ความเสี่ยงที่เพิ่มขึ้นในการผิดนัดชำระหนี้ นอกจากนี้ ตัวแปรอื่นๆ เช่น “จำนวนครั้งที่จ่ายบางส่วนใน เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

3 เดือนล่าสุด” และ “จำนวนครั้งที่ชำระมากที่สุด ใน 12 เดือนล่าสุด” ก็ปรากฏอยู่ในลำดับต้นๆ ของความสำคัญด้วยเช่นกัน ในส่วนของการวิเคราะห์ SHAP Value ตามที่แสดงในรูปที่ 4.21 พบว่า ตัวแปรที่มีค่า SHAP Value สูงสุด ได้แก่ “จำนวนครั้งที่ติดต่อไม่ได้ใน 9 เดือนล่าสุด” ซึ่งเมื่อมีค่าฟีเจอร์สูง (สีแดง) จะมีแนวโน้มผลักดันค่าทำนายของแบบจำลองไปในทิศทางที่สะท้อนถึงความเสี่ยงในการผิดนัดชำระหนี้ที่เพิ่มขึ้น ในขณะที่ค่าตัวแปรต่ำ (สีน้ำเงิน) มีแนวโน้มลดระดับความเสี่ยงที่แบบจำลองคาดการณ์ นอกจากนี้ตัวแปรเช่น “จำนวนครั้งที่จ่ายบางส่วนใน 6 เดือนล่าสุด” และ “การจัดอันดับกลุ่มความเสี่ยง” ก็แสดงผลกระทบที่สอดคล้องกัน คือ ค่าสูงมีแนวโน้มเพิ่มความเสี่ยง และค่าต่ำมีแนวโน้มลดความเสี่ยง

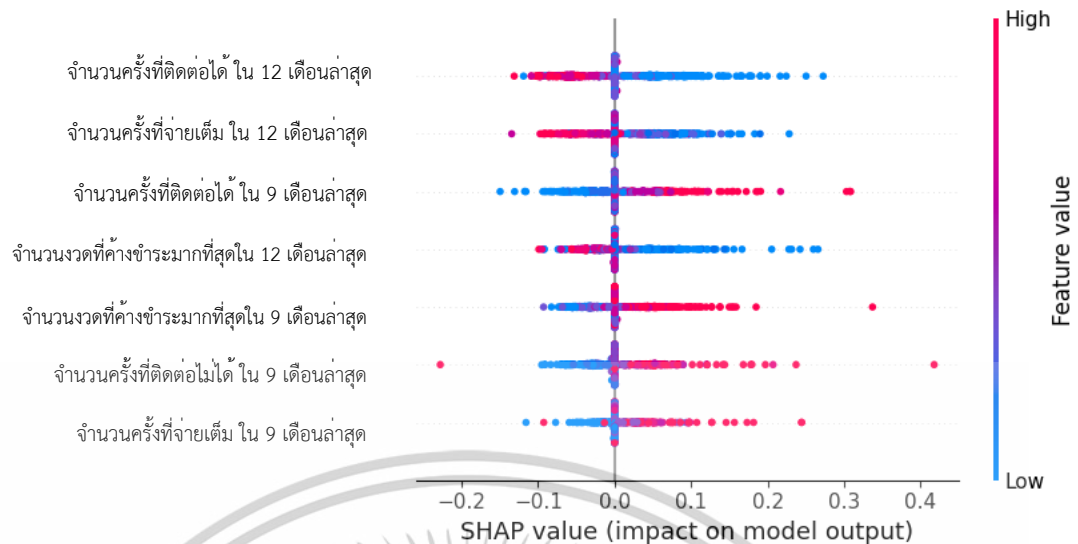
5) แบบจำลอง Bidirectional LSTM

ในการศึกษาปัจจัยที่มีผลต่อประสิทธิภาพการทำนายการผิดนัดชำระหนี้ของแบบจำลอง Bidirectional LSTM ได้มีการวิเคราะห์ความสำคัญของตัวแปรโดยใช้ทั้งวิธี Permutation Importance และ SHAP ซึ่งทั้งสองเทคนิคนี้ช่วยให้สามารถตีความเชิงลึกถึงปัจจัยที่มีอิทธิพลมากที่สุดต่อการตัดสินใจของแบบจำลอง



รูปที่ 4.22 ผลจากการทำ Permutation Importance ของแบบจำลอง Bidirectional LSTM

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้



รูปที่ 4.23 ค่าอิทธิพลของตัวแปรด้วยวิธี SHAP ของแบบจำลอง Bidirectional LSTM

จากรูปที่ 4.22 ซึ่งแสดงผลการวิเคราะห์ Permutation Importance ของแบบจำลอง Bidirectional LSTM พบว่า ตัวแปรที่เกี่ยวข้องกับพฤติกรรมเครดิตและการชำระเงินย้อนหลัง ยังคงมีบทบาทสำคัญสูงต่อผลลัพธ์ของแบบจำลอง โดยเฉพาะตัวแปร “จำนวนครั้งที่ติดต่อได้ใน 12 เดือนล่าสุด” ที่มีค่า Permutation Importance สูงที่สุด แสดงให้เห็นว่าความสามารถในการติดต่อกับลูกค้าอย่างสม่ำเสมอเป็นปัจจัยสำคัญที่ช่วยแยกแยะกลุ่มลูกค้าที่มีความเสี่ยงผิดนัดได้อย่างมีประสิทธิภาพ พีเจอร์สำคัญลำดับถัดมาคือ “จำนวนงวดที่ทำการชำระมากที่สุดใน 3 เดือนล่าสุด” และ “จำนวนงวดที่ทำการชำระมากที่สุดใน 9 เดือนล่าสุด” ซึ่งแสดงถึงความสม่ำเสมอและความเข้มข้นของพฤติกรรมการชำระเงินในช่วงเวลาที่ผ่านไป เมื่อพิจารณาผลการวิเคราะห์ค่า SHAP Value ในรูปที่ 4.23 จะเห็นได้ว่าตัวแปร “จำนวนครั้งที่ติดต่อได้ใน 12 เดือนล่าสุด” มีค่า SHAP ที่กระจายกว้าง โดยค่าสูง (สีแดง) มีแนวโน้มลดความเสี่ยงในการผิดนัดชำระ ขณะที่ค่าต่ำ (สีน้ำเงิน) ผลักดันค่าทำนายไปในทิศทางของความเสี่ยงที่เพิ่มขึ้นในแบบจำลอง นอกจากนี้ ตัวแปร “จำนวนครั้งที่จ่ายเต็มใน 12 เดือนล่าสุด” และ “จำนวนครั้งที่ติดต่อได้ใน 9 เดือนล่าสุด” ก็สะท้อนผลกระทบที่สอดคล้องกัน กล่าวคือ ค่าที่สูงส่งผลเชิงบวกต่อสถานะของลูกค้า ในขณะที่ค่าที่ต่ำอาจเป็นสัญญาณของความเสี่ยง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ตารางที่ 4.16 ผลวิเคราะห์ Permutation Importance และ SHAP Value ของแต่ละแบบจำลอง

แบบจำลอง	ตัวแปรที่ Permutation Importance สูงสุด	ตัวแปรที่ SHAP Value สูงสุด
Logistic Regression	จำนวนครั้งที่ผิดนัดชำระใน 3 เดือน ล่าสุด	จำนวนครั้งที่ผิดนัดชำระใน 3 เดือนล่าสุด
Decision Tree	จำนวนครั้งที่ผิดนัดชำระใน 9 เดือน ล่าสุด	จำนวนครั้งที่ผิดนัดชำระใน 9 เดือนล่าสุด
Neural Network	จำนวนครั้งที่ผิดนัดชำระใน 3 เดือน ล่าสุด	จำนวนครั้งที่ผิดนัดชำระใน 3 เดือนล่าสุด
LSTM	จำนวนครั้งที่ผิดนัดชำระใน 9 เดือน ล่าสุด	จำนวนครั้งที่ผิดนัดชำระใน 9 เดือนล่าสุด
Bidirectional LSTM	จำนวนครั้งที่ผิดนัดชำระใน 12 เดือน ล่าสุด	จำนวนครั้งที่ผิดนัดชำระใน 12 เดือนล่าสุด

จากตารางที่ 4.22 พบว่าแบบจำลองที่ใช้ในงานวิจัยนี้ โดยเฉพาะ Logistic Regression และ Neural Network ให้ความสำคัญกับพฤติกรรมผิดนัดชำระในช่วง 3 เดือนล่าสุดเป็นหลัก ในขณะที่แบบจำลองที่มีศักยภาพในการเรียนรู้ข้อมูลเชิงลำดับเวลามากขึ้น ได้แก่ LSTM และ Bidirectional LSTM มีแนวโน้มให้ความสำคัญกับข้อมูลผิดนัดสะสมในระยะยาว (9-12 เดือนล่าสุด) สะท้อนถึงการนำข้อมูลประวัติย้อนหลังที่ลึกซึ้งมาใช้ประโยชน์ในการประเมินความเสี่ยง

4.4 ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 2 ความเสี่ยงต่ำ (Low Risk)

ในหัวข้อนี้ ผู้วิจัยได้วิเคราะห์กลุ่มลูกค้าในกลุ่มที่ 2 ซึ่งจัดอยู่ในประเภทกลุ่มความเสี่ยงต่ำ (Low Risk) ได้จากการคัดเลือกลูกค้าที่เคยมีประวัติผิดนัดชำระหนี้เพียง 1 ครั้งภายในระยะเวลา 2 เดือนที่ผ่านมา ซึ่งสะท้อนถึงพฤติกรรมความเสี่ยงในระดับเริ่มต้น อาจเกิดจากปัจจัยชั่วคราว เช่น การลืมนำชำระหนี้ โดยมีวัตถุประสงค์เพื่อศึกษาความเป็นไปได้ที่ลูกค้าในกลุ่มนี้จะกลับมาผิดนัดชำระหนี้ซ้ำในเดือนถัดไป เพื่อให้สามารถนำข้อมูลมาใช้ในการวางมาตรการเชิงป้องกันและสนับสนุนวินัยทางการเงินได้อย่างเหมาะสม ก่อนที่จะพัฒนาไปสู่การผิดนัดซ้ำ ลูกค้ากลุ่มนี้มีจำนวนทั้งสิ้น 22,896 บัญชี แบ่งออกเป็น 2 กลุ่มย่อย คือ กลุ่มที่ผิดนัดชำระหนี้ซ้ำ (Default) จำนวน 17,668 บัญชี และกลุ่มที่กลับมาชำระหนี้ปกติ (Non-Default) จำนวน 10,628 บัญชี ซึ่งพบว่าข้อมูลมีลักษณะไม่สมดุล (Imbalanced) ซึ่งส่งผลกระทบต่อความแม่นยำ เพื่อแก้ไขปัญหาดังกล่าว ผู้วิจัยได้นำเทคนิค SMOTE (Synthetic Minority Over-sampling Technique) มาใช้ในการสร้างตัวอย่างข้อมูลสำหรับกลุ่มกลับมาชำระหนี้ปกติ (Non-Default) เพื่อปรับสมดุลของข้อมูลก่อนนำไปพัฒนาแบบจำลอง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

1) การกำหนด Hyperparameter ที่เหมาะสมให้แก่แบบจำลอง

การปรับแต่งค่าพารามิเตอร์ด้วยเทคนิค Grid Search ซึ่งเป็นการค้นหาค่าที่เหมาะสมที่สุดสำหรับแต่ละแบบจำลอง สามารถสรุปได้ดังตารางต่อไปนี้

ตารางที่ 4.17 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลอง ในกลุ่มที่ 2 ความเสี่ยงต่ำ (Low Risk)

แบบจำลอง	ค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters)
Logistic Regression	C = 0.1, penalty = 'l1', solver = 'liblinear', class_weight = None
Decision Tree	criterion = 'gini', max_depth = 5
Neural Network	activation = 'relu', alpha = 2, early_stopping = True, hidden_layer_sizes = (50, 30), learning_rate_init = 0.005
LSTM	batch_size = 16, epochs = 20, activation = 'relu', learning_rate = 0.001, units = 128
Bidirectional LSTM	batch_size = 32, epochs = 50, learning_rate = 0.0005, units = 128

จากตารางที่ 4.17 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองสำหรับกลุ่มความเสี่ยงต่ำ สามารถสังเกตแนวโน้มการเลือกค่าพารามิเตอร์ที่เหมาะสมในแต่ละแบบจำลองได้ดังนี้

Logistic Regression แบบจำลองนี้เลือกใช้ค่า C ในระดับต่ำ (0.1) เพื่อควบคุมความซับซ้อนของแบบจำลองและลดโอกาสการเกิด Overfitting โดยใช้ Penalty แบบ L1 เพื่อช่วยในการคัดเลือกเฉพาะตัวแปรสำคัญ และเลือก Solver แบบ 'liblinear' ที่เหมาะกับ L1 Regularization รวมทั้งไม่ตั้งค่า class_weight เพิ่มเติม เนื่องจากได้ปรับสมดุลข้อมูลด้วย SMOTE แล้ว

Decision Tree กำหนด criterion เป็น 'gini' เพื่อความรวดเร็วในการประมวลผล พร้อมทั้งกำหนด max_depth ที่ 5 เพื่อป้องกันไม่ใ้แบบจำลองมีความซับซ้อนมากเกินไป ลดความเสี่ยงในการเกิด Overfitting และช่วยให้ต้นไม้ตัดสินใจได้อย่างเหมาะสมกับโครงสร้างข้อมูล

Neural Network เลือกใช้ Activation Function แบบ 'relu' เพื่อเพิ่มประสิทธิภาพในการเรียนรู้ความสัมพันธ์เชิงซ้อน โดยกำหนด Alpha ในระดับปานกลางและใช้ early_stopping เพื่อตรวจสอบการหยุดเรียนรู้หากไม่มีพัฒนาการ พร้อมทั้งออกแบบ hidden_layer_sizes เป็นสองชั้น (50, 30) เพื่อรองรับความซับซ้อนของข้อมูล และใช้ learning_rate_init ในระดับที่เหมาะสม

LSTM ตั้งค่า batch_size ที่ 16 เพื่อการอัปเดตน้ำหนักอย่างต่อเนื่องและลดความเสี่ยง Overfitting กำหนดจำนวน Epochs ไม่สูงเกินไป (20) ใช้ Activation แบบ 'relu' รวมถึง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

learning_rate ที่สมดุล (0.001) และ Units จำนวน 128 เพื่อรองรับการจับลักษณะลำดับของข้อมูลได้เพียงพอ

Bidirectional LSTM เลือกใช้ batch_size ที่สูงขึ้น (32) และ Epochs สูง (50) เพื่อให้สอดคล้องกับความซับซ้อนของข้อมูลและรองรับการเรียนรู้ทั้งจากอดีตและอนาคต กำหนด learning_rate ที่ต่ำ (0.0005) เพื่อลดความเสี่ยงจากการอัปเดตน้ำหนักแบบก้าวกระโดด และกำหนด Units เท่ากับ LSTM เพื่อเปรียบเทียบประสิทธิภาพโดยตรง

2) ผลการทดลองและการเปรียบเทียบประสิทธิภาพแบบจำลอง

ตารางที่ 4.18 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 2 ความเสี่ยงต่ำ (Low Risk)

แบบจำลอง	ค่าความเที่ยงตรง	ค่าความแม่นยำ	ค่าความระลึก	ค่าประสิทธิภาพโดยรวม
Logistic Regression	0.7984	0.8161	0.7984	0.8013
Decision Tree	0.7966	0.8177	0.7966	0.7977
Neural Network	0.7978	0.8169	0.7978	0.8008
LSTM	0.8008	0.8178	0.8008	0.8037
Bidirectional LSTM	0.8024	0.8147	0.8024	0.8049

ผลการทดสอบกับข้อมูลชุดทดสอบ (Test Set) พบว่าแบบจำลอง Bidirectional LSTM ได้ผลลัพธ์ที่ดีที่สุด โดยให้ค่าความระลึก (Recall) เท่ากับร้อยละ 80.24 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับ ร้อยละ 80.49 รองลงมาคือแบบจำลอง LSTM ที่มีค่าความระลึก (Recall) เท่ากับ ร้อยละ 80.08 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับ ร้อยละ 80.37 ขณะที่แบบจำลองอื่นๆ ให้ค่าความระลึก (Recall) และค่าประสิทธิภาพโดยรวม (F1-Score) อยู่ในช่วงประมาณ ร้อยละ 79.66 - 79.84 และ ร้อยละ 79.77 - 80.13 ตามลำดับ สรุปได้ว่าแบบจำลอง Bidirectional LSTM ได้ผลการทำนายที่ดีที่สุดในกลุ่มลูกค้าที่เคยผิณฑ์ชำระหนี้เพียงครั้งเดียวในช่วง 2 เดือนที่ผ่านมา ทั้งในแง่ของความแม่นยำและความสามารถในการตรวจจับกลุ่มที่มีแนวโน้มจะผิณฑ์ชำระหนี้ซ้ำ

4.5 ผลการทำนายการผิณฑ์ชำระหนี้ในกลุ่มที่ 3 ความเสี่ยงปานกลาง (Moderate Risk)

ในหัวข้อนี้ ผู้วิจัยได้ดำเนินการวิเคราะห์กลุ่มลูกค้าในกลุ่มที่ 3 ซึ่งจัดอยู่ในประเภทกลุ่มความเสี่ยงปานกลาง (Moderate Risk) โดยกลุ่มนี้ได้จากการคัดเลือกลูกค้าที่มีพฤติกรรมผิณฑ์ชำระหนี้ในเดือนล่าสุด ซึ่งเป็นสัญญาณของความเสี่ยงที่อาจเกิดขึ้นต่อเนื่องในระยะเวลาอันใกล้ พฤติกรรมเอกสารนี้เป็นเอกสารที่ส่งวนเวียนสำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปเผยแพร่โดยไม่ผ่านการอนุมัติทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ดังกล่าวสะท้อนถึงภาวะความไม่มั่นคงทางการเงินในช่วงเวลาปัจจุบัน ซึ่งอาจมีสาเหตุมาจากความสามารถในการชำระหนี้ที่ลดลง ความไม่แน่นอนทางรายได้ หรือปัจจัยทางเศรษฐกิจอื่นๆ ที่ส่งผลกระทบต่อสภาพคล่องของลูกค้า เนื่องจากพฤติกรรมผิดนัดชำระหนี้เกิดขึ้นล่าสุดอาจทำให้ไม่สามารถจัดการภาระหนี้ได้อย่างเหมาะสมในเดือนถัดไป จึงมีความจำเป็นอย่างยิ่งที่จะต้องใช้แบบจำลองที่สามารถเรียนรู้ลำดับเหตุการณ์ได้ดีในการคาดการณ์ ลูกค้าในกลุ่มนี้มีทั้งหมด 16,451 บัญชี โดยจำแนกออกเป็น 2 กลุ่ม ได้แก่ กลุ่มที่ผิดนัดชำระหนี้อีกครั้ง (Default) มีจำนวนลูกค้าที่ผิดนัดชำระซ้ำจำนวน 10,628 บัญชี และกลุ่มลูกค้าที่สามารถกลับมาชำระหนี้ได้ตามปกติ (Non-Default) จำนวน 5,823 บัญชี ซึ่งพบว่าข้อมูลมีลักษณะไม่สมดุล (Imbalanced) ซึ่งอาจส่งผลกระทบต่อความแม่นยำของแบบจำลอง เพื่อแก้ไขปัญหาดังกล่าว จึงได้นำเทคนิค SMOTE (Synthetic Minority Over-sampling Technique) มาประยุกต์ใช้ในการสร้างตัวอย่างข้อมูลสำหรับกลุ่มกลับมาชำระหนี้ปกติ (Non-Default) เพื่อปรับสมดุลของข้อมูลก่อนนำไปพัฒนาแบบจำลอง ส่งผลให้แบบจำลองสามารถเรียนรู้ลักษณะของกลุ่มชำระหนี้ปกติได้ดียิ่งขึ้น

1) การกำหนด Hyperparameter ที่เหมาะสมให้แก่แบบจำลอง

การปรับแต่งค่าพารามิเตอร์ด้วยเทคนิค Grid Search ซึ่งเป็นการค้นหาค่าที่เหมาะสมที่สุดสำหรับแต่ละแบบจำลอง สามารถสรุปได้ดังตารางต่อไปนี้

ตารางที่ 4.19 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองในกลุ่มที่ 3 ความเสี่ยงปานกลาง (Moderate Risk)

แบบจำลอง	ค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters)
Logistic Regression	{'C': 100, 'class_weight': 'balanced', 'penalty': 'l1', 'solver': 'liblinear'}
Decision Tree	{'criterion': 'gini', 'max_depth': 5}
Neural Network	{'activation': 'relu', 'alpha': 2, 'early_stopping': True, 'hidden_layer_sizes': (30,),'learning_rate_init': 0.001}
LSTM	{'batch_size': 32, 'epochs': 50, 'model__learning_rate': 0.0005, 'model__units': 256}
Bidirectional LSTM	{'batch_size': 32, 'epochs': 50, 'model__learning_rate': 0.001, 'model__units': 256}

จากตารางที่ 4.19 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองสำหรับกลุ่มความเสี่ยงปานกลาง เลือกค่าพารามิเตอร์ที่เหมาะสมได้ดังนี้

Logistic Regression แบบจำลองนี้เลือกใช้ค่า C สูง (100) เพื่อเพิ่มความยืดหยุ่นในการสร้างเส้นแบ่งเขตการตัดสินใจ พร้อมทั้งกำหนด penalty แบบ L1 เพื่อช่วยในการเลือกตัวแปรสำคัญ เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

และใช้ Solver 'liblinear' ซึ่งเหมาะสมกับข้อมูลขนาดกลาง นอกจากนี้ยังตั้งค่า class_weight เป็น 'balanced' เพื่อชดเชยความไม่สมดุลของข้อมูลระหว่างกลุ่ม Default และ Non-Default

Decision Tree กำหนด Criterion เป็น 'gini' เพื่อประสิทธิภาพในการคำนวณ และกำหนด max_depth ที่ 5 เพื่อควบคุมความซับซ้อนของโครงสร้างต้นไม้ ช่วยลดความเสี่ยงจากการเกิด Overfitting และรักษาความสามารถในการจำแนกกลุ่มลูกค้าได้อย่างเหมาะสม

Neural Network เลือกใช้ Activation Function แบบ 'relu' เพื่อเพิ่มประสิทธิภาพในการเรียนรู้ความสัมพันธ์ที่ซับซ้อน โดยกำหนด Alpha เป็น 2 เพื่อป้องกัน Overfitting ใช้ early_stopping เพื่อตรวจสอบและหยุดการเรียนรู้เมื่อไม่มีการปรับปรุงประสิทธิภาพ โครงสร้าง hidden_layer_sizes กำหนดเป็นชั้นเดียวขนาด 30 และ learning_rate_init ที่ 0.001 เพื่อความสมดุลในการอัปเดตน้ำหนัก

LSTM ตั้งค่า batch_size ที่ 32 และ Epochs สูง (50) เพื่อรองรับการเรียนรู้ข้อมูลลำดับเวลาได้ต่อเนื่องและละเอียด กำหนด learning_rate ต่ำ (0.0005) เพื่อควบคุมการอัปเดตน้ำหนักให้ค่อยเป็นค่อยไป และใช้ Units ขนาดใหญ่ (256) เพื่อเพิ่มศักยภาพในการเรียนรู้ลำดับข้อมูล

Bidirectional LSTM ใช้ batch_size เท่ากับ LSTM (32) และ Epochs สูง (50) เพื่อรองรับการเรียนรู้ที่ซับซ้อนทั้งจากอดีตและอนาคต กำหนด learning_rate ที่ 0.001 และ Units ขนาด 256 เช่นเดียวกับ LSTM เพื่อเสริมความสามารถในการจับลักษณะเฉพาะของข้อมูลลำดับเวลา

2) ผลการทดลองและการเปรียบเทียบประสิทธิภาพแบบจำลอง

ตารางที่ 4.20 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 3 ความเสี่ยงปานกลาง (Moderate Risk)

แบบจำลอง	ค่าความเที่ยงตรง	ค่าความแม่นยำ	ค่าความระลึก	ค่าประสิทธิภาพโดยรวม
Logistic Regression	0.7935	0.8362	0.7935	0.8031
Decision Tree	0.7933	0.8344	0.7933	0.8027
Neural Network	0.7902	0.8364	0.7902	0.8003
LSTM	0.7981	0.8376	0.7981	0.8072
Bidirectional LSTM	0.8016	0.8347	0.8016	0.8099

ผลการทดสอบกับข้อมูลชุดทดสอบ (Test Set) พบว่าแบบจำลอง Bidirectional LSTM ได้ผลลัพธ์ที่ดีที่สุด โดยให้ค่าความระลึก (Recall) เท่ากับร้อยละ 80.16 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 80.99 รองลงมาคือ LSTM ที่มีค่าความระลึก (Recall) เท่ากับร้อยละ

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

79.81 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 80.72 ส่วนแบบจำลองอื่น เช่น Decision Tree และ Logistic Regression ให้ค่าความระลึก (Recall) ประมาณ ร้อยละ 79.33 - 79.35 และค่าประสิทธิภาพโดยรวม (F1-Score) ใกล้เคียงกันที่ประมาณ ร้อยละ 80.27 - 80.31 สรุปได้ว่า แบบจำลอง Bidirectional LSTM เหมาะสมอย่างยิ่งสำหรับการนำไปใช้ในการประเมินความเสี่ยงของลูกค้ากลุ่มนี้ เนื่องจากสามารถคาดการณ์แนวโน้มผิดนัดชำระหนี้ในระยะใกล้ได้อย่างแม่นยำ และสามารถนำไปประยุกต์ใช้ในการวางแผนกลยุทธ์เชิงป้องกัน เพื่อป้องกันการเข้าสู่กลุ่มผิดนัดชำระหนี้ถาวรในอนาคต

4.6 ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 4 ความเสี่ยงสูง (High Risk)

ในหัวข้อนี้ ผู้วิจัยได้ดำเนินการวิเคราะห์กลุ่มลูกค้าในกลุ่มที่ 4 ซึ่งจัดอยู่ในประเภทกลุ่มความเสี่ยงสูง (High Risk) โดยกลุ่มนี้ได้จากการคัดเลือกลูกค้าที่มีพฤติกรรมผิดนัดชำระหนี้อย่างต่อเนื่องเป็นระยะเวลา 2 เดือนติดต่อกัน พฤติกรรมดังกล่าวสะท้อนถึงภาวะความเสี่ยงที่อยู่ในระดับค่อนข้างสูง และแสดงให้เห็นถึงแนวโน้มที่ลูกค้าอาจไม่สามารถฟื้นกลับมาชำระหนี้ได้ตามปกติในระยะสั้น ซึ่งอาจเกิดจากปัญหาทางการเงินที่สะสมต่อเนื่อง หรือความสามารถในการชำระหนี้ที่ลดลงและมีแนวโน้มที่จะผิดนัดชำระหนี้ซ้ำในอนาคตหากไม่ได้รับการจัดการหรือสนับสนุนอย่างเหมาะสม ลูกค้าในกลุ่มนี้มีทั้งหมด 15,779 บัญชี แบ่งเป็นกลุ่มที่ผิดนัดชำระหนี้ซ้ำ (Default) จำนวน 13,196 บัญชี และกลุ่มที่สามารถกลับมาชำระหนี้ได้ (Non-Default) จำนวน 2,583 บัญชี ซึ่งอัตราส่วน (Ratio) นี้ชี้ให้เห็นถึงอัตราการผิดนัดชำระหนี้ที่ค่อนข้างสูง คิดเป็นประมาณร้อยละ 83.64 ของกลุ่มลูกค้าทั้งหมด ซึ่งพบว่าข้อมูลมีลักษณะไม่สมดุล (Imbalanced) ซึ่งอาจส่งผลกระทบต่อความแม่นยำของแบบจำลองในการจำแนกกลุ่มลูกค้าความเสี่ยง เพื่อแก้ไขปัญหาดังกล่าว ผู้วิจัยได้นำเทคนิค SMOTE (Synthetic Minority Over-sampling Technique) มาประยุกต์ใช้ในการสร้างตัวอย่างข้อมูลสำหรับกลุ่มกลับมาชำระหนี้ปกติ (Non-Default) เพื่อปรับสมดุลของข้อมูลก่อนนำไปพัฒนาแบบจำลอง

1) การกำหนด Hyperparameter ที่เหมาะสมให้แก่แต่ละแบบจำลอง

การปรับแต่งค่าพารามิเตอร์ด้วยเทคนิค Grid Search ซึ่งเป็นการค้นหาค่าที่เหมาะสมที่สุดสำหรับแต่ละแบบจำลอง โดยรายละเอียดของค่าพารามิเตอร์ที่เลือกใช้ในกลุ่มนี้ สามารถสรุปได้ดังตารางต่อไปนี้

ตารางที่ 4.21 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลอง ในกลุ่มที่ 4 ความเสี่ยงสูง (High Risk)

แบบจำลอง	ค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters)
Logistic Regression	{'C': 120, 'class_weight': None, 'penalty': 'l2', 'solver': 'lbfgs'}
Decision Tree	{'criterion': 'gini', 'max_depth': 7}
Neural Network	{'activation': 'relu', 'alpha': 2, 'early_stopping': True, 'hidden_layer_sizes': (50, 30), 'learning_rate_init': 0.005}
LSTM	{'batch_size': 32, 'epochs': 50, 'model__learning_rate': 0.001, 'model__units': 128}
Bidirectional LSTM	{'batch_size': 32, 'epochs': 30, 'model__learning_rate': 0.001, 'model__units': 256}

จากตารางที่ 4.21 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองสำหรับกลุ่มความเสี่ยงสูง สามารถสังเกตแนวโน้มการเลือกค่าพารามิเตอร์ที่เหมาะสมได้ดังนี้

Logistic Regression แบบจำลองนี้เลือกใช้ค่า C สูง (120) เพื่อเพิ่มความยืดหยุ่นในการหาขอบเขตการจำแนกกลุ่มข้อมูล โดยใช้ Penalty แบบ L2 และ Solver 'lbfgs' ที่เหมาะสมกับข้อมูลที่มีขนาดใหญ่ขึ้นและซับซ้อน

Decision Tree กำหนด Criterion เป็น 'gini' เพื่อความรวดเร็วในการคำนวณ และเพิ่ม max_depth เป็น 7 เพื่อให้สามารถจับโครงสร้างที่ซับซ้อนมากขึ้นในข้อมูลกลุ่มนี้ ช่วยให้การตัดสินใจของแบบจำลองมีความละเอียดและแม่นยำยิ่งขึ้น

Neural Network เลือกใช้ Activation Function แบบ 'relu' เพื่อประสิทธิภาพในการเรียนรู้ความสัมพันธ์เชิงซ้อน ตั้งค่า Alpha ที่ 2 เพื่อควบคุม Regularization ใช้ early_stopping เพื่อตรวจสอบการหยุดเรียนรู้หากไม่มีการปรับโครงสร้าง hidden_layer_sizes

ถูกออกแบบเป็นสองชั้น (50, 30) สำหรับรองรับความซับซ้อนของข้อมูล และ learning_rate_init ที่ 0.005 เพื่อความเหมาะสมในการอัปเดตน้ำหนัก

LSTM ตั้งค่า batch_size ที่ 32 และ Epochs สูง (50) เพื่อรองรับการเรียนรู้ข้อมูลลำดับเวลาที่ต่อเนื่อง กำหนด learning_rate ที่ 0.001 และ Units ขนาด 128 เพื่อช่วยให้แบบจำลองจับ pattern เชิงลำดับในกลุ่มข้อมูลที่มีความซับซ้อนระดับกลางได้ดี

Bidirectional LSTM ใช้ batch_size เท่ากับ LSTM (32) แต่กำหนด Epochs ที่ 30 ให้เหมาะสมกับขนาดและโครงสร้างของข้อมูล กำหนด learning_rate ที่ 0.001 และ Units ขนาด 256 เพื่อเพิ่มศักยภาพในการเรียนรู้ข้อมูลลำดับเวลาแบบสองทิศทาง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

2) ผลการทดลองและการเปรียบเทียบประสิทธิภาพแบบจำลอง

ตารางที่ 4.22 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 4 ความเสี่ยงสูง (High Risk)

แบบจำลอง	ค่าความเที่ยงตรง	ค่าความแม่นยำ	ค่าความระลึก	ค่าประสิทธิภาพโดยรวม
Logistic Regression	0.7724	0.8663	0.7724	0.7985
Decision Tree	0.8054	0.8639	0.8054	0.8237
Neural Network	0.7788	0.8649	0.7788	0.8033
LSTM	0.7994	0.8612	0.7994	0.8187
Bidirectional LSTM	0.8165	0.8638	0.8165	0.8320

ผลการทดสอบกับข้อมูลชุดทดสอบ (Test Set) พบว่าแบบจำลอง Bidirectional LSTM ได้ผลลัพธ์ที่ดีที่สุด โดยมีค่าความระลึก (Recall) เท่ากับร้อยละ 81.65 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 83.20 รองลงมาคือแบบจำลอง LSTM ที่มีค่าความระลึก (Recall) เท่ากับร้อยละ 79.94 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 81.87 ขณะที่แบบจำลอง Decision Tree ให้ค่าความระลึก (Recall) เท่ากับร้อยละ 80.54 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 82.37 ซึ่งอยู่ในระดับใกล้เคียงกัน ส่วนแบบจำลอง Neural Network และ Logistic Regression มีค่าความระลึกต่ำกว่าที่ร้อยละ 77.88 และ 77.24 ตามลำดับ สรุปได้ว่าแบบจำลอง Bidirectional LSTM เหมาะสมอย่างยิ่งสำหรับการนำไปใช้เป็นเครื่องมือประกอบการตัดสินใจในการบริหารความเสี่ยงด้านสินเชื่อของสถาบันการเงินสำหรับกลุ่มลูกค้าที่ผิดนัดชำระหนี้ต่อเนื่อง 2 เดือน

4.7 ผลการทำนายการผิดนัดชำระหนี้ในกลุ่มที่ 5 ความเสี่ยงวิกฤต (Critical Risk)

ในหัวข้อนี้ ผู้วิจัยได้ดำเนินการวิเคราะห์กลุ่มลูกค้าในกลุ่มที่ 5 ซึ่งจัดอยู่ในประเภทกลุ่มความเสี่ยงวิกฤต (Critical Risk) ได้จากการคัดเลือกลูกค้าที่มีพฤติกรรมผิดนัดชำระหนี้อย่างต่อเนื่องเป็นระยะเวลา 3 เดือน ซึ่งถือเป็นกลุ่มที่มีความเสี่ยงสูงที่สุดในบรรดากลุ่มลูกค้าทั้งหมด สะท้อนถึงปัญหาทางการเงินที่มีความรุนแรงและต่อเนื่อง ทั้งยังบ่งชี้ถึงภาวะความเสี่ยงเชิงโครงสร้างที่อาจยากต่อการฟื้นตัวในระยะสั้น ซึ่งมีความสำคัญอย่างยิ่งต่อการบริหารจัดการความเสี่ยงเชิงนโยบายของสถาบันการเงิน ทั้งในด้านการติดตามอย่างใกล้ชิด การวางมาตรการช่วยเหลืออย่างเร่งด่วน ลูกค้ากลุ่มนี้มีจำนวนทั้งสิ้น 11,532 บัญชี แบ่งออกเป็น 2 กลุ่ม ได้แก่ กลุ่มที่ผิดนัดชำระหนี้ซ้ำ (Default) จำนวน 10,569 บัญชี และกลุ่มที่สามารถกลับมาชำระหนี้ได้ตามปกติ (Non-Default) จำนวน 963 บัญชี ซึ่ง

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

พบว่าข้อมูลมีลักษณะไม่สมดุล (Imbalanced) อย่างชัดเจน กล่าวคือ กลุ่มที่ผิคนัดซ้ำมีจำนวนมากกว่ากลุ่มที่กลับมาชำระหนี้ตามปกติอย่างมีนัยสำคัญ ซึ่งอาจส่งผลกระทบต่อประสิทธิภาพในการจำแนกกลุ่มเสี่ยงของแบบจำลอง ดังนั้น ผู้วิจัยจึงนำเทคนิค SMOTE (Synthetic Minority Over-sampling Technique) มาประยุกต์ใช้เพื่อสร้างตัวอย่างข้อมูลเทียมในกลุ่ม Non-Default ให้มีจำนวนสมดุลกับกลุ่ม Default ส่งผลให้แบบจำลองสามารถเรียนรู้พฤติกรรมของลูกค้ำทั้งสองกลุ่มได้อย่างมีประสิทธิภาพมากขึ้น

1) การกำหนด Hyperparameter ที่เหมาะสมให้แก่แต่ละแบบจำลอง

การปรับแต่งค่าพารามิเตอร์ด้วยเทคนิค Grid Search ซึ่งเป็นการค้นหาค่าที่เหมาะสมที่สุดสำหรับแต่ละแบบจำลอง สามารถสรุปได้ดังตารางต่อไปนี้

ตารางที่ 4.23 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองในกลุ่มที่ 5 กลุ่มความเสี่ยงวิกฤต (Critical Risk)

แบบจำลอง	ค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters)
Logistic Regression	{'C': 0.05, 'class_weight': None, 'penalty': 'l2', 'solver': 'liblinear'}
Decision Tree	{'criterion': 'gini', 'max_depth': 7}
Neural Network	{'activation': 'relu', 'alpha': 3, 'early_stopping': True, 'hidden_layer_sizes': (50, 50), 'learning_rate_init': 0.001}
LSTM	{'batch_size': 32, 'epochs': 30, 'model__learning_rate': 0.001, 'model__units': 128}
Bidirectional LSTM	{'batch_size': 32, 'epochs': 30, 'model__learning_rate': 0.001, 'model__units': 128}

จากตารางที่ 4.23 สรุปค่าพารามิเตอร์ที่เหมาะสมที่สุด (Best Hyperparameters) ของแต่ละแบบจำลองสำหรับกลุ่มความเสี่ยงวิกฤต สามารถสังเกตแนวโน้มการเลือกค่าพารามิเตอร์ได้ดังนี้

Logistic Regression แบบจำลองนี้เลือกใช้ค่า C ที่ค่อนข้างต่ำ (0.05) เพื่อควบคุมไม่ให้แบบจำลองเกิด Overfitting โดยใช้ Penalty แบบ L2 และ Solver 'liblinear' ซึ่งเหมาะกับข้อมูลขนาดใหญ่ ส่วน class_weight ไม่ได้ตั้งค่าเพิ่มเติม เนื่องจากลักษณะของข้อมูลกลุ่มนี้ไม่จำเป็นต้องชดเชยความไม่สมดุล

Decision Tree กำหนด Criterion เป็น 'gini' เพื่อความรวดเร็วในการประมวลผล และเพิ่มค่า max_depth เป็น 7 เพื่อให้สามารถแบ่งแยกข้อมูลกลุ่มที่มีความซับซ้อนสูงอย่างต่อเนื่องได้อย่างมีประสิทธิภาพ โดยยังคงรักษาสมดุลระหว่างความแม่นยำและการป้องกัน Overfitting

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

Neural Network เลือกใช้ Activation Function แบบ 'relu' เพื่อการเรียนรู้ที่มีประสิทธิภาพ กำหนดค่า Alpha ที่ 3 เพื่อควบคุม Regularization พร้อมทั้งใช้ early_stopping เพื่อลดโอกาสเกิด Overfitting โครงสร้าง hidden_layer_sizes ออกแบบเป็นสองชั้น (50, 50) เพื่อรองรับความซับซ้อนของข้อมูล และใช้ learning_rate_init ที่ 0.001 เพื่อให้การปรับมีความเสถียร

LSTM ตั้งค่า batch_size ที่ 32 และ Epochs ที่ 30 เพื่อความสมดุลระหว่างการเรียนรู้และการป้องกัน Overfitting กำหนด learning_rate ที่ 0.001 และ Units ขนาด 128 เพื่อเสริมศักยภาพในการจับลักษณะลำดับของข้อมูลที่ต่อเนื่องกันหลายเดือน

Bidirectional LSTM ตั้งค่า batch_size และ epochs เท่ากับ LSTM (32, 30) พร้อมทั้ง learning_rate ที่ 0.001 และ Units ขนาด 128 เช่นเดียวกัน เพื่อเปรียบเทียบประสิทธิภาพการเรียนรู้ข้อมูล Sequence ทั้งสองทิศทางได้อย่างมีประสิทธิภาพ

2) ผลการทดลองและการเปรียบเทียบประสิทธิภาพแบบจำลอง

ตารางที่ 4.24 ผลการทำนายของข้อมูลชุดทดสอบในกลุ่มที่ 5 กลุ่มความเสี่ยงวิกฤต (Critical Risk)

แบบจำลอง	ค่าความเที่ยงตรง	ค่าความแม่นยำ	ค่าความระลึก	ค่าประสิทธิภาพโดยรวม
Logistic Regression	0.7927	0.9224	0.7927	0.8352
Decision Tree	0.8218	0.9179	0.8218	0.8550
Neural Network	0.8352	0.9161	0.8352	0.8639
LSTM	0.8604	0.9162	0.8604	0.8810
Bidirectional LSTM	0.8556	0.9195	0.8556	0.8785

ผลการทดสอบกับข้อมูลชุดทดสอบ (Test Set) พบว่าแบบจำลอง LSTM ได้ผลลัพธ์ที่ดีที่สุด โดยมีค่าความระลึก (Recall) เท่ากับร้อยละ 86.04 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 88.10 รองลงมาคือแบบจำลอง Bidirectional LSTM ที่มีค่าความระลึก (Recall) เท่ากับร้อยละ 85.56 และค่าประสิทธิภาพโดยรวม (F1-Score) เท่ากับร้อยละ 87.85

สรุปได้ว่าแบบจำลอง LSTM มีความเหมาะสมอย่างยิ่งสำหรับการนำไปใช้เป็นเครื่องมือประกอบการตัดสินใจในการบริหารความเสี่ยงด้านสินเชื่อของสถาบันการเงิน โดยเฉพาะในกลุ่มลูกค้าที่มีพฤติกรรมผิดนัดชำระหนี้อย่างต่อเนื่องเป็นระยะเวลา 3 เดือน ซึ่งจัดเป็นกลุ่มความเสี่ยงวิกฤต (Critical Risk) หรือและกำลังจะเป็น NPL (Non-Performing Loan) เนื่องจากแบบจำลอง LSTM ให้ผลลัพธ์ที่ดีที่สุดในแง่ของค่าความระลึก (Recall) และค่าประสิทธิภาพโดยรวม (F1-Score) เมื่อเปรียบเทียบกับแบบจำลอง แม้ว่าแบบจำลอง Bidirectional LSTM จะมีโครงสร้างที่สามารถเรียนรู้

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ข้อมูลได้ทั้งจากอดีตและอนาคต แต่ในกรณีของกลุ่มที่มีพฤติกรรมผิวนัดต่อเนื่อง ซึ่งมีลักษณะลำดับข้อมูลที่ เป็นไปในทิศทางเดียว (จากอดีตสู่ปัจจุบัน) แบบจำลอง LSTM ที่เน้นการเรียนรู้ตามลำดับเหตุการณ์จริงจึงสามารถจับรูปแบบพฤติกรรมได้อย่างมีประสิทธิภาพมากกว่า อีกทั้งยังลดความซับซ้อนของโมเดลลง ทำให้สามารถเรียนรู้ได้ดีขึ้นภายใต้เงื่อนไขของข้อมูลเฉพาะกลุ่มนี้

4.8 อภิปรายผล

จากผลการศึกษาพบว่าแบบจำลอง LSTM และ Bidirectional LSTM ที่เป็นอัลกอริทึมการเรียนรู้เชิงลึก มีประสิทธิภาพสูงกว่าแบบจำลอง Logistic Regression, Decision Tree และ Neural Network ที่เป็นอัลกอริทึมการเรียนรู้ของเครื่อง โดยพิจารณาจากค่าความระลึก (Recall) และค่าประสิทธิภาพโดยรวม (F1-score) ในทุกกลุ่มลูกค้าที่มีพฤติกรรมผิวนัดชำระหนี้ที่แตกต่างกัน โดยเฉพาะอย่างยิ่งในกลุ่มลูกค้าที่มีความเสี่ยงสูงและกลุ่มความเสี่ยงวิกฤตที่มีค่าประสิทธิภาพสูงกว่าแบบจำลองทั้ง 3 แบบ อย่างเห็นได้ชัด

ผลการศึกษาพบว่าแบบจำลอง Bidirectional LSTM มีประสิทธิภาพสูงที่สุดในกลุ่มลูกค้าทั่วไป กลุ่มเสี่ยงต่ำ เสี่ยงกลาง และเสี่ยงสูง ซึ่งสอดคล้องกับการค้นพบในงานวิจัยของ Ala'raj et al. (2021) ที่ระบุว่าแบบจำลอง Bidirectional LSTM สามารถเรียนรู้ลักษณะข้อมูลที่มีความซับซ้อนและต่อเนื่องตามเวลาได้ดีกว่าแบบจำลองทั่วไป เพราะ Bidirectional LSTM เป็นแบบจำลองการเรียนรู้โดยใช้โครงข่ายประสาทเทียมขนาดใหญ่ (Neural Networks) โดยเฉพาะแบบหลายชั้น (Deep Neural Networks) อีกทั้งสามารถประมวลผลได้สองทิศทาง จึงเหมาะสมกับข้อมูลที่มีความซับซ้อนสูง เช่น ข้อมูลพฤติกรรมผิวนัดชำระหนี้ในกลุ่มลูกค้าที่มีลักษณะการผิวนัดชำระหนี้ที่ไม่ชัดเจน ข้อมูลมีลักษณะที่หลากหลาย และอาจมีตัวแปรพฤติกรรมย้อนหลังบางตัวที่อาจมีความสัมพันธ์กันเองอีกด้วย

ในขณะที่แบบจำลอง LSTM มีประสิทธิภาพสูงสุดในกลุ่มลูกค้าความเสี่ยงวิกฤต ซึ่งแตกต่างจาก 4 กลุ่มลูกค้าก่อนหน้านี้ที่แบบจำลอง Bidirectional LSTM มีประสิทธิภาพสูงที่สุด ทั้งนี้เนื่องมาจากโครงสร้างของ LSTM ที่ประมวลผลข้อมูลในทิศทางเดียว (Forward Direction) จึงมีความซับซ้อนน้อยกว่า Bidirectional LSTM ทำให้สามารถติดตามแนวโน้มการผิวนัดจากอดีตสู่ปัจจุบันได้อย่างชัดเจนและเฉพาะเจาะจงยิ่งขึ้น เนื่องจากในกลุ่มลูกค้าความเสี่ยงวิกฤต คือกลุ่มที่มีพฤติกรรมผิวนัดชำระหนี้อย่างต่อเนื่องมาแล้ว 3 เดือน ดังนั้นจึงมีแนวโน้มค่อนข้างชัดเจนว่าจะผิวนัดชำระหนี้ในเดือนถัดไป สอดคล้องกับงานวิจัยของ Lai et al. (2017) ที่ระบุว่าแบบจำลอง Bidirectional RNN ที่มีโครงสร้างสอดคล้องกับ Bidirectional LSTM (Graves, 2012) อาจไม่เหมาะสมสำหรับการทำนายข้อมูลล่วงหน้าที่มีแนวโน้มและลำดับที่ชัดเจน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

บทที่ 5

สรุปผลการวิจัยและข้อเสนอแนะ

5.1 สรุปผลการวิจัย

งานวิจัยนี้มีวัตถุประสงค์หลัก 2 ประการ (1) สร้างแบบจำลองที่เหมาะสมสำหรับกลุ่มลูกค้าที่มีพฤติกรรมการผิदनัดชำระหนี้ที่แตกต่างกัน และ (2) วิเคราะห์ปัจจัยที่มีอิทธิพลต่อการผิदनัดชำระหนี้ โดยสามารถสรุปแบบจำลองที่เหมาะสมที่สุดสำหรับกลุ่มลูกค้าได้ดังนี้

1) การสร้างแบบจำลองที่เหมาะสมสำหรับกลุ่มลูกค้าที่มีพฤติกรรมการผิदनัดชำระหนี้ที่แตกต่างกัน

ผลการวิเคราะห์การเปรียบเทียบประสิทธิภาพของแบบจำลองทั้ง 5 เทคนิค ได้แก่ Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM ภายใต้การแบ่งกลุ่มลูกค้าที่มีพฤติกรรมการผิदनัดชำระหนี้ที่แตกต่างกัน 5 กลุ่ม แสดงการเปรียบเทียบค่าประสิทธิภาพโดยรวม (F1-Score) และ ค่าความระลึก (Recall) ดังตารางที่ 5.1 และ 5.2 ตามลำดับ

ตารางที่ 5.1 ผลเปรียบเทียบค่าประสิทธิภาพโดยรวม (F1-Score) ตามกลุ่มลูกค้าของทั้ง 5 เทคนิค

แบบจำลอง	ค่าประสิทธิภาพโดยรวม				
	กลุ่มที่ 1	กลุ่มที่ 2	กลุ่มที่ 3	กลุ่มที่ 4	กลุ่มที่ 5
Logistic Regression	0.7969	0.8013	0.8031	0.7985	0.8352
Decision Tree	0.7959	0.7977	0.8027	0.8237	0.8550
Neural Network	0.8045	0.8008	0.8003	0.8033	0.8639
LSTM	0.8039	0.8037	0.8072	0.8187	0.8810
Bidirectional LSTM	0.8056	0.8049	0.8099	0.8320	0.8785

ตารางที่ 5.2 ผลเปรียบเทียบค่าความระลึก (Recall) ตามกลุ่มลูกค้าของทั้ง 5 เทคนิค

แบบจำลอง	ค่าความระลึก				
	กลุ่มที่ 1	กลุ่มที่ 2	กลุ่มที่ 3	กลุ่มที่ 4	กลุ่มที่ 5
Logistic Regression	0.7972	0.7984	0.7935	0.7724	0.7927
Decision Tree	0.7958	0.7966	0.7933	0.8054	0.8218
Neural Network	0.8053	0.7978	0.7902	0.7788	0.8352
LSTM	0.8052	0.8008	0.7981	0.7994	0.8604
Bidirectional LSTM	0.8068	0.8024	0.8016	0.8165	0.8556

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้าไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

จาก ตารางที่ 5.1 ซึ่งแสดงค่าประสิทธิภาพโดยรวม (F1-Score) พบว่า **Bidirectional LSTM** ให้ค่าที่สูงที่สุดใน 4 กลุ่ม ได้แก่ กลุ่มที่ 1 ร้อยละ 80.56, กลุ่มที่ 2 ร้อยละ 80.49, กลุ่มที่ 3 ร้อยละ 80.99 และกลุ่มที่ 4 ร้อยละ 83.20 ขณะที่ **LSTM** ให้ค่าที่สูงที่สุดในกลุ่มที่ 5 ที่ร้อยละ 88.10 สำหรับตารางที่ 5.2 ซึ่งแสดงค่าความระลึก (Recall) พบว่า **Bidirectional LSTM** ให้ค่าที่สูงที่สุดใน 4 กลุ่ม ได้แก่ กลุ่มที่ 1 ร้อยละ 80.68, กลุ่มที่ 2 ร้อยละ 80.24, กลุ่มที่ 3 ร้อยละ 80.16 และกลุ่มที่ 4 ร้อยละ 81.65 ขณะที่ **LSTM** ให้ค่าที่สูงที่สุดในกลุ่มที่ 5 ที่ร้อยละ 86.04

โดยสรุปค่าประสิทธิภาพโดยรวม (F1-Score) ให้ผลสอดคล้องกับค่าความระลึก (Recall) และจากผลดังกล่าวสรุปได้ว่าแบบจำลอง Bidirectional LSTM มีความเหมาะสมในกลุ่มลูกค้าทั่วไป ความเสี่ยงต่ำ ความเสี่ยงปานกลาง และ ความเสี่ยงสูง และ LSTM เหมาะสมในกลุ่มลูกค้าความเสี่ยงวิกฤต

นอกจากนี้ผู้วิจัยยังได้แสดงตัวอย่างการวิเคราะห์การทดสอบสมมติฐานด้วย McNemar's Test เฉพาะในกลุ่มที่ 1 ซึ่งพบว่าแบบจำลอง Bidirectional LSTM มีความแตกต่างจากแบบจำลองอื่นอย่างมีนัยสำคัญทางสถิติที่ระดับนัยสำคัญ 0.05

2) การวิเคราะห์ปัจจัยที่มีอิทธิพลต่อการผิดนัดชำระหนี้

การวิเคราะห์ความสำคัญของตัวแปร (Feature Importance) โดยใช้เทคนิค Permutation Importance และ SHAP ในกลุ่มลูกค้าทั่วไป (กลุ่มที่ 1) สำหรับแบบจำลองทั้ง 5 แบบ ได้แก่ Logistic Regression, Decision Tree, Neural Network, LSTM และ Bidirectional LSTM พบว่า ปัจจัยที่มีอิทธิพลสูงสุดต่อการผิดนัดชำระหนี้ที่มีลักษณะร่วมกัน คือ พฤติกรรมการชำระหนี้ย้อนหลัง ซึ่งประกอบด้วยจำนวนครั้งที่ชำระบางส่วน จำนวนวันค้างชำระ และความถี่ในการติดต่อกับลูกค้า ปัจจัยเหล่านี้สะท้อนพฤติกรรมทางการเงินที่ไม่สม่ำเสมอ หรือมีสัญญาณเตือนล่วงหน้าถึงภาวะเสี่ยงต่อการผิดนัด ซึ่งสามารถยืนยันได้จากค่า SHAP ที่แสดงให้เห็นถึงผลกระทบของแต่ละตัวแปรต่อผลการทำนายของแบบจำลองอย่างชัดเจน โดยมีรายละเอียดในแต่ละแบบจำลองดังนี้

Logistic Regression ปัจจัยสำคัญที่สุดคือ "จำนวนครั้งที่จ่ายบางส่วนใน 3 เดือนล่าสุด" ซึ่งมีผลกระทบชัดเจนต่อความเสี่ยง โดยจากค่า SHAP พบว่าค่าตัวแปรที่สูงจะเพิ่มโอกาสผิดนัดชำระอย่างชัดเจน

Decision Tree ปัจจัยที่มีความสำคัญสูงสุดคือ "จำนวนครั้งที่จ่ายบางส่วนใน 9 เดือนล่าสุด" และ "จำนวนวันค้างชำระ" โดยเฉพาะอย่างยิ่ง "จำนวนวันค้างชำระ" ซึ่งแสดงผลกระทบจากค่า SHAP ว่าค่าที่สูงจะเพิ่มความเสี่ยงในการผิดนัดชำระหนี้อย่างเด่นชัด

Neural Network ตัวแปรที่มีผลกระทบมากที่สุดคือ "จำนวนวันค้างชำระ" รองลงมาคือ "จำนวนครั้งที่จ่ายเต็มใน 3 และ 6 เดือนล่าสุด" จากการวิเคราะห์ค่า SHAP พบว่าการมีจำนวนวันค้างชำระที่สูงจะเพิ่มความเสี่ยงต่อการผิดนัดชำระได้อย่างชัดเจน

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

LSTM ปัจจัยที่มีอิทธิพลสูงคือ "จำนวนงวดที่ทำการชำระมากที่สุด ใน 9 เดือนล่าสุด" และ "จำนวนครั้งที่ติดต่อไม่ได้ใน 9 เดือนล่าสุด" จากค่า SHAP ยืนยันว่าการที่ลูกค้ามีการติดต่อไม่ได้หรือชำระหนี้ไม่สม่ำเสมอ ส่งผลให้ความเสี่ยงผิดนัดเพิ่มสูงขึ้น

Bidirectional LSTM ปัจจัยสำคัญคือ "จำนวนครั้งที่ติดต่อได้ใน 12 เดือนล่าสุด" ซึ่งจากค่า SHAP แสดงให้เห็นอย่างชัดเจนว่าการติดต่อที่สม่ำเสมอลดความเสี่ยงต่อการผิดนัดชำระหนี้ได้อย่างมีนัยสำคัญ ขณะที่การขาดการติดต่อหรือการชำระหนี้ไม่สม่ำเสมอเพิ่มความเสี่ยงขึ้นชัดเจน

โดยสรุปการวิเคราะห์เชิงลึกด้วยค่า SHAP และ Permutation Importance ชี้ชัดว่า พฤติกรรมการชำระหนี้ย้อนหลัง โดยเฉพาะการจ่ายบางส่วน จำนวนวันค้างชำระ และความสามารถในการติดต่อกับลูกค้า เป็นปัจจัยหลักที่ส่งผลต่อการผิดนัดชำระหนี้มากที่สุดในทุกแบบจำลอง การนำข้อมูลเหล่านี้ไปใช้ในเชิงนโยบายจะช่วยให้สถาบันการเงินสามารถดำเนินกลยุทธ์บริหารความเสี่ยงเชิงรุกได้อย่างมีประสิทธิภาพยิ่งขึ้น เช่น การคัดกรองลูกค้ากลุ่มเสี่ยง การแจ้งเตือนล่วงหน้า และการออกแบบมาตรการสนับสนุนเฉพาะรายอย่างเหมาะสม

5.2 ข้อเสนอแนะ

จากผลการวิจัยที่ได้ดำเนินการเปรียบเทียบประสิทธิภาพของแบบจำลองการทำนายการผิดนัดชำระหนี้ในสินเชื่อย่อย ผู้วิจัยขอเสนอข้อเสนอแนะเชิงนโยบายและเชิงปฏิบัติการที่สามารถนำไปประยุกต์ใช้ในการพัฒนาระบบบริหารความเสี่ยงของสถาบันการเงินได้อย่างเป็นรูปธรรม ดังนี้

1) การศึกษาแนวทางการสร้างชุดข้อมูลสำหรับแบบจำลองลำดับเวลา เนื่องจากแบบจำลอง LSTM และ Bidirectional LSTM เป็นเทคนิคที่เหมาะสมกับข้อมูลที่มีลักษณะต่อเนื่องตามเวลา (Sequential Data) จึงแนะนำให้มีการปรับปรุงวิธีการจัดรูปแบบข้อมูลจากลักษณะ Snapshot รายเดือน เพื่อนำมาเปรียบเทียบกับแนวทางการศึกษาในงานวิจัยนี้ ทั้งนี้เพื่อเพิ่มศักยภาพในการเรียนรู้และเพิ่มความแม่นยำในการทำนาย

2) การประยุกต์ใช้แบบจำลองในการบริหารความเสี่ยงเชิงนโยบายและเชิงปฏิบัติ ผลการทดลองที่แสดงให้เห็นถึงความเหมาะสมของแบบจำลองแต่ละประเภทในกลุ่มลูกค้าที่แตกต่างกัน สามารถนำไปสู่การออกแบบแนวทางการบริหารความเสี่ยงที่มีประสิทธิภาพ โดยสามารถประยุกต์ใช้ได้หลากหลายมิติ ดังนี้

2.1) การคัดกรองลูกค้าตามระดับความเสี่ยงก่อนการอนุมัติสินเชื่อ สถาบันการเงินสามารถนำผลการศึกษานี้มาออกแบบขั้นตอนการประเมินความเสี่ยงของลูกค้าในเบื้องต้นได้อย่างชัดเจน โดยกำหนดให้ลูกค้าผ่านการประเมินเบื้องต้นด้วยแบบจำลองที่ไม่ซับซ้อน เช่น Logistic Regression หรือ Decision Tree สำหรับลูกค้ากลุ่มที่เริ่มมีสัญญาณการผิดนัดชำระหนี้ในระดับต่ำหรือปานกลาง ซึ่งช่วยลดภาระค่าใช้จ่ายและเวลาในการวิเคราะห์ข้อมูลเบื้องต้นลงได้อย่างมีประสิทธิภาพ

2.2) การติดตามลูกค้าเชิงรุกด้วยแบบจำลองลำดับเวลา สำหรับกลุ่มลูกค้าที่มีประวัติการผิดนัดต่อเนื่องหรือมีพฤติกรรมเสี่ยงที่ชัดเจนแล้ว สถาบันการเงินควรนำแบบจำลองที่มีศักยภาพสูงในเอกสารนี้เป็นเอกสารที่ส่งมอบไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น เมื่อนำไปใช้จริงบนระบบสารสนเทศทางการเงินไม่ว่ากรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้คัดลอกเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

การจัดการข้อมูลลำดับเวลา ได้แก่ LSTM หรือ Bidirectional LSTM มาใช้ติดตามพฤติกรรมลูกค้าแบบเรียลไทม์หรือรายเดือน เพื่อให้สามารถตรวจจับแนวโน้มและความเปลี่ยนแปลงของลูกค้าแต่ละรายได้เร็วขึ้น ซึ่งจะช่วยให้สามารถดำเนินการแก้ไขหรือลดความเสี่ยงได้ทันทั่วทั้งที่ เช่น การแจ้งเตือนลูกค้า การให้คำปรึกษาทางการเงิน หรือการปรับเงื่อนไขการชำระคืนตามความเหมาะสม

2.3) การพัฒนาระบบรายงานและแดชบอร์ด (Dashboard) สถาบันการเงินสามารถนำผลลัพธ์และการตีความผลจากแบบจำลองที่เหมาะสมที่สุดมาสร้างระบบรายงานหรือ Dashboard เพื่อให้เจ้าหน้าที่สามารถตรวจสอบสถานะลูกค้าแต่ละรายได้สะดวกและรวดเร็ว เช่น Dashboard ที่แสดงความเสี่ยงรายบุคคลที่คำนวณจากแบบจำลอง พร้อมคำอธิบายที่ชัดเจนเกี่ยวกับปัจจัยสำคัญที่นำไปสู่การประเมินความเสี่ยง เพื่อให้เจ้าหน้าที่สามารถใช้ข้อมูลนี้ประกอบการตัดสินใจในการอนุมัติสินเชื่อหรือการติดตามลูกค้าได้อย่างมีประสิทธิภาพมากขึ้น

2.4) การประยุกต์ใช้ข้อมูลจากแบบจำลองเพื่อวางแผนกลยุทธ์สินเชื่อเชิงนโยบาย สถาบันสามารถนำข้อมูลเชิงลึกจากการวิเคราะห์แบบจำลองไปประยุกต์ใช้ในการวางแผนนโยบายการให้สินเชื่อ เช่น การกำหนดเงื่อนไขการอนุมัติที่รัดกุมสำหรับกลุ่มลูกค้าที่มีความเสี่ยงสูง หรือการกำหนดนโยบายการติดตามที่เข้มข้นขึ้นในกลุ่มที่มีพฤติกรรมผิดนัดต่อเนื่อง ซึ่งจะช่วยให้การจัดสรรทรัพยากรและการบริหารจัดการลูกค้าในเชิงรุกมีประสิทธิภาพสูงสุด ลดโอกาสการเกิดหนี้เสียในอนาคต

3) การใช้ผลการวิเคราะห์ความสำคัญของตัวแปร (Feature Importance) เพื่อเสริมสร้างความเข้มแข็งในการบริหารความเสี่ยง แนะนำให้สถาบันการเงินนำผลการวิเคราะห์ความสำคัญของตัวแปร (Feature Importance) ไปประยุกต์ใช้ในการออกแบบเกณฑ์การอนุมัติสินเชื่อหรือกำหนดมาตรการเชิงป้องกัน เพื่อเสริมสร้างความเข้มแข็งในการบริหารความเสี่ยง ช่วยให้สามารถวางแผนบริหารจัดการและรับมือกับปัญหาหนี้เสีย (NPL) ได้อย่างเหมาะสมยิ่งขึ้น

เอกสารอ้างอิง

- ชลลดา ม่วงฉิ่งและคณะ. 2564. “การพัฒนาแบบจำลองการพิจารณาให้คะแนนสินเชื่อโดยใช้เทคนิคเหมืองข้อมูล.” วารสารวิทยาศาสตร์และเทคโนโลยี มหาวิทยาลัยหัวเฉียวเฉลิมพระเกียรติ.7(1): 82-92.
- ธนาคารกรุงไทย. 2566. หนี้เสียคืออะไร แก้ไขยังไง ก่อนเป็นหนี้ท่วมตัว. [Online]. Available: <https://krungthai.com/th/financial-partner/learn-financial/1571>
- ธนาคารแห่งประเทศไทย. 2566. สรุปภาพรวมธนาคารพาณิชย์ไตรมาส 2 ปี 2566. [Online]. Available:<https://www.bot.or.th/content/dam/bot/documents/th/news-and-media/news/2023/news-th-20230822.pdf>
- สงกรานต์ สมบุญ. 2562. การพัฒนาระบบจัดอันดับความเสี่ยงพอร์ตสินเชื่อสหกรณ์การเกษตรธนาคารเพื่อการเกษตรและสหกรณ์การเกษตร. [Online]. Available: <https://so01.tcithaijo.org/index.php/stouagjournal/article/view/246565/166656>
- สุพริศร์ สุวรรณิก. 2566. หนี้ครัวเรือนไทย: เพราะเหตุใดจึงต้องกังวล?. [Online]. Available: https://www.bot.or.th/content/dam/bot/documents/th/research-and-publications/articles-and-publications/articles/pdf/Article_20Mar2023_01.pdf
- สถาบันคุ้มครองเงินฝาก. 2567. งบการเงินธนาคารพาณิชย์ต่างจากงบการเงินทั่วไปอย่างไร (ตอนที่ 1). [Online]. Available: <https://www.dpa.or.th/articles/commercial-bank-financial-statements-1>
- สถาบันวิจัยเศรษฐกิจป๋วย อึ๊งภากรณ์. 2566. ปัญหาหนี้ในระบบในประเทศไทย. [Online]. Available: <https://www.pier.or.th/abridged/2022/08/#ปัญหาหนี้ในระบบในประเทศไทย>
- อโนทัย พุทธารีย์และคณะ. 2561. Credit Scoring Model: เครื่องมือในการประเมินคุณภาพสินเชื่อ. [Online]. Available: <https://content.botlc.or.th/mminfo/BOTCollection/BOTFAQ/FAQ132.pdf>
- อุไรพรรณ เจริญรัล และ ภาสกร ตาปสนันท์. 2566. นาโนไฟแนนซ์ คืออะไร ใครรู้บ้าง?. [Online]. Available: https://www.bot.or.th/th/research-and-publications/articles-and-publications/articles/Article_23Apr2015.html

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เอกสารอ้างอิง (ต่อ)

- Adisa, J, et al. 2022. "Credit Score Prediction using Genetic Algorithm-LSTM Technique. 2022 Conference on Information Communications Technology and Society (ICTAS).1-6. DOI: 10.1109/ICTAS53252.2022.9744714.
- Ala'raj, M., Abbod, M. F., & Majdalawieh, M. (2021). Modelling customers credit card behaviour using bidirectional LSTM neural networks. *Journal of Big Data*, 8, Article 61. <https://doi.org/10.1186/s40537-021-00461-7>
- Bergstra, J., & Bengio, Y. 2012. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(Feb), 281-305.
- Bhagat A. 2018. "Predicting Loan Defaults using Machine Learning Techniques." Master's Thesis of University of California State Northridge.
- Breiman, L., Friedman, J., Stone, C.J., & Olshen, R.A. 1984. *Classification and Regression Trees*. Wadsworth.
- Breiman, L. 2001. Random forests. *Machine Learning*, 45(1), 5-32. Quinlan, J.R. 1986. Induction of decision trees. *Machine Learning*, 1(1), 81-106.
- Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T., Shearer, C., & Wirth, R. 2000. *CRISP-DM 1.0: Step-by-step data mining guide*. SPSS Inc.
- Cooper, A. 2021. *Explaining Machine Learning Models: A Non-Technical Guide to Interpreting SHAP Analyses*.
- Coralogix. 2023. *Permutation Importance (PI): Explain Machine Learning Predictions*. Retrieved June 3, 2025, from https://coralogix.com/ai-blog/permutation-importance-pi-explain-machine-learning-predictions/?utm_source=chatgpt.com
- Coser, A., Mazzer, M. M., & Albu, C. 2019. Predictive models for loan default risk assessment. *Economic Computation and Economic Cybernetics Studies and Research*, 53(2), 149–165.
- Duchessnay, É. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
- Dietterich, T. G. 1998. Approximate Statistical Tests for Comparing Supervised. Classification Learning Algorithms. *Neural Computation*, 10(7), 1895–1923. <https://doi.org/10.1162/089976698300017197>

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เอกสารอ้างอิง (ต่อ)

- Fawcett, T. 2006. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8), 861–874. <https://doi.org/10.1016/j.patrec.2005.10.010>
- Géron, A. 2019. *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow* (2nd ed.). O'Reilly Media.
- Gonçalves, E. B., and Gouvêa, M. A. 2016. “Collection Score and the opportunities for non-performing loans market.” Department of Marketing, Sao Paulo State University.
- Gonçalves, E. B., and Gouvêa, M. A. 2021. “Credit Risk Analysis Applying Logistic Regression, Neural Networks and Genetic Algorithms Models.” *International Journal of Advanced Engineering Research and Science*. 8(9): 198-209.
- Goodfellow, I., Bengio, Y., & Courville, A. 2016. *Deep Learning*. MIT Press.
- Graves, A. 2012. *Supervised sequence labelling with recurrent neural networks*. Springer. <https://doi.org/10.1007/978-3-642-24797-2>
- Graves, A., & Schmidhuber, J. 2005. Framewise phoneme Classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, 18(5-6), 602-610.
- Graves, A., Fernández, S., Gomez, F., & Schmidhuber, J. 2013. Connectionist temporal Classification: labelling unsegmented sequence data with recurrent neural networks. *ICML*.
- Guillaume Chevalier. 2017. *The LSTM Cell*. Wikimedia Commons. Retrieved from https://commons.wikimedia.org/wiki/File:The_LSTM_cell.png
- Guyon, I., & Elisseeff, A. 2003. An Introduction to Variable and Feature Selection. *Journal of Machine Learning Research*, 3, 1157–1182.
- Han, J., Kamber, M., & Pei, J. 2011. *Data Mining: Concepts and Techniques*. Morgan Kaufmann.
- Hastie, T., Tibshirani, R., & Friedman, J. 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer.
- Hochreiter, S., & Schmidhuber, J. 1997. Long short-term memory. *Neural Computation*, 9(8), 1735-1780.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เอกสารอ้างอิง (ต่อ)

- Hosmer, D.W., Lemeshow, S., & Sturdivant, R.X. 2013. Applied Logistic Regression (3rd ed.). Wiley.
- Kohavi, R. 1995. A study of Cross-Validation and bootstrap for accuracy estimation and model selection. In Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI) (pp. 1137-1143).
- Kuhn, M., & Johnson, K. 2013. Applied Predictive Modeling. Springer.
- Lemaître, G., Nogueira, F., & Aridas, C. K. 2017. Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. *Journal of Machine Learning Research*, 18, 1-5.
- Lai, G., Chang, W. C., Yang, Y., & Liu, H. (2017). Modeling Long- and Short-Term Temporal Patterns with Deep Neural Networks. arXiv preprint arXiv:1703.07015. Retrieved from <https://arxiv.org/abs/1703.07015>
- LeCun, Y., Bengio, Y., & Hinton, G. 2015. Deep learning. *Nature*, 521(7553), 436-444.
- Lessmann, S., Baesens, B., Seow, H.V., & Thomas, L.C. 2015. Benchmarking state-of-the-art Classification algorithms for credit scoring: An update of research. *European Journal of Operational Research*, 247(1), 124-136.
- Ilyas Idrisovich Ismagilov et al. 2020. "Collection Scoring Models Development and Research Based on the Deductor Analytical Platform." *Nexo Revista Científica*. 33(02): 608-615.
- McNemar, Q. (1947). Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika*, 12(2), 153-157. <https://doi.org/10.1007/BF02295996>
- Molnar, C. 2022. Interpretable Machine Learning. <https://christophm.github.io/interpretable-ml-book/>
- Ng, A.Y. 2004. Feature selection, L1 vs. L2 Regularization, and rotational inVariance. Proceedings of the twenty-first international conference on Machine learning.
- Noriega, J. P., Rivera, L. A., & Herrera, J. A. (2023). Machine learning for credit risk prediction: A systematic literature review. *Data*, 8(11), 169. <https://doi.org/10.3390/data8110169>

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

เอกสารอ้างอิง (ต่อ)

- Powers, D. M. W. (2011). Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness & Correlation. *Journal of Machine Learning Technologies*, 2(1), 37–63.
- Refaeilzadeh, P., Tang, L., & Liu, H. 2009. Cross-Validation. In *Encyclopedia of Database Systems* (pp. 532-538). Springer.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. 1986. Learning representations by back-propagating errors. *Nature*, 323(6088), 533-536.
- Sabato, G. 2010. Credit Risk Scoring Models. [Online]. Available: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=1546347
- Safavian, S.R., & Landgrebe, D. 1991. A survey of decision tree classifier methodology. *IEEE Transactions on Systems, Man, and Cybernetics*, 21(3), 660-674.
- Saito, T., & Rehmsmeier, M. (2015). The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLOS ONE*, 10(3), e0118432. <https://doi.org/10.1371/journal.pone.0118432>
- Schuster, M., & Paliwal, K.K. 1997. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11), 2673-2681.
- Shearer, C. 2000. The CRISP-DM model: The new blueprint for data mining *Journal of Data Warehousing*, 5(4), 13-22.
- Tavakoli, N. (2020). Locality Sensitive Hashing-based Sequence Alignment Using Deep Bidirectional LSTM Models. *arXiv preprint arXiv:2004.02094*
- Wirth, R., & Hipp, J. 2000. CRISP-DM: Towards a standard process model for data mining. In *Proceedings of the 4th International Conference on the Practical Applications of Knowledge Discovery and Data Mining* (pp. 29-39).
- Wiso. (n.d.). *File:Neural network example.svg*. Wikimedia Commons. Retrieved from https://commons.wikimedia.org/wiki/File:Neural_network_example.svg
- Zhou, L., & Li, Y. 2019. Credit risk evaluation using machine learning: An application to Chinese consumer credit. *Emerging Markets Finance and Trade*, 55(10), 2255-2273.
- Zhou, Y. 2022. Loan Default Prediction Based on Machine Learning Methods. School of Statistics, Beijing Normal University.

เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ตัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้

ประวัติผู้เขียน

ชื่อ นางสาวเปรมสุตา ปัดสาโย
 วัน เดือน ปีเกิด 21 ตุลาคม 2531
 ที่อยู่ปัจจุบัน 8/52 The connect 37 ซ.ช่างอากาศอุทิศ แยก1-2 แขวงดอนเมือง เขต
 ดอนเมือง กรุงเทพมหานคร 10200
 ประวัติการศึกษา 2555 วิทยาศาสตรบัณฑิต สาขาวิทยาการคอมพิวเตอร์
 มหาวิทยาลัยมหาสารคาม



เอกสารนี้เป็นเอกสารที่สงวนไว้สำหรับการใช้งานเพื่อการศึกษาเท่านั้น ไม่อนุญาตให้นำไปใช้ประโยชน์ด้านการค้า
 ไม่ว่าจะกรณีใดๆทั้งสิ้น อีกทั้งห้ามมิให้ดัดแปลงเนื้อหา และต้องอ้างอิงถึงเจ้าของเอกสารทุกครั้งที่มีการนำไปใช้