



การวิเคราะห์การถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของ
การเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์

A Multi-Output Regression Approach for Predicting Frequency and
Severity of Auto Insurance Claims

นางสาวชฎารัตน์	อัมสารพวงค์	รหัสประจำตัว	643021198-6
นายปารเมศ	ศิริพรรณนท	รหัสประจำตัว	643020446-8

รายงานนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต

หลักสูตรสถิติและวิทยาการข้อมูล วิชาเอกสถิติศาสตร์ คณะวิทยาศาสตร์

มหาวิทยาลัยขอนแก่น

ปีการศึกษา 2567

การวิเคราะห์การถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของ
การเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์

A Multi-Output Regression Approach for Prediction Frequency and
Severity of Auto Insurance Claims

นางสาวชฎารัตน์	อิมสารพวงค์	รหัสประจำตัว	643021198-6
นายปารเมศ	ศิริพรรณนธ์	รหัสประจำตัว	643020446-8

รายงานนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาวิทยาศาสตรบัณฑิต

หลักสูตรสถิติและวิทยาการข้อมูล วิชาเอกสถิติศาสตร์ คณะวิทยาศาสตร์

มหาวิทยาลัยขอนแก่น

ปีการศึกษา 2567

หัวข้อโครงงานวิจัย	การวิเคราะห์การถดถอยพหุผลสัมฤทธิ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์	
นักศึกษา	นางสาวชฎารัตน์ อิ่มสารพวงค์	รหัสประจำตัว 643021198-6
	นายปารเมศ ศิริพรรณนนท์	รหัสประจำตัว 643020446-8
อาจารย์ที่ปรึกษา	ผศ. ดร.ธิปไตย พงษ์ศาสตร์	
	ดร.พิชญ์ วรัชโชติเสถียร	

สาขาวิชาสถิติ คณะวิทยาศาสตร์ มหาวิทยาลัยขอนแก่น อนุมัติให้รายงานฉบับนี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรวิทยาศาสตรบัณฑิต (สถิติและวิทยาการข้อมูล)

.....อาจารย์ที่ปรึกษา

(ผศ. ดร.ธิปไตย พงษ์ศาสตร์)

วันที่ เดือน..... พ.ศ. 2568

.....อาจารย์ที่ปรึกษา

(ดร.พิชญ์ วรัชโชติเสถียร)

วันที่ เดือน..... พ.ศ. 2568

.....หัวหน้าสาขาวิชาสถิติ

(ผศ. ดร.สุกัญญา เรืองสุวรรณ)

วันที่ เดือน..... พ.ศ. 2568

หัวข้อโครงงานวิจัย	การวิเคราะห์การถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์		
นักศึกษา	นางสาวชฎารัตน์ อิ่มสารพวงค์	รหัสประจำตัว	643021198-6
	นายปารเมศ ศิริพรรณนนท์	รหัสประจำตัว	643020446-8
อาจารย์ที่ปรึกษา	ผศ. ดร.ธิปไตย พงษ์ศาสตร์ ดร.พิชญา วิรัชโชติเสถียร		

บทคัดย่อ

การคาดการณ์ความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนจากผู้เอาประกันภัย มีความสำคัญต่อการบริหารจัดการความเสี่ยงและการกำหนดอัตราเบี้ยประกันภัยของบริษัทประกันภัยรถยนต์ หากการคาดการณ์ไม่มีความแม่นยำอาจส่งผลกระทบต่อเสถียรภาพทางการเงินของบริษัท งานวิจัยนี้มีวัตถุประสงค์เพื่อสร้างและเปรียบเทียบประสิทธิภาพของแบบจำลองการถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนด้วยการเรียนรู้ของเครื่อง และเพื่อพัฒนาเว็บแอปพลิเคชันที่ช่วยในการทำนายความถี่และความรุนแรงในการเรียกร้องค่าสินไหมทดแทน โดยข้อมูลที่ใช้ในงานวิจัยเป็นข้อมูลกรมธรรม์ประกันภัยรถยนต์จากบริษัทประกันภัยแห่งหนึ่งในประเทศสเปน ตั้งแต่ปี พ.ศ. 2549 – 2558 จำนวน 80,924 ราย แบบจำลองการถดถอยพหุผลลัพธ์ที่ศึกษา ได้แก่ Random Forest, Extreme Gradient Boosting, Light Gradient Boosting Machine และ Artificial Neural Network ทำการวัดประสิทธิภาพของแบบจำลองด้วยค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (MAE) ค่ารากที่สองของความคลาดเคลื่อนกำลังสองเฉลี่ย (RMSE) และค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (SMAPE) ข้อมูลถูกแบ่งเป็นชุดฝึกและชุดทดสอบในอัตราส่วน 80:20 และฝึกแบบจำลองด้วย 10-fold Cross Validation ($K = 10$) ทำการปรับแต่งพารามิเตอร์ที่เหมาะสมสำหรับแต่ละแบบจำลองด้วยวิธี GridSearch และใช้ Shapley Additive Explanations (SHAP) เพื่อวิเคราะห์ปัจจัยที่มีอิทธิพลต่อการทำนาย นอกจากนี้ยังมีการเปรียบเทียบระหว่าง Multi-Output Regression ที่มีผลประสิทธิภาพดีที่สุด กับ Single-Output Regression

ผลการศึกษาพบว่า แบบจำลอง Multi-Output Random Forest Regression มีประสิทธิภาพสูงสุดในการทำนายความถี่ของการเรียกร้องค่าสินไหมทดแทน ($MAE = 0.2427$, $RMSE = 0.7832$, $SMAPE = 28.8311$) และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน ($MAE = 228.9800$, $RMSE = 864.6607$, $SMAPE = 29.2370$) การวิเคราะห์ Shapley

Additive Explanations (SHAP) แสดงให้เห็นว่าปัจจัยสำคัญที่มีผลต่อการทำนายมากที่สุดสามอันดับแรก คือ จำนวนกรรมธรรมที่ผู้ถือกรรมธรรมมี อายุของใบอนุญาตขับขี่ และอายุของผู้ถือกรรมธรรม ผลการเปรียบเทียบระหว่าง Multi-Output Regression และ Single-Output Regression พบว่าประสิทธิภาพโดยรวมใกล้เคียงกัน แต่ Multi-Output Regression มีข้อได้เปรียบด้านความสะดวก เนื่องจากสามารถทำนายหลายผลลัพธ์พร้อมกัน ช่วยลดขั้นตอนในการฝึกแบบจำลอง และเพื่อรองรับการใช้งานจริงผู้วิจัยยังได้พัฒนาเว็บแอปพลิเคชันที่เชื่อมต่อกับ multi-output random forest regression ผ่าน Flask ช่วยให้บริษัทประกันภัยสามารถป้อนข้อมูลกรรมธรรม และทราบผลการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนได้แบบเรียลไทม์

คำสำคัญ: การถดถอยพหุผลลัพธ์ จำนวนการเรียกร้องค่าสินไหมทดแทน ความรุนแรงของการเรียกร้องค่าสินไหมทดแทน การเรียนรู้ของเครื่อง การประกันภัยรถยนต์

สาขาวิชาสถิติ

ลายมือชื่อนักศึกษา.....

ปีการศึกษา 2567

(นางสาวชฎารัตน์ อิ่มสารพวงค์)

ลายมือชื่อนักศึกษา.....

(นายปารเมศ ศิริพรรณนนท์)

ลายมือชื่อนักศึกษา.....

(ผศ. ดร.ธิปไตย พงษ์ศาสตร์)

ลายมือชื่อนักศึกษา.....

(ดร.พิชญ์ วิรัชโชติเสถียร)

Title	A Multi-Output Regression Approach for Prediction Frequency and Severity of Auto Insurance Claims	
Student	Miss. Chadarat Imsarapang	Student ID 643021198-6
	Mister. Parames Siripathanon	Student ID 643020446-8
Project Advisor	Asst. Prof. Dr. Tippatai Pongsart	
	Dr. Pitchaya Wiratchotisation	

ABSTRACT

Predicting insurance claim frequency and severity is crucial for risk management and premium rate determination in automobile insurance. Inaccurate predictions may adversely affect the financial stability of insurance companies. The aims of the study is to develop and evaluate the performance of multi-output regression models for predicting the frequency and severity of insurance claims using machine learning techniques, to develop a web application to support real-time prediction of claim frequency and severity. The dataset used in this research consists of automobile insurance policy records from an insurance company in Spain, covering the period from 2006 to 2015, with a total of 80,924 records. The multi-output regression models examined include Random Forest, Extreme Gradient Boosting, Light Gradient Boosting Machine, and Artificial Neural Network. Model performance is assessed using Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Symmetric Mean Absolute Percentage Error (SMAPE). The dataset is divided into training and testing sets in an 80:20 ratio, with models trained using 10-fold cross-validation ($K = 10$). Additionally, the optimal parameters for each model were fine-tuned using GridSearch. The Shapley Additive Explanation (SHAP) analysis is utilized in order to identify the most influential factors in claim prediction. Additionally, the performance of the best Multi-Output Regression is compared with a Single-Output Regression.

The results indicated that the Multi-Output Random Forest Regression model achieves the highest predictive performance for both claim frequency (MAE = 0.2427, RMSE = 0.7832, SMAPE = 28.8311) and claim severity (MAE = 228.9800, RMSE = 864.6607, SMAPE = 29.2370). The Shapley Additive Explanations (SHAP) analysis identified the three most influential factors in claim prediction: the number of policies held by the policyholder, the age of the driver's license, and the age of the policyholder. A comparison between Multi-Output Regression and Single-Output Regression demonstrated that while their overall performance was similar, Multi-Output Regression offers a practical advantage by predicting multiple outputs simultaneously, helps reduce the steps in model training. To support real-world applications, a Flask-based web application integrating the Multi-Output Random Forest Regression model was developed, enabling insurers to input policyholder data and obtain real-time claim predictions.

Keywords: Multi-Output Regression, Claims Frequency, Claims Severity, Machine Learning, Auto Insurance

Department of Statistics
Academic year 2024

Signature of student.....
(Miss. Chadarat Imsarapang)

Signature of student.....
(Mister. Parames Siripathanon)

Signature of project advisor.....
(Asst. Prof. Dr. Tippatai Pongsart)

Signature of project advisor.....
(Dr. Pitchaya Wiratchotisation)

กิตติกรรมประกาศ

การวิจัยเรื่อง “การวิเคราะห์การถดถอยพหุผลสัมฤทธิ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์” ฉบับนี้ ดำเนินการวิจัยจนสำเร็จจุล่งตามวัตถุประสงค์ของการศึกษา ทั้งนี้ผู้วิจัยขอขอบคุณอาจารย์ที่ปรึกษาการวิจัย ผู้ช่วยศาสตราจารย์ ดร.ธิปไตย พงษ์ศาสตร์ และ ดร.พิชญา วิรัชโชติเสถียร ที่ได้ให้คำปรึกษา แนะนำอย่างละเอียดรอบคอบ และช่วยแก้ไขปัญหในทุกขั้นตอนของการศึกษาด้วยความใส่ใจ ทำให้ผู้วิจัยสามารถดำเนินงานวิจัยได้อย่างมีประสิทธิภาพและบรรลุเป้าหมายที่วางไว้ ขอขอบคุณคณะกรรมการรองศาสตราจารย์ ดร.วิชุดา ไชยสีวามงคล และ ดร.ธนพงศ์ อินทระ ที่ได้ชี้แนะแนวทางในการปรับปรุงแก้ไข เพื่อให้งานวิจัยมีความสมบูรณ์และมีคุณค่ามากยิ่งขึ้น

ขอขอบคุณครอบครัวและเพื่อนอันเป็นที่รัก ที่คอยอยู่เคียงข้าง สนับสนุน และเป็นกำลังใจที่ช่วยให้ผู้วิจัยสามารถทำงานวิจัย ไม่ว่าจะเป็นคำพูดให้กำลังใจ การรับฟังเมื่อยามเหนื่อยล้า หรือการอยู่เคียงข้างในช่วงเวลาที่หนักหน่วง และขอขอบคุณตัวเองที่มุ่งมั่น ไม่ย่อท้อ และพยายามอย่างเต็มที่ในทุกขั้นตอนของการศึกษา

สุดท้ายนี้ ผู้วิจัยขอแสดงความซาบซึ้งใจและขอบคุณทุกท่าน ทั้งที่ได้กล่าวนามและไม่ได้เอ่ยนาม ที่มีส่วนร่วมในการสนับสนุน ให้คำแนะนำ และมอบความเมตตาในรูปแบบต่าง ๆ ตลอดระยะเวลาการดำเนินงานวิจัย ผู้วิจัยหวังเป็นอย่างยิ่งว่างานวิจัยฉบับนี้จะเป็นประโยชน์ต่อผู้ที่ศึกษาเกี่ยวกับการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ หากมีข้อผิดพลาดประการใด คณะผู้วิจัยขออภัยไว้ ณ ที่นี้

คณะผู้วิจัย

มีนาคม 2568

สารบัญ (ต่อ)

	หน้า
2.1.5.3 การถดถอยพหุผลลัพธ์แบบโลทเกรเดียนส์บูสติง (Multi-Output Light Gradient Boosting Regression)	22
2.1.5.4 การถดถอยพหุผลลัพธ์โครงข่ายประสาทเทียม (Multi-Output Artificial Neural Network)	23
2.1.6 การฝึกแบบจำลอง (Model train)	25
2.1.7 การปรับแต่งพารามิเตอร์ (Hyperparameter Tuning)	26
2.1.8 การประเมินและเปรียบเทียบประสิทธิภาพของแบบจำลอง (Model evaluation and comparison)	27
2.1.9 Shapley Additive Explanation (SHAP)	30
2.2 งานวิจัยที่เกี่ยวข้อง	31
2.3 กรอบแนวคิดการวิจัย	34
บทที่ 3 วิธีการดำเนินงานวิจัย	35
3.1 การทำความเข้าใจธุรกิจ (Business Understanding)	36
3.2 การศึกษาและทำความเข้าใจข้อมูล (Data Understanding)	36
3.3 การเตรียมข้อมูล (Data Preparation)	37
3.3.1 การทำความสะอาดข้อมูล (Data Cleansing)	37
3.3.2 การแปลงข้อมูล (Data Transformation)	38
3.4 การสร้างแบบจำลองการถดถอยพหุผลลัพธ์โดยใช้การเรียนรู้ของเครื่อง (Modeling)	39
3.4.1 การแบ่งข้อมูล (Data partitioning)	39
3.4.2 การสร้างแบบจำลองการเรียนรู้ของเครื่อง (Model)	40
3.5 การประเมินและเปรียบเทียบประสิทธิภาพแบบจำลอง (Model Evaluation and Comparison)	45
3.6 การนำแบบจำลองไปใช้งานจริง (Deployment)	46
บทที่ 4 ผลการวิจัย	48
4.1 ผลการศึกษาธุรกิจประกันภัยรถยนต์	49
4.2 ผลการศึกษาข้อมูลการเรียกร้องค่าสินไหมทดแทน	50
4.2.1 รายละเอียดของข้อมูล	50

สารบัญ (ต่อ)

	หน้า
4.2.2 ผลการวิเคราะห์ข้อมูลเชิงสำรวจ	52
4.3 ผลการเตรียมข้อมูลก่อนการวิเคราะห์	59
4.3.1 ผลการทำความสะอาดข้อมูล	59
4.3.2 ผลการแปลงข้อมูล	61
4.4 ผลการสร้างแบบจำลอง	62
4.4.1 ผลการแบ่งข้อมูล	62
4.4.2 ผลการสร้างแบบจำลองฐาน	63
4.4.3 ผลการปรับแต่งพารามิเตอร์ของแต่ละแบบจำลอง	63
4.5 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลอง	65
4.5.1 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการ ถดถอยพหุผลลัพ์	65
4.5.2 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการ ถดถอยพหุผลลัพ์กับการถดถอยผลลัพ์เดียว	67
4.5.3 ผลการอธิบายแบบจำลองด้วยเทคนิค SHAP	69
4.6 ผลการพัฒนาเว็บแอปพลิเคชัน	71
บทที่ 5 สรุปผลการวิจัย	74
5.1 สรุปผลการวิจัย	75
5.2 อภิปรายผลการวิจัย	76
5.3 ประโยชน์ของสถิติ/สารสนเทศสถิติที่ใช้ในการวิจัย	77
5.3.1 การแปลงข้อมูลก่อนการวิเคราะห์	77
5.3.2 การเรียนรู้ของเครื่อง	78
5.3.3 การพัฒนาเว็บแอปพลิเคชัน	78
5.4 ข้อเสนอแนะ	78
5.4.1 ข้อเสนอแนะในการนำผลการวิจัยไปใช้	78
5.4.2 ข้อเสนอแนะในการทำวิจัยครั้งต่อไป	79
เอกสารอ้างอิง	81
ภาคผนวก	85
ประวัติย่อผู้วิจัย	108

สารบัญตาราง

	หน้า
ตารางที่ 1.1 ข้อมูลประกันภัยรถยนต์ของบริษัทประกันภัยแห่งหนึ่งในประเทศสเปน	3
ตารางที่ 2.1 แบบจำลองและวิธีการที่ใช้ในการสร้างแบบจำลองจากการทบทวนวรรณกรรม	33
ตารางที่ 2.2 ตัวแปรที่ส่งผลต่อการทำนายการเรียกร้องค่าสินไหมทดแทนจากการทบทวนวรรณกรรม	34
ตารางที่ 3.1 library ที่ใช้สำหรับวิเคราะห์ข้อมูล	41
ตารางที่ 3.2 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-RFR	42
ตารางที่ 3.3 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-XGBoost	43
ตารางที่ 3.4 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-LightGBM	43
ตารางที่ 3.5 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-ANN	43
ตารางที่ 3.6 พารามิเตอร์สำหรับ GridSearchCV	44
ตารางที่ 4.1 รายละเอียดข้อมูลการประกันภัยรถยนต์	51
ตารางที่ 4.2 ตัวอย่างข้อมูลที่มีการบันทึกผิดพลาด	59
ตารางที่ 4.3 ตัวอย่างข้อมูลการประกันภัยรถยนต์ที่ทำการแปลงเรียบร้อยแล้ว	61
ตารางที่ 4.4 ตัวอย่างการแปลงข้อมูล Has_claim	62
ตารางที่ 4.5 ร้อยละของข้อมูลชุดฝึกจำแนกตามการเรียกร้องค่าสินไหมทดแทน	62
ตารางที่ 4.6 ร้อยละข้อมูลชุดทดสอบจำแนกตามการเรียกร้องค่าสินไหมทดแทน	63
ตารางที่ 4.7 ค่าทำนายพื้นฐานจากค่าเฉลี่ยของตัวแปรตามในชุดฝึกสอน	63
ตารางที่ 4.8 ขอบเขตของไฮเปอร์พารามิเตอร์กำหนดใน GridSearchCV ของแต่ละแบบจำลอง	64
ตารางที่ 4.9 พารามิเตอร์ที่ดีที่สุดจากการค้นหาด้วย GridsearchCV ของแต่ละแบบจำลอง	65
ตารางที่ 4.10 ผลการทดสอบประสิทธิภาพการทำนายของแต่ละแบบจำลอง	66
ตารางที่ 4.11 ค่าจริงและค่าจากการทำนายโดยใช้แบบจำลอง M-RFR จากข้อมูลชุดทดสอบ	67
ตารางที่ 4.12 ค่าพารามิเตอร์จาก GridSearchCV ของ Single-Output Random Forest Regression	68
ตารางที่ 4.13 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการถดถอยพหุผลลัพธ์กับการถดถอยผลลัพธ์เดียว	68

สารบัญรูปภาพ

	หน้า
ภาพที่ 2.1 ประเภทการเรียนรู้ของเครื่อง	8
ภาพที่ 2.2 ลักษณะการทำงานของการเรียนรู้แบบมีผู้สอน	9
ภาพที่ 2.3 ตัวอย่างประเภทแบบจำลองการเรียนรู้แบบมีผู้สอน	10
ภาพที่ 2.4 ตัวอย่างประเภทแบบจำลองการเรียนรู้แบบไม่มีผู้สอน	10
ภาพที่ 2.5 ความสัมพันธ์ระหว่าง artificial intelligence, machine learning และ deep learning	11
ภาพที่ 2.6 การทำนายพหุผลลัพธ์	12
ภาพที่ 2.7 ลักษณะการทำงานของ Decision Tree	13
ภาพที่ 2.8 ตัวอย่างการทำงานของ Decision Tree	15
ภาพที่ 2.9 ลักษณะการทำงานของ Random Forest	16
ภาพที่ 2.10 แผนผังเซลล์ประสาทมนุษย์	24
ภาพที่ 2.11 ตัวอย่างการทำงานแบบจำลอง ANN	24
ภาพที่ 2.12 กระบวนการทำงานของ 4-fold Cross-validation	26
ภาพที่ 2.13 กรอบแนวคิดการวิจัย	34
ภาพที่ 3.1 ขั้นตอนการดำเนินงานวิจัย	36
ภาพที่ 3.2 ลักษณะของ Box plot	38
ภาพที่ 4.1 ผลการวิเคราะห์สถิติเบื้องต้นของข้อมูล	52
ภาพที่ 4.2 การกระจายตัวของข้อมูลความถี่ของการเรียกร้อยค่าสินไหมทดแทน	54
ภาพที่ 4.3 การกระจายตัวของข้อมูลความรุนแรงในการเรียกร้อยค่าสินไหมทดแทน	54
ภาพที่ 4.4 การกระจายตัวของข้อมูลจำนวนกรมธรรม์ประกันภัยรถยนต์	55
ภาพที่ 4.5 การกระจายตัวของข้อมูลอายุผู้ถือกรมธรรม์	55
ภาพที่ 4.6 การกระจายตัวของข้อมูลอายุใบขับขี่	56
ภาพที่ 4.7 การกระจายตัวของข้อมูลจำนวนกรมธรรม์ประกันภัยอื่น	56
ภาพที่ 4.8 การตรวจสอบการกระจายของข้อมูลเชิงคุณภาพ	57
ภาพที่ 4.9 Scatter plot แสดงความสัมพันธ์ระหว่างตัวแปรเชิงปริมาณกับเชิงคุณภาพ	58
ภาพที่ 4.10 ผลการตรวจสอบค่าสูญหาย (Missing Value)	60
ภาพที่ 4.11 ผลการตรวจสอบค่านอกเกณฑ์	61
ภาพที่ 4.12 ประสิทธิภาพการทำนายความถี่ของการเรียกร้อยค่าสินไหมทดแทน	66

สารบัญรูปภาพ (ต่อ)

ภาพที่ 4.13 ประสิทธิภาพการทำนายความความรุนแรงของการเรียกร้องค่าสินไหมทดแทน	66
ภาพที่ 4.14 ผล SHAP ของแบบจำลอง M-RFR ในการทำนายความถี่ของการเรียกร้องค่าสินไหมทดแทน	69
ภาพที่ 4.15 ผล SHAP ของแบบจำลอง M-RFR ในการทำนายความรุนแรงของการเรียกร้องค่าสินไหมทดแทน	70
ภาพที่ 4.16 หน้าต่างเว็บแอปพลิเคชัน	71
ภาพที่ 4.17 หน้าต่างเว็บแอปพลิเคชันสำหรับกรอกข้อมูลผู้เอาประกันภัย (1)	72
ภาพที่ 4.18 หน้าต่างเว็บแอปพลิเคชันหลังจากกรอกข้อมูลผู้เอาประกันภัย (2)	72
ภาพที่ 4.19 หน้าต่างเว็บแอปพลิเคชันที่แสดงผลการทำนาย	73
ภาพที่ 4.20 QR code: Web Application	73
ภาพที่ 4.21 QR code: GitHub	73

บทที่ 1

บทนำ

1.1 ความเป็นมาและความสำคัญของปัญหา

การประกันภัยรถยนต์ช่วยคุ้มครองความสูญเสียหรือความเสียหายที่เกิดจากการใช้รถยนต์แก่ผู้เอาประกันภัย ไม่ว่าจะเป็นชีวิต ร่างกาย ทรัพย์สินของบุคคลภายนอก รวมถึงบุคคลที่โดยสารในรถยนต์ โดยบริษัทประกันภัยจะประเมินความเสี่ยงด้วยปัจจัยต่างๆ เช่น ประวัติการขับขี่ อายุ ประสบการณ์การขับขี่และอื่น ๆ หากเกิดอุบัติเหตุทางรถยนต์และอยู่ในเงื่อนไขความคุ้มครองของกรมธรรม์ บริษัทประกันภัยรถยนต์จะจ่ายค่าสินไหมทดแทนตามเงื่อนไขในกรมธรรม์เมื่อเกิดความเสียหาย ซึ่งอาจรวมถึงค่ารักษาพยาบาล ค่าซ่อมที่ศูนย์บริการ อุบัติเหตุ หรือค่าธรรมเนียมอื่นๆ ทางกฎหมาย (บริษัทรู้ใจ จำกัด, 2566) หมายความว่าผู้เอาประกันภัยจะทำประกันภัยเพื่อถ่ายโอนความเสี่ยงของความสูญเสีย และในทางกลับกันบริษัทประกันภัยจะต้องรับความเสี่ยงนั้นไว้

จากรายงานล่าสุด พบว่าผลประกอบการของธุรกิจประกันวินาศภัยในปี พ.ศ. 2566 มีมูลค่าเบี้ยประกันภัยรวมทั้งสิ้น 210,141 ล้านบาท โดยเบี้ยประกันภัยรับตรงจากการประกันภัยรถยนต์มีมูลค่ารวม 118,419 ล้านบาท คิดเป็น 56.35% และมีมูลค่าการเรียกร้องค่าสินไหมทดแทนสูงถึง 141,060 ล้านบาท ขณะที่รายงานของสมาคมประกันวินาศภัยไทยปี พ.ศ. 2567 ระบุว่าเบี้ยประกันภัยรับตรงจากประกันภัยรถยนต์ในปีดังกล่าวมีมูลค่า 116,909 ล้านบาท คิดเป็น 55.92% ของผลประกอบการธุรกิจประกันวินาศภัยทั้งหมด ซึ่งอยู่ที่ 209,060 ล้านบาท หมายความว่าบริษัทประกันภัยต้องจ่ายค่าสินไหมทดแทนประมาณ 55.9 บาทต่อเบี้ยประกันภัยทุก ๆ 100 บาทที่ได้รับ (สมาคมประกันวินาศภัยไทย, 2567) แสดงให้เห็นว่าประกันภัยรถยนต์มีความต้องการสูงและเกี่ยวข้องกับคนจำนวนมาก นอกจากนี้จากข้อมูลของศูนย์อุบัติเหตุทางถนน พบว่ามีอัตราการเกิดอุบัติเหตุจราจรเพิ่มขึ้นจะทำให้มีการเรียกร้องค่าสินไหมทดแทนที่สูงขึ้น ส่งผลให้บริษัทประกันวินาศภัยต้องประเมินความเสี่ยงและปรับเบี้ยประกันภัยเพิ่มขึ้น (Inn Why, 2567) ดังนั้นบริษัทประกันภัยจำเป็นต้องมีการเตรียมตัวรับมือกับการเรียกร้องค่าสินไหมทดแทนในอนาคต เนื่องจากอัตราความเสียหายที่สูงอาจส่งผลกระทบต่อกิจการของบริษัท

หลังจากวิกฤติโควิด-19 ในปี พ.ศ. 2565 พฤติกรรมการซื้อประกันภัยรถยนต์ในประเทศไทยเปลี่ยนแปลงไปเนื่องจากผู้เอาประกันภัยรู้สึกไม่คุ้มกับราคาเบี้ยประกันภัยที่ต้องจ่าย ทำให้ส่วนมากซื้อประกันรถยนต์แบบตามความต้องการของผู้ใช้ (Personalize insurance) เพื่อตอบโจทย์และเหมาะสมกับแต่ละบุคคลอย่างแท้จริง (Amarin TV, 2565) การนำอัลกอริทึมมาช่วยในการวิเคราะห์ข้อมูลของลูกค้าแต่ละราย ทำให้บริษัทสามารถเสนอแผนประกันภัยที่เหมาะสมกับพฤติกรรม

การขับเคลื่อนหรือความต้องการด้านความคุ้มครอง และยังสร้างกลยุทธ์การตลาดที่มีประสิทธิภาพ (Binariks, 2024) โดยการประยุกต์ใช้เทคนิคการเรียนรู้ของเครื่อง (Machine learning) เพื่อการคาดการณ์ความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน ช่วยให้บริษัทสามารถประเมินความเสี่ยงที่สอดคล้องกับผู้เอาประกันภัยแต่ละราย และใช้ในการพิจารณาเบี้ยประกันภัยที่เหมาะสมให้กับผู้เอาประกันภัย ทำให้บริษัทสามารถบริหารจัดการความเสี่ยงได้ดีขึ้น

ดังนั้น เพื่อคาดการณ์ความถี่และมูลค่าความเสียหายของการเรียกร้องค่าสินไหมทดแทนจากผู้เอาประกันภัย ผู้วิจัยจะใช้การเรียนรู้ของเครื่อง (Machine learning) ซึ่งเป็นองค์ประกอบหนึ่งของปัญญาประดิษฐ์ (Artificial Intelligence) ที่ช่วยให้คอมพิวเตอร์สามารถเรียนรู้และปรับปรุงประสิทธิภาพได้ด้วยตนเอง (Athiwat, 2019) มาใช้ในการทำนายโดยใช้การทำนายแบบพหุผลลัพธ์ (Multi-output Regression) ที่มีความสามารถในการทำนายค่าผลลัพธ์หลายตัวแปรพร้อมกัน และคำนึงถึงความสัมพันธ์ระหว่างตัวแปรตาม

1.2 วัตถุประสงค์ของงานวิจัย

1.2.1 เพื่อสร้างและเปรียบเทียบประสิทธิภาพของแบบจำลองการถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนด้วยการเรียนรู้ของเครื่อง

1.2.2 เพื่อพัฒนาเว็บแอปพลิเคชันที่ช่วยในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน

1.3 สมมติฐานการวิจัย

แบบจำลองการถดถอยพหุผลลัพธ์แบบเอ็กซ์ตรีมเกรเดียนต์บูสติง (Multi-Output Extreme Gradient Boosting Regression) สามารถทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ได้อย่างมีประสิทธิภาพสูงสุด เนื่องจากการรวมต้นไม้หลายต้น (Decision Tree) โดยแต่ละต้นจะลดความคลาดเคลื่อนจากต้นก่อนหน้า และสามารถประมวลผลแบบขนาน (parallelization) ทำให้มีความเร็วในการประมวลผลสูงมาก ส่งผลให้สามารถจัดการกับข้อมูลที่ซับซ้อนได้ดี และยังลดความคลาดเคลื่อนในการทำนายได้อย่างมีประสิทธิภาพ

1.4 ขอบเขตการวิจัย

1.4.1 ขอบเขตด้านข้อมูล

ข้อมูลที่ใช้ในการวิจัยเป็นข้อมูลการประกันรถยนต์ของบริษัทประกันภัยแห่งหนึ่งในประเทศสเปน ตั้งแต่ปี ค.ศ. 2006 ถึง ค.ศ. 2015 เป็นระยะเวลา 10 ปี รวมทั้งสิ้น 80,924 ราย

ประกอบด้วยตัวแปรอิสระ 8 ตัวแปร และตัวแปรตาม 2 ตัวแปร ได้มาจากการวิจัยของ Catalina Bolance and Raluca Vernic ปี ค.ศ. 2019 แสดงรายละเอียดดังตารางที่ 1.1

ตารางที่ 1.1 ข้อมูลประกันภัยรถยนต์ของบริษัทประกันภัยแห่งหนึ่งในประเทศสเปน

ตัวแปรอิสระ	
ตัวแปร	คำอธิบายตัวแปร
1. npol_auto (X_1)	จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี
2. client_sex (X_2)	เพศของผู้ถือกรมธรรม์ (0=ชาย, 1=หญิง)
3. client_age (X_3)	อายุของผู้ถือกรมธรรม์
4. lic_age (X_4)	อายุของใบอนุญาตขับขี่
5. city (X_5)	พื้นที่อยู่อาศัย (0=อื่นๆ, 1=ในเมืองใหญ่)
6. north (X_6)	อาศัยอยู่ภาคเหนือหรือไม่ (0=ไม่ใช่, 1=ใช่)
7. rest (X_7)	อาศัยอยู่พื้นที่อื่นๆของประเทศ (0=ไม่ใช่, 1=ใช่)
8. client_nother (X_8)	จำนวนกรมธรรม์ประเภทอื่นในบริษัท
ตัวแปรตาม	
1. nclaims_md (Y_1)	ความถี่ในการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ที่มีความเสียหายต่อทรัพย์สิน
2. cost_md (Y_2)	ความรุนแรงในการเรียกร้องค่าสินไหมทดแทนทั้งหมดสำหรับประกันภัยรถยนต์ที่มีความเสียหายต่อทรัพย์สิน (หน่วย: ยูโร)

1.4.2 ขอบเขตด้านแบบจำลองและตัวชี้วัดประสิทธิภาพของแบบจำลอง

การศึกษาในงานวิจัยนี้ มีวัตถุประสงค์เพื่อสร้างแบบจำลองและเปรียบเทียบประสิทธิภาพการทำนายพหุผลลัพธ์ด้วยวิธีการเรียนรู้ของเครื่อง (Machine learning) ประกอบด้วยแบบจำลองดังนี้

1. การถดถอยพหุผลลัพธ์แบบการสุ่มป่าไม้ (Multi-Output Random Forest Regression: M-RFF)
2. การถดถอยพหุผลลัพธ์แบบเอ็กซ์ตรีมเกรเดียนต์บูสติง (Multi-Output Extreme Gradient Boosting Regression: M-XGBoost)
3. การถดถอยพหุผลลัพธ์แบบไลต์เกรเดียนต์บูสติง (Multi-Output Light Gradient Boosting Regression: M-LightGBM)

4. การถดถอยพหุผลลัพธ์แบบโครงข่ายประสาทเทียม (Multi-Output Artificial Neural Network: M-ANN)

ทำการวัดประสิทธิภาพของแบบจำลอง จากค่าต่อไปนี้

1. ค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (Mean Absolute Error: MAE)
2. ค่ารากที่สองของความคลาดเคลื่อนกำลังสองเฉลี่ย (Root Mean Square Error: RMSE)
3. ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (Symmetric Mean Absolute Percentage Error: SMAPE)

1.5 นิยามศัพท์เฉพาะ

1.5.1 ศัพท์เฉพาะ

ความถี่ในการเรียกร้องค่าสินไหมทดแทน คือ จำนวนการเรียกร้องค่าสินไหมทดแทนของประกันภัยรถยนต์ด้านทรัพย์สิน (property damage) จากผู้เอาประกันภัยแต่ละรายในช่วงปี ค.ศ. 2006 ถึง ค.ศ. 2015

ความรุนแรงในการเรียกร้องค่าสินไหมทดแทน คือ มูลค่าการเรียกร้องค่าสินไหมทดแทนของประกันภัยรถยนต์ด้านทรัพย์สิน (property damage) จากผู้เอาประกันภัยแต่ละรายในช่วงปี ค.ศ. 2006 ถึง ค.ศ. 2015

การเรียกร้องค่าสินไหมทดแทน คือ ความถี่หรือความรุนแรงในการเรียกร้องค่าสินไหมทดแทน

1.5.2 คำย่อ

RFR	Random Forest Regression
ANN	Artificial Neural Network
XGboost	Extreme Gradient Boosting
LightGBM	Light Gradient Boosting Machine Regression
M-RFR	Multi-Output Random Forest Regression
M-XGBoost	Multi-Output Extreme Gradient Boosting Regression
M-LightGBM	Multi-Output Light Gradient Boosting Machine
M-ANN	Multi-Output Artificial Neural Network
S-RFR	Single-Output Random Forest Regression
RSS	Residual Sum of Squares
MAE	Mean Absolute Error
RMSE	Root Mean Square Error

MAPE	Mean Absolute Percentage Error
SMAPE	Symmetric Mean Absolute Percentage Error
EDA	Exploratory Data Analysis

1.6 ประโยชน์ที่คาดว่าจะได้รับจากการวิจัย

บริษัทประกันภัยรถยนต์สามารถนำแบบจำลองที่ดีที่สุดและคุณลักษณะที่มีผลกระทบต่อการทำนายความถี่และความรุนแรงในการเรียกร้องค่าสินไหมทดแทน มาประยุกต์ใช้กับข้อมูลของบริษัทประกันภัยเพื่อทำนายการเรียกร้องค่าสินไหมทดแทนและการกำหนดเบี้ยประกันภัยที่เหมาะสมให้กับผู้เอาประกันภัยได้อย่างมีประสิทธิภาพ และพัฒนาเว็บแอปพลิเคชันเพื่อให้บริษัทประกันภัยรถยนต์สามารถเข้าถึงและคาดการณ์ความถี่รวมถึงความรุนแรงของการเรียกร้องค่าสินไหมทดแทนได้อย่างสะดวก รวดเร็ว และมีประสิทธิภาพ

บทที่ 2

ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

ในการสร้างแบบจำลองพหุผลลัพธ์เพื่อทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ ได้อาศัยแนวคิด ทฤษฎี และงานวิจัยที่เกี่ยวข้องดังนี้

2.1 แนวคิดและทฤษฎีที่เกี่ยวข้อง

2.1.1 การประกันภัยรถยนต์

2.1.2 การเรียนรู้ของเครื่อง

2.1.3 การเรียนรู้เชิงลึก

2.1.4 การทำนายพหุผลลัพธ์

2.1.5 แบบจำลองที่ใช้ในการวิเคราะห์การถดถอยพหุผลลัพธ์

2.1.5.1 การถดถอยพหุผลลัพธ์แบบการสุ่มป่าไม้

2.1.5.2 การถดถอยพหุผลลัพธ์แบบเอ็กตริมเกรเดียนต์บูตติง

2.1.5.3 การถดถอยพหุผลลัพธ์แบบโลทเกรเดียนต์บูตติง

2.1.5.4 การถดถอยพหุผลลัพธ์โครงข่ายประสาทเทียม

2.1.6 การฝึกแบบจำลอง

2.1.7 การปรับแต่งพารามิเตอร์

2.1.8 การประเมินและเปรียบเทียบประสิทธิภาพของแบบจำลอง

2.1.9 Shapley Additive Explanation (SHAP)

2.2 งานวิจัยที่เกี่ยวข้อง

2.3 กรอบแนวคิดการวิจัย

โดยมีรายละเอียดแต่ละหัวข้อดังต่อไปนี้

2.1 แนวคิดทฤษฎีที่เกี่ยวข้อง

2.1.1 การประกันภัยรถยนต์

การประกันภัยรถยนต์คือคุ้มครองความสูญเสียหรือความเสียหายที่เกิดกับรถยนต์ แก่ผู้เอาประกันภัยไม่ว่าจะเป็นชีวิต ร่างกาย ทรัพย์สินของบุคคลภายนอก รวมถึงบุคคลที่โดยสารในรถยนต์ โดยบริษัทประกันภัย จะรับผิดชอบค่าใช้จ่ายบางส่วนหรือทั้งหมดตามเงื่อนไขในกรมธรรม์ รายงานองค์การอนามัยโลก (WHO) ระบุว่าอุบัติเหตุจราจรทางถนนส่งผลให้มีผู้เสียชีวิตประมาณ 1.19 ล้านคนต่อปี และมีผู้ได้รับบาดเจ็บระหว่าง 20 ถึง 50 ล้านคน โดยกลุ่มเสี่ยงหลัก ได้แก่ คนเดินถนน นักปั่นจักรยาน และผู้ขับขี่จักรยานยนต์ ตัวเลขเหล่านี้สะท้อนให้เห็นถึงความรุนแรงของปัญหาอุบัติเหตุทางถนน ซึ่งทำให้การประกันภัยรถยนต์เป็นปัจจัยสำคัญที่ช่วยลดภาระทางการเงินจากความเสียหาย และเสริมสร้างความมั่นใจแก่ผู้ขับขี่ (WHO, 2024) นอกจากนี้ การทำประกันภัยยังเป็นแนวทางสำคัญในการบริหารความเสี่ยงทางการเงิน โดยช่วยบรรเทาผลกระทบจากอุบัติเหตุหรือเหตุการณ์ไม่คาดคิด เมื่อเกิดความเสียหาย บริษัทประกันภัยจะเข้ามาชดเชยค่าใช้จ่ายตามเงื่อนไขของกรมธรรม์ ส่งผลให้ปัจจุบันผู้บริโภคให้ความสนใจทำประกันภัยมากขึ้น โดยพิจารณาปัจจัยต่าง ๆ เช่น ระดับความคุ้มครอง ค่าเบี้ยประกัน และความสะดวกในการเรียกร้องค่าสินไหมทดแทน (นพัตถ วิระมิตรชัย, 2562)

ปัจจัยหลักที่เกี่ยวข้องกับการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์มักมีความคล้ายคลึงกันในหลายประเทศ แม้จะมีความแตกต่างด้านกฎหมายหรือเงื่อนไขเฉพาะ แต่ปัจจัยพื้นฐาน เช่น ประวัติการเกิดอุบัติเหตุ อายุและเพศของผู้ขับขี่ ประวัติการเรียกร้องค่าสินไหมทดแทน ความถี่ในการขับขี่ พื้นที่ที่ใช้ขับขี่ และสภาพเศรษฐกิจ ล้วนเป็นปัจจัยสำคัญในการพิจารณาความเสี่ยงและการกำหนดเบี้ยประกันภัย

ในประเทศสเปน พบว่าปัจจัยทางเศรษฐกิจและสภาพภูมิอากาศของแต่ละภูมิภาค ส่งผลต่ออัตราอุบัติเหตุและการเรียกร้องค่าสินไหมทดแทน โดยผู้ขับขี่อายุต่ำกว่า 25 ปีมักถูกจัดอยู่ในกลุ่มเสี่ยงสูง ขณะที่ผู้ขับขี่ที่มีประสบการณ์และไม่มีประวัติอุบัติเหตุจะได้รับอัตราเบี้ยประกันที่ต่ำกว่า นอกจากนี้ ประเภทและอายุของรถยนต์ ยังเป็นอีกปัจจัยที่มีอิทธิพลต่อการกำหนดความเสี่ยง (Johnson, 2024)

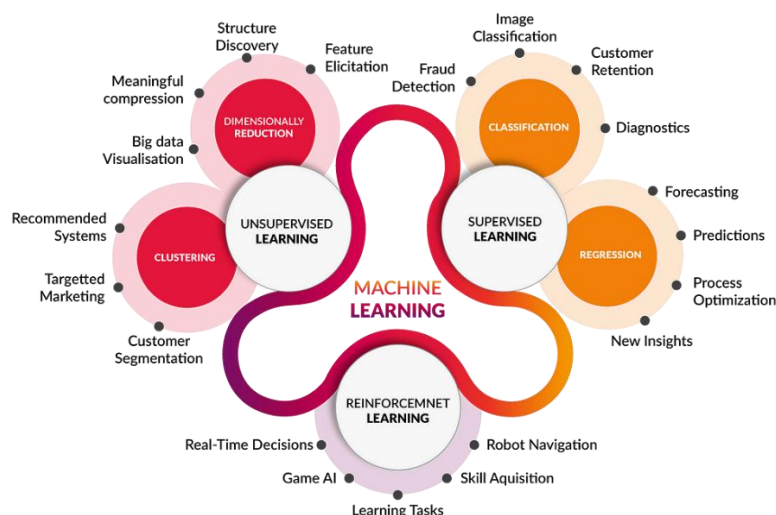
สำหรับประเทศไทย ซึ่งมีอัตราการเกิดอุบัติเหตุสูงที่สุดในเอเชีย ปัจจัยที่ส่งผลต่อการเรียกร้องค่าสินไหมทดแทน ได้แก่ อายุของผู้ถือกรมธรรม์ อายุของรถยนต์ และเพศของผู้ขับขี่ ซึ่งถูกนำมาพิจารณาในการประเมินความเสี่ยงและการกำหนดค่าเบี้ยประกันภัย นอกจากนี้

สภาพถนนและพฤติกรรมการขับขี่ยังเป็นปัจจัยเสริมที่ส่งผลต่อความถี่ในการเกิดอุบัติเหตุ (ปวีศา สุขเรื่อย และสำรวม จงเจริญ, 2561)

เนื่องจากปัจจัยที่ส่งผลต่อการเรียกร้องค่าสินไหมทดแทนมีลักษณะคล้ายคลึงกัน ในหลายประเทศ การวิเคราะห์ข้อมูลจากประเทศหนึ่งสามารถใช้เป็นข้อมูลอ้างอิงหรือแนวทาง ในการพัฒนาระบบการเรียกร้องค่าสินไหมทดแทนในประเทศอื่น ๆ ได้ ซึ่งจะช่วยเพิ่ม ประสิทธิภาพในการประเมินความเสี่ยงและการกำหนดอัตราเบี้ยประกันภัยอย่างเหมาะสม

2.1.2 การเรียนรู้ของเครื่อง (Machine Learning)

การเรียนรู้ของเครื่อง เป็นอัลกอริทึมที่ถูกพัฒนาขึ้นเพื่อจำลองการทำงานของสติปัญญา มนุษย์โดยอาศัยหลักการจากหลากหลายศาสตร์ เช่น ปัญญาประดิษฐ์ สถิติ ความน่าจะเป็น วิทยาการ คอมพิวเตอร์ ทฤษฎีสารสนเทศ จิตวิทยา ทฤษฎีการควบคุม และปรัชญา อัลกอริทึมเหล่านี้สามารถ ทำนายผลลัพธ์ได้โดยไม่ต้องเขียนโปรแกรมคอมพิวเตอร์ที่ซับซ้อน โดยอาศัยการเรียนรู้จากชุด ข้อมูลที่จัดเตรียมไว้ และประสบการณ์จากการทำซ้ำเพื่อทำนายผลลัพธ์ กระบวนการเรียนรู้นี้ เรียกว่าการฝึกอบรม ซึ่งช่วยให้อัลกอริทึมสามารถปรับปรุงประสิทธิภาพได้ด้วยตนเอง ตลอดเวลา และยังสามารถเรียนรู้รูปแบบที่ซับซ้อนและละเอียดอ่อนได้มากกว่ามนุษย์ (Naqa et al., 2022) ประเภทของการเรียนรู้ของเครื่อง (Machine Learning) สามารถแบ่งออกเป็น 3 ประเภท ได้แก่ 1) การเรียนรู้แบบมีผู้สอน (Supervised Learning) 2) การเรียนรู้แบบไม่มีผู้สอน (Unsupervised Learning) และ 3) การเรียนรู้แบบเสริมกำลัง (Reinforcement Learning) แสดงดังภาพที่ 2.1

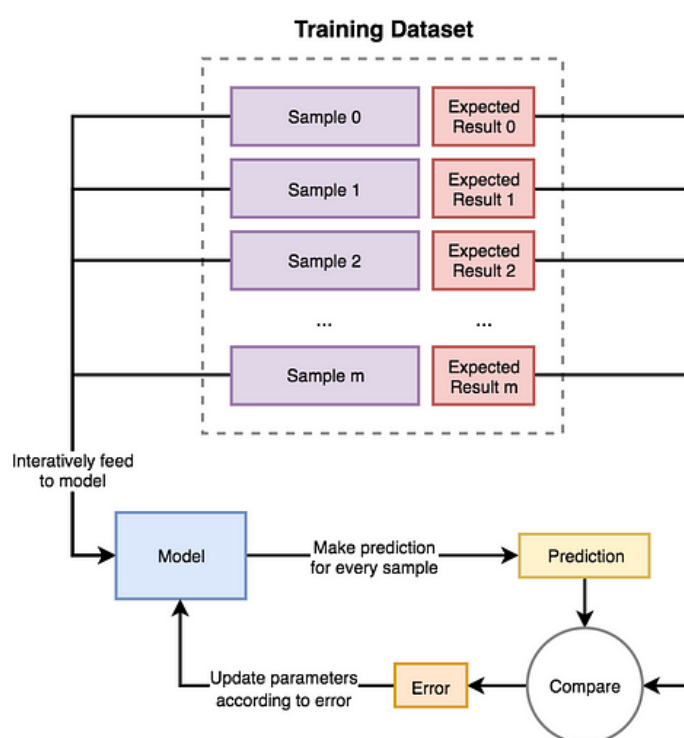


ภาพที่ 2.1 ประเภทการเรียนรู้ของเครื่อง

หมายเหตุ. จาก <https://resources.experfy.com/ai-ml/coding-deep-learning-for-beginners-types-of-machine-learning/>

1) การเรียนรู้แบบมีผู้สอน (Supervised Learning)

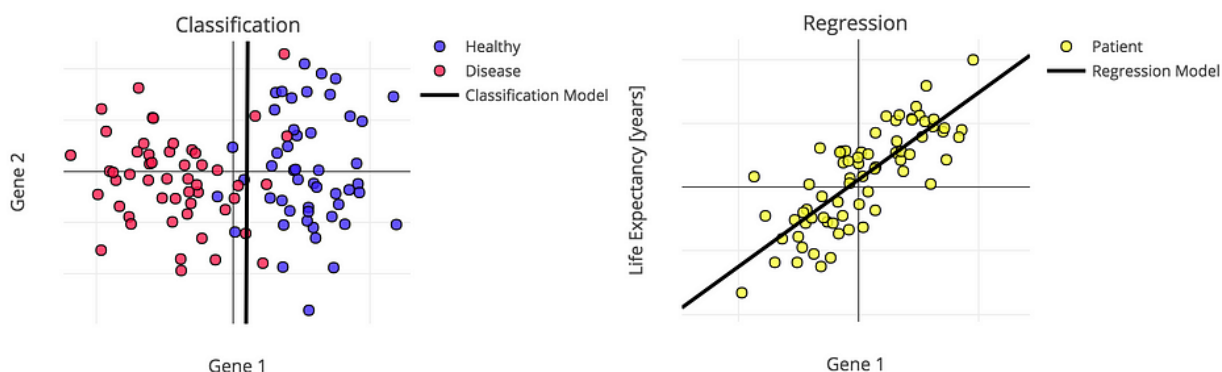
การเรียนรู้แบบมีผู้สอนเป็นเทคนิคที่อัลกอริทึมเรียนรู้จากชุดข้อมูลตัวอย่างที่ประกอบด้วยข้อมูลนำเข้า (input data) และผลลัพธ์ (label) ที่ถูกกำหนดไว้ล่วงหน้า แบบจำลองการเรียนรู้แบบมีผู้สอนจะหาค่าพารามิเตอร์ที่ทำนายผลลัพธ์ ในกระบวนการฝึกอบรมแบบจำลองจะหาค่าสำหรับแต่ละตัวอย่างและเปรียบเทียบกับที่ได้กับผลลัพธ์ที่กำหนดไว้ หากมีความแตกต่างกัน (error) แบบจำลองจะเรียนรู้และปรับปรุงพารามิเตอร์เพื่อลดความผิดพลาดในการทำนายข้อมูลที่ไม่รู้จัก แสดงดังภาพที่ 2.2



ภาพที่ 2.2 ลักษณะการทำงานของการเรียนรู้แบบมีผู้สอน

หมายเหตุ. จาก <https://resources.experfy.com/ai-ml/coding-deep-learning-for-beginners-types-of-machine-learning/>

ประเภทของผลลัพธ์ในเทคนิคนี้แบ่งออกเป็น 2 ประเภท คือ การจำแนกประเภท (Classification) เป็นการกำหนดหมวดหมู่หรือประเภท เช่น การทำนายว่าเป็นโรคเบาหวานหรือไม่ การตรวจจับการฉ้อโกง และอีกประเภทหนึ่งคือ การถดถอย (Regression) ซึ่งผลการทำนายจะเป็นค่าต่อเนื่อง (Krzyk, 2023) แสดงดังภาพที่ 2.3

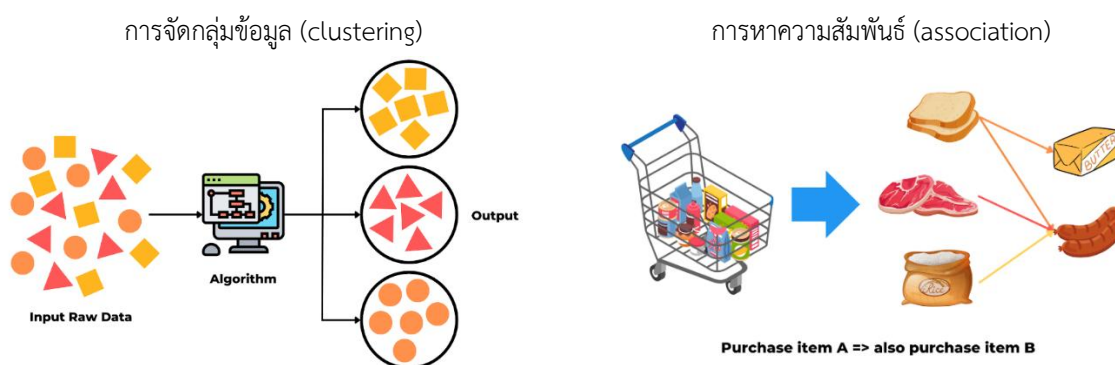


ภาพที่ 2.3 ตัวอย่างประเภทแบบจำลองการเรียนรู้แบบมีผู้สอน

หมายเหตุ. จาก <https://resources.experfy.com/ai-ml/coding-deep-learning-for-beginners-types-of-machine-learning/>

2) การเรียนรู้แบบไม่มีผู้สอน (Unsupervised Learning)

การเรียนรู้แบบไม่มีผู้สอนเป็นเทคนิคที่อัลกอริทึมสามารถเรียนรู้และสำรวจชุดข้อมูลโดยที่ไม่มีการกำหนดผลลัพธ์ (label) ไว้ล่วงหน้า แบบจำลองจะตีความโครงสร้างและความสัมพันธ์ด้วยตนเอง การเรียนรู้แบบไม่มีผู้สอนถูกใช้ในงานต่างๆ เช่น สำหรับการจัดกลุ่มข้อมูล (clustering) การลดมิติข้อมูล (dimensionality reduction) การหาความสัมพันธ์ (association) (Eastgate Software, 2024) แสดงดังภาพที่ 2.4



ภาพที่ 2.4 ตัวอย่างประเภทแบบจำลองการเรียนรู้แบบไม่มีผู้สอน

หมายเหตุ. จาก <https://eastgate-software.com/what-is-unsupervised-learning/>

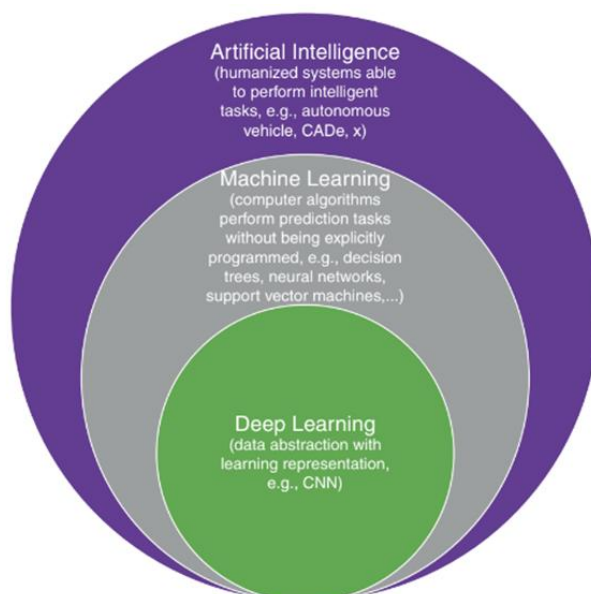
3) การเรียนรู้แบบเสริมกำลัง (Reinforcement Learning)

การเรียนรู้แบบเสริมกำลังเป็นเทคนิคที่อัลกอริทึมเรียนรู้ผ่านการลองผิดลองถูกและการทำซ้ำ โดยจะได้รับรางวัล (Reward) ตามผลลัพธ์ของการกระทำ (Action) ที่ส่งผลต่อสิ่งแวดล้อม

(Environment) ผู้กระทำจะเรียนรู้จากข้อผิดพลาดในอดีตเพื่อปรับปรุงและพัฒนาตนเองให้สามารถได้รับรางวัลที่ดีที่สุด ตัวอย่างการใช้งาน ได้แก่ ระบบขับรถยนต์อัตโนมัติ (Self-driving car) หรือการซื้อขายหุ้นเพื่อให้ได้ผลตอบแทนสูงสุด (Stock Trading Optimization) (BDI, 2020)

2.1.3 การเรียนรู้เชิงลึก (Deep Learning)

การเรียนรู้เชิงลึก เป็นสาขาย่อยของการเรียนรู้ของเครื่อง โดยมีจุดเด่นที่การเรียนรู้จากการแสดงข้อมูล (Representation Learning) จากข้อมูลพื้นฐานไปจนถึงข้อมูลที่ซับซ้อน ช่วยให้ระบบสามารถเรียนรู้จากข้อมูลดิบได้อย่างอัตโนมัติและสามารถจับรูปแบบเชิงลึกภายในข้อมูลได้โดยไม่ต้องมีการกำหนดคุณลักษณะด้วยมืออย่างชัดเจน ด้วยการใช้โครงข่ายประสาทเทียม (Neural Networks) (Naqa, Murphy and Li, 2022) สามารถนำมาใช้งานได้หลากหลาย เช่น การจดจำภาพ (Image Recognition) การประมวลผลภาษาธรรมชาติ (Natural Language Processing) การรู้จำเสียง (Speech Recognition) และการวิเคราะห์ข้อมูลขนาดใหญ่ (Big Data Analytics) เป็นต้น มีแผนภาพแสดงดังภาพที่ 2.5

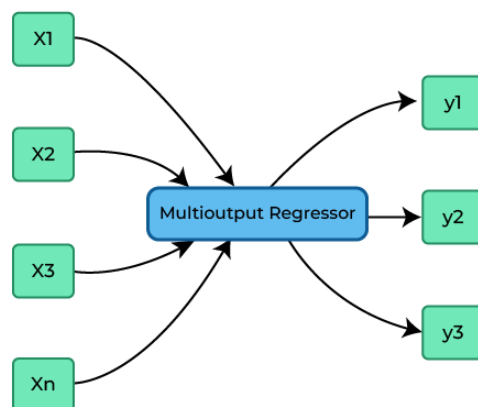


ภาพที่ 2.5 ความสัมพันธ์ระหว่าง artificial intelligence, machine learning และ deep learning
หมายเหตุ. จาก What Are Machine and Deep Learning (p. 3-15), by Naqa et al., 2022

2.1.4 การทำนายพหุผลลัพธ์ (Multi-output regression)

การทำนายพหุผลลัพธ์ (Multi-output regression) หรือที่เรียกว่า multi-target regression หรือ multi-variate regression หรือ multi-response regression ซึ่งเป็นประเภทของการเรียนรู้แบบมีผู้สอน (Supervised Learning) จัดเป็นเทคนิคหนึ่งของการเรียนรู้ของเครื่องช่วย

สร้างแบบจำลองสำหรับทำนายพหุผลลัพธ์โดยใช้ชุดข้อมูลตัวแปรทำนายเดียวกัน เทคนิคนี้สามารถแก้ปัญหาความซับซ้อนในโลกความเป็นจริงได้อย่างมีประสิทธิภาพมากกว่าการทำนายผลลัพธ์ทีละตัว เนื่องจากแบบจำลองจะพิจารณาทั้งความสัมพันธ์ระหว่างคุณลักษณะ (Feature) กับตัวแปรตาม (Output) และความสัมพันธ์ระหว่างตัวแปรตามด้วยกันเอง ทำให้แบบจำลองที่สร้างขึ้นสามารถวิเคราะห์ความซับซ้อนได้ดีและมีประสิทธิภาพสูงในการทำนาย (Borchani et al., 2015) แสดงดังภาพที่ 2.6



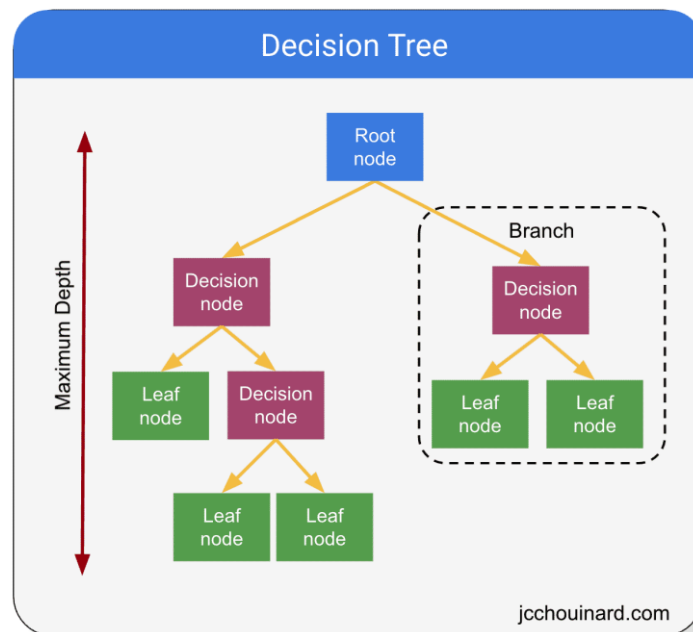
ภาพที่ 2.6 การทำนายพหุผลลัพธ์

หมายเหตุ. จาก <https://www.geeksforgeeks.org/multioutput-regression-in-machine-learning/>

2.1.5 แบบจำลองที่ใช้ในการวิเคราะห์การถดถอยพหุผลลัพธ์ (Model)

แบบจำลองในงานวิจัยฉบับนี้ เป็นการเรียนรู้แบบมีผู้สอน (Supervised Learning) ประเภทการถดถอย (Regression) ที่มีผลลัพธ์ของการทำนายเป็นค่าที่ต่อเนื่อง โดยแบบจำลองที่ใช้ในงานวิจัยนี้เป็นแบบจำลองที่มีโครงสร้างของต้นไม้การตัดสินใจ (Decision Tree) มีรายละเอียดดังนี้

การถดถอยต้นไม้ตัดสินใจ (Decision Tree Regression) เป็นแบบจำลองการเรียนรู้ของเครื่องประเภทมีผู้สอน โดยมีเป้าหมายในการสร้างแบบจำลองเพื่อทำนายค่าผลลัพธ์ด้านการถดถอยหรือค่าที่ต่อเนื่อง การทำงานของ Decision Tree Regression จะใช้โครงสร้างต้นไม้ (tree structure) ในการสร้างแบบจำลองการทำนาย โดยการแบ่งข้อมูลออกเป็นกลุ่มย่อยตามลำดับของเงื่อนไขที่กำหนด โดยมีส่วนประกอบแสดงดังภาพที่ 2.7



ภาพที่ 2.7 ลักษณะการทำงานของ Decision Tree

หมายเหตุ. จาก <https://www.jcchouinard.com/decision-trees-in-machine-learning/>

จากภาพที่ 2.8 สามารถอธิบายส่วนประกอบได้ดังนี้

1. Root Node คือ โหนดเริ่มต้นของการตัดสินใจที่ประกอบด้วยข้อมูลทั้งหมด โหนดนี้จะทำการวิเคราะห์และแบ่งข้อมูลออกเป็นกลุ่มย่อย ซึ่งเมื่อทำการแบ่งแล้วโหนดรากจะสร้างโหนดย่อย (child nodes) เพื่อดำเนินการแบ่งข้อมูลในขั้นตอนถัดไป
2. Child nodes คือ โหนดย่อยที่เกิดจากการแบ่งของโหนดตัดสินใจ โดยแต่ละโหนดย่อยจะมีข้อมูลที่แตกต่างกันตามการแบ่งข้อมูลของโหนดก่อนหน้า
3. Decision node คือ โหนดที่ทำหน้าที่ในการแบ่งข้อมูลออกเป็นกลุ่มย่อยๆ เพื่อเพิ่มความแม่นยำในการทำนายค่าผลลัพธ์ โหนดนี้จะพิจารณาคุณลักษณะของข้อมูลและตัดสินใจว่าจะทำการแบ่งข้อมูลอย่างไรเพื่อให้ได้ผลลัพธ์ที่ดีที่สุด
4. Leaf Node or Terminal Node คือ โหนดสุดท้ายที่ไม่สามารถแบ่งย่อยได้อีก ค่าในโหนดนี้จะเป็นค่าที่ใช้ในการทำนาย (prediction)
5. Branches or Sub-Tree คือ กิ่งของต้นไม้การตัดสินใจที่เชื่อมต่อกับโหนดตัดสินใจและโหนดย่อย โดยจำนวนชั้นของกิ่งที่มากที่สุดจากโหนดรากไปยังโหนดสุดท้ายเรียกว่า "Maximum Depth" ของต้นไม้ (Pathmind, n.d.)

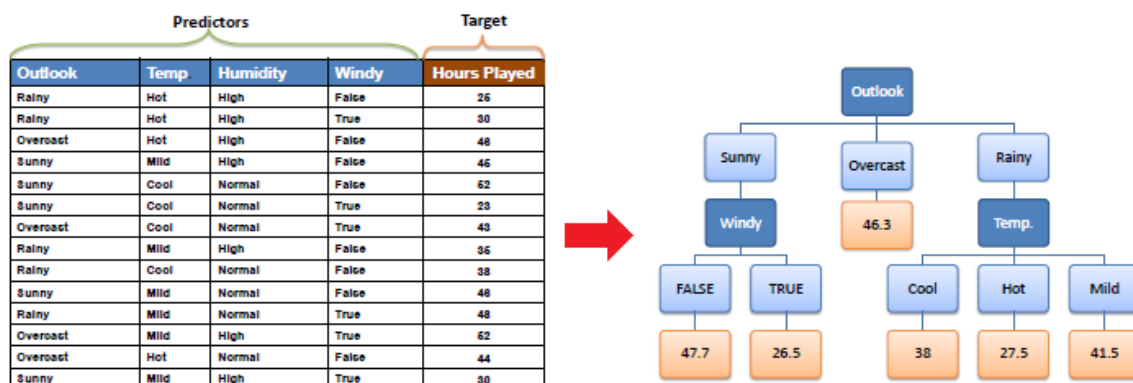
การแบ่งโหนด (Splitting) กระบวนการในการแบ่งโหนด ทำได้ตามเกณฑ์หรือเงื่อนไขต่าง ๆ ซึ่งการถดถอยต้นไม้ตัดสินใจ จะเลือกคุณลักษณะ (feature) หนึ่งตัวเพื่อใช้ตัดสินใจในการแบ่งกลุ่มข้อมูลออกเป็นสองกลุ่มในลักษณะการแบ่งแบบไบนารีที่ซ้ำไปเรื่อย ๆ (recursive binary split) โดยจะเลือกจุดแบ่งที่ทำให้ค่าคลาดเคลื่อนกำลังสอง (Residual Sum of Squares: RSS) ต่ำที่สุด หลักการแบ่งข้อมูลในแต่ละโหนดสำหรับข้อมูลที่มี k ตัวแปร และ n ข้อมูล มีดังนี้

1. เลือกคุณลักษณะ (feature) คือ เลือก 1 ตัวแปรจากตัวแปรทั้งหมด k ตัวแปร แล้วเรียงลำดับค่าของข้อมูลตามคุณลักษณะที่เลือก
2. เลือกจุดแบ่ง (split point) คือ พิจารณาจุดแบ่งที่เป็นไปได้ทั้งหมดจากข้อมูล n ข้อมูล โดยจุดแบ่งจะเป็นตำแหน่งที่อยู่ระหว่างข้อมูลที่เรียงตามลำดับ จำนวนจุดแบ่งจะเท่ากับ $n - \text{min sample}$ จุด เมื่อ min sample คือจำนวนข้อมูลขั้นต่ำที่ต้องการสำหรับแต่ละโหนด
3. คำนวณค่า Residual Sum of Squares (RSS) คือ สำหรับแต่ละจุดแบ่งที่เป็นไปได้ จะคำนวณค่า RSS เพื่อวัดความแตกต่างระหว่างค่าจริงและค่าที่คาดการณ์ จุดแบ่งที่ดีที่สุดคือจุดที่ทำให้ค่า RSS ต่ำที่สุด เนื่องจากแสดงว่าการแบ่งข้อมูลที่จุดนั้นทำให้ข้อมูลในแต่ละกลุ่มมีความแตกต่างกันน้อยที่สุด มีสูตรคำนวณดังนี้

$$\sum_{j=1}^i \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

เมื่อ	R_j	คือ แต่ละกลุ่มของ observation ที่ถูกแบ่งออกมา ทั้งหมด j กลุ่ม
	y_i	คือ ค่าของตัวแปรตาม
	\hat{y}_{R_j}	คือ ค่าทำนายผลลัพธ์ในแต่ละกลุ่ม คำนวณจากค่าเฉลี่ยของตัวแปรตามกลุ่มนั้นๆ

4. ทำซ้ำ โดยการตัดสินใจและการแบ่งข้อมูลจะดำเนินการไปเรื่อย ๆ จนกว่าจะถึงเงื่อนไขสิ้นสุด เช่น ความสูงของต้นไม้ (max depth) หรือจำนวนข้อมูลใน leaf node ต่ำกว่าเกณฑ์ที่กำหนด (min sample)
5. ทำนายค่า คือ เมื่อสร้างต้นไม้เสร็จแล้ว การทำนายค่าสำหรับข้อมูลใหม่จะทำโดยการคำนวณค่าเฉลี่ยของค่าจริงของตัวแปรผลลัพธ์ที่อยู่ใน leaf node นั้น ๆ (Daroontham, 2020)



ภาพที่ 2.8 ตัวอย่างการทำงานของ Decision Tree

หมายเหตุ. จาก https://www.saedsayad.com/decision_tree_reg.htm#:~:text=Decision%20tree%20builds%20regression%20or,decision%20nodes%20and%20leaf%20nodes

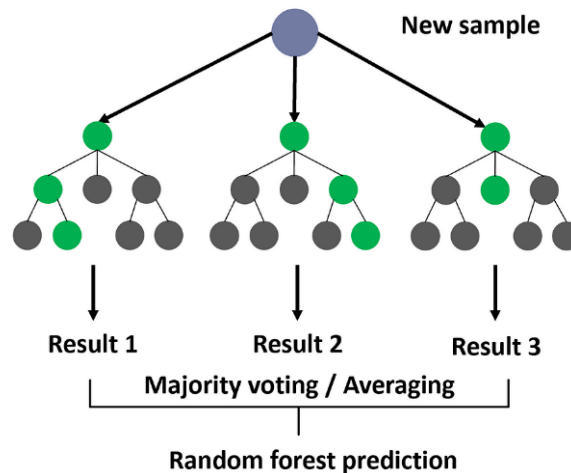
จากภาพที่ 2.8 เป็นการทำนาย Hours Played โดยเริ่มจากโหนดรากคือ Outlook แยกโหนดย่อยได้เป็น Sunny Overcast Rainy และในแต่ละรายการจะมีการแยกย่อยออกไป สุดท้ายผลการทำนายจะอยู่ที่โหนดสุดท้ายที่ไม่มีสามารถแยกย่อยไปอีก ซึ่งเป็ค่าทำนายเรียกว่า Leaf Node ตัวอย่างเช่น Outlook = Rainy แล้วมี Temp. = Hot ผลลัพธ์ของการทำนายจะได้ Hours Played = 27.5

2.1.5.1 การถดถอยพหุผลลัพธ์แบบการสุ่มป่าไม้ (Multi-Output Random Forest Regression: M-RFR)

การถดถอยพหุผลลัพธ์แบบการสุ่มป่าไม้ เป็นวิธีการสร้างแบบจำลองที่สามารถจัดการกับ ข้อมูลที่มีหลายผลลัพธ์ในเวลาเดียวกัน โดยเฉพาะในกรณีที่ผลลัพธ์มีความสัมพันธ์กัน พัฒนามาจากการถดถอยแบบป่าสุ่ม (Random Forest Regressor) แบบผลลัพธ์เดียว สามารถใช้ได้ทั้งในงานที่เกี่ยวข้องกับการจำแนกประเภทและการถดถอย

1) การถดถอยแบบป่าสุ่ม (Random Forest Regressor) เป็นวิธีการรวมแบบจำลองที่ใช้ การถดถอยต้นไม้ตัดสินใจ (Decision Tree Regression) หลายต้นในการทำนายผลลัพธ์ โดยใช้การสุ่มตัวอย่างแบบมีการคืนตัวอย่าง (bootstrap sampling) และการสุ่มเลือกคุณลักษณะเพื่อสร้างความหลากหลายของต้นไม้ในป่า ทำให้ข้อมูลเดียวกันสามารถถูกเลือกซ้ำได้ในแต่ละครั้ง วิธีนี้ช่วยสร้างความหลากหลายของชุดข้อมูลที่ใช้ในการฝึกแต่ละต้นไม้ ซึ่งส่งผลให้แบบจำลองมีความแม่นยำและความเสถียรภาพเพิ่มขึ้น นอกจากนี้การสุ่มเลือกคุณลักษณะและการรวมผลทำนายจากต้นไม้หลายต้นยังช่วยลดการเกิด overfitting ได้ ซึ่งทำให้แบบจำลองสามารถ

ทำนายได้แม่นยำยิ่งขึ้นเมื่อถูกนำไปใช้กับข้อมูลชุดทดสอบหรือข้อมูลใหม่ ผลลัพธ์สุดท้ายของการทำนายจะได้รับการเฉลี่ยค่าทำนายจากต้นไม้แต่ละต้นในป่า (Abebe et al., 2020)



ภาพที่ 2.9 ลักษณะการทำงานของ Random Forest

หมายเหตุ. จาก <https://blog.devgenius.io/learning-random-forest-classification-using-iris-dataset-eeb930612e0e>

2) การถดถอยพหุผลลัพธ์แบบการสุ่มป่าไม้ (Multi-Output Random Forest Regression) เป็นวิธีการสร้างแบบจำลองที่สามารถจัดการกับข้อมูลที่มีหลายผลลัพธ์ โดย Segal (1992) ได้เสนอการสร้างต้นไม้การถดถอยที่สามารถทำนายค่าผลลัพธ์หลายค่าได้พร้อมกันในแต่ละโหนด ซึ่งต้นไม้ถดถอยเหล่านี้มีพื้นฐานมาจากฟังก์ชันการแบ่ง (split function) แบบกำลังสองน้อยที่สุด (least squares) มีหลักการดังนี้

$$\phi(s, t) = SS(t) - SS(t_L) - SS(t_R)$$

เมื่อ $SS(t)$ คือ ผลรวมของค่าความคลาดเคลื่อนกำลังสองในโหนดที่ t

$$SS(t) = \sum_{i \in t} (y_i - \bar{y}(t))^2$$

หลักจากนั้นเพิ่มค่าถ่วงน้ำหนักความแปรปรวนเข้ากับค่าความคลาดเคลื่อนกำลังสอง เพื่อให้การสร้างโหนดมีความเป็นเนื้อเดียวกัน (homogeneous clusters) โดยพิจารณาจากชุดของผลลัพธ์ที่ต้องการทำนาย

$$SS(t) = \sum_{i \in t} (y_i - \bar{y}(t))' V^{-1}(t, \eta) (y_i - \bar{y}(t))$$

โดยที่ $V(t)$ เป็นเมทริกซ์ความแปรปรวนของโหนด t และ η เป็นพารามิเตอร์ที่กำหนดโครงสร้างการถ่วงน้ำหนักความแปรปรวน ดังนั้นการถดถอยพหุผลลัพธ์แบบการสุ่มป่าไม้ สามารถจัดการกับการทำนายหลายตัวแปรได้โดยการเปลี่ยนต้นไม้ที่เป็นการทำนายแบบผลลัพธ์เดียวให้เป็นหลายผลลัพธ์ได้ (Linusson, 2013)

ผลลัพธ์สุดท้ายจะเป็นค่าที่ทำนายสำหรับแต่ละตัวแปรผลลัพธ์ในชุดข้อมูลมาจากค่าเฉลี่ย ของข้อมูลใน Leaf node โดยในงานวิจัยนี้ใช้ Python library ‘sklearn.ensemble’ คำสั่ง MultiOutput Regressor(RandomForstRegressor) ของ scikit-learn ในการทำนายพหุผลลัพธ์ สำหรับพารามิเตอร์หลักในการสร้างแบบจำลองมีดังนี้

- n_estimators: จำนวนต้นไม้ทั้งหมดที่ใช้ในการสร้างแบบจำลอง โดยทั่วไปหากมีค่ามากทำให้การทำนายแม่นยำมาก แต่จะใช้เวลาในสร้าง
- max_depth: ความลึกสูงสุดของต้นไม้ในป่าแต่ละต้น
- min_samples_split: จำนวนตัวอย่างขั้นต่ำที่ต้องใช้ในการแบ่งที่โหนด
- min_samples_leaf: จำนวนตัวอย่างขั้นต่ำที่ต้องใช้ในโหนดใบ
- random_state: ค่าที่ใช้ในการสุ่มชุดข้อมูลในการฝึกสอนและชุดข้อมูลทดสอบ

2.1.5.2 การถดถอยพหุผลลัพธ์แบบเอ็กซ์ตรีมเกรเดียนต์บูสติง (Multi-Output Extreme Gradient Boosting Regression: M-XGBoost)

การถดถอยพหุผลลัพธ์แบบเอ็กซ์ตรีมเกรเดียนต์บูสติง เป็นการทำนายหลายผลลัพธ์โดยใช้แบบจำลอง Extreme Gradient Boosting หรือ XGBoost ที่ได้รับการพัฒนาขึ้นจาก Gradient Boosting โดยเน้นการเพิ่มประสิทธิภาพและความสามารถในการประมวลผลของแบบจำลองที่ใช้ต้นไม้ตัดสินใจ (Decision Trees) แต่ละต้นไม้จะพยายามแก้ไขข้อผิดพลาดของต้นไม้ก่อนหน้า ทำให้แบบจำลองมีความแม่นยำมากขึ้น โดยใช้ Loss Function (Bentejac et al., 2019) มีขั้นตอนการทำงานของแบบจำลองการถดถอย XGBoost ดังนี้

1. กำหนดฟังก์ชันเริ่มต้น $f_0(x) = \arg \min_{\beta} \sum_{i=1}^n L(y_i, \beta)$ ซึ่งเป็นฟังก์ชันที่ทำให้การสูญเสีย $L(y_i, \beta)$ ต่ำที่สุดเท่าที่จะเป็นไปได้จากข้อมูลทั้งหมด

$$L_{xgb}(\hat{f}_m) = \sum_{i=1}^N L(y_i, \hat{f}_m(x_i))$$

เมื่อ $L_{xgb}(\hat{f}_m)$ คือ ฟังก์ชันการสูญเสียที่วัดจากความต่างระหว่างค่าทำนายกับค่าจริงในรอบที่ m
 N คือ จำนวนข้อมูลทั้งหมดในชุดฝึก

$L(y_i, \hat{f}_m(x_i))$ คือ ค่าความแตกต่างระหว่างค่าจริง y_i กับค่าที่ทำนายจากแบบจำลอง สำหรับข้อมูลลำดับที่ i

β คือ ค่าทำนาย

โดยต้นไม้ต้นถัดไปจะพยายามแก้ไขข้อผิดพลาดจากต้นก่อนหน้า เมื่อตัวแปร m คือจำนวนรอบที่ได้ทำการปรับปรุงแบบจำลองไปแล้ว หรือเป็นลำดับของการเพิ่มต้นไม้การตัดสินใจ

2. เรียนรู้ซ้ำ สำหรับแต่ละรอบ $m = 1, \dots, M$:

$$L_{xgb}(\hat{f}_m) = \sum_{i=1}^N (L(y_i, \hat{f}_{m-1}(x_i)) + h_m(x_i))$$

เมื่อ $L(y_i, \hat{f}_{m-1}(x_i))$ คือ ค่าความแตกต่างระหว่างค่าจริง y_i กับค่าที่ทำนายจากแบบจำลอง สำหรับข้อมูลลำดับที่ i ในรอบก่อนหน้า ($m - 1$)

$h_m(x_i)$ คือ ค่าการปรับปรุงที่แบบจำลองทำในรอบที่ m เพื่อลดข้อผิดพลาดที่เหลือจากรอบก่อนหน้า

ซึ่ง XGboost ใช้การขยายเทย์เลอร์ (Taylor expansion) ของฟังก์ชันการสูญเสียรอบ \hat{f}_m โดยพิจารณาพหุนามเทย์เลอร์ลำดับที่สองเพื่อช่วยประมาณค่าจุดต่ำที่สุด

$$L_{xgb}(\hat{f}_m) \approx \sum_{i=1}^N (L(y_i, \hat{f}_{m-1}(x_i)) + G_m(x_i) \cdot h_m(x_i) + \frac{1}{2} H_m(x_i) \cdot h_m^2(x_i))$$

2.1 คำนวณค่าอนุพันธ์อันดับหนึ่งและอนุพันธ์อันดับสองของฟังก์ชันการสูญเสีย

เมื่อ $G_m(x_i)$ คือ อนุพันธ์อันดับหนึ่งของฟังก์ชันการสูญเสีย และ $H_m(x_i)$ คือ อนุพันธ์อันดับสองของฟังก์ชันการสูญเสีย โดยสองค่านี้ช่วยให้ทราบถึงทิศทางและความโค้งของฟังก์ชันการสูญเสียที่ใช้ในการปรับปรุงการทำนาย มีสมการดังนี้

$$\text{Gradient} \quad G_m(x_i) = \left[\frac{\partial L(y_i, \hat{f}(x_i))}{\partial \hat{f}(x_i)} \right]_{\hat{f}(x_i) = \hat{f}_{m-1}(x_i)}$$

$$\text{Hessian} \quad H_m(x_i) = \left[\frac{\partial^2 L(y_i, \hat{f}(x_i))}{\partial \hat{f}(x_i)^2} \right]_{\hat{f}(x_i) = \hat{f}_{m-1}(x_i)}$$

2.2 สร้างแบบจำลองการถดถอยใหม่ h_m โดยใช้ชุดข้อมูล x_i จาก Gradient และ Hessian ที่คำนวณได้ เพื่อให้ต้นไม้สามารถปรับปรุงฟังก์ชันการพยากรณ์ได้ดีขึ้น

$$\left\{ \left(x_i - \frac{\sum_{i \in I_j} G_m(x_i)}{\sum_{i \in I_j} H_m(x_i) + \lambda} \right) : i = 1, \dots, n \right\}$$

เมื่อ $\sum_{i \in I_j} G_m(x_i)$ คือ ผลรวมของ Gradient สำหรับข้อมูลทั้งหมดที่อยู่ในใบ j

$\sum_{i \in I_j} H_m(x_i)$ คือ ผลรวมของ Hessian สำหรับข้อมูลทั้งหมดที่อยู่ในใบ j

λ คือ ค่าพารามิเตอร์ regularization ที่ช่วยควบคุมไม่ให้ค่าน้ำหนักที่กำหนดให้กับแต่ละใบของต้นไม้ (Leaf Weights) ไม่ให้สูงเกินไป

j คือ จำนวนใบของต้นไม้

ในแต่ละโหนดของต้นไม้จะมีการเลือกจุดแบ่งที่เพิ่มความแม่นยำสูงสุด ซึ่ง Gain เป็นค่าที่ใช้วัดประสิทธิภาพของการแบ่งใบการแบ่ง หากจุดใดมีค่าสูงจะถือว่าความสามารถในการลดความสูญเสียมาก มีสมการดังนี้

$$\text{Gain} = \frac{1}{2} \left[\frac{(\sum_{i \in I_l} G_m(x_i))^2}{\sum_{i \in I_l} H_m(x_i) + \lambda} + \frac{(\sum_{i \in I_r} G_m(x_i))^2}{\sum_{i \in I_r} H_m(x_i) + \lambda} - \frac{(\sum_{i \in I} G_m(x_i))^2}{\sum_{i \in I} H_m(x_i) + \lambda} \right] - \gamma$$

เมื่อ I_l และ I_r คือ กลุ่มของข้อมูลที่อยู่ในกลุ่มซ้ายและกลุ่มขวาหลังจากการแบ่งตามเกณฑ์ที่กำหนด

I คือ กลุ่มข้อมูลในใบก่อนการแบ่ง

γ คือ ค่าพารามิเตอร์ regularization ที่กำหนดค่า minimum Gain ที่ต้องการสำหรับการสร้างใบใหม่ ถ้า Gain ของการแบ่งมีค่าน้อยกว่าค่า γ แบบจำลอง จะไม่สร้างใบใหม่ในส่วนนั้นเพื่อลดความซับซ้อนของต้นไม้

หลังจากนั้นจะทำการอัปเดตฟังก์ชันการทำนาย $\hat{f}_m = \hat{f}_{m-1} + \eta h_m$ เมื่อ η คือ ค่าพารามิเตอร์ควบคุมขนาดการเปลี่ยนแปลงที่เกิดจากต้นไม้ใหม่ หรือ learning rate เป็นค่าคงที่ที่ใช้ควบคุมอัตราการเรียนรู้ เพื่อไม่ให้เกิดการอัปเดตแต่ละรอบมีขนาดใหญ่เกินไป

3. เมื่อทำซ้ำทุกรอบ M แล้ว ผลลัพธ์คือฟังก์ชันการทำนาย \hat{f}_M ที่รวมการเรียนรู้จากทุกต้นไม้เข้าด้วยกัน (Schagen, 2023)

การถดถอยพหุผลลัพธ์แบบเอ็กตรีเมียร์เกรเดียนต์บูสติง สามารถทำได้โดยใช้แบบจำลองรวม (Combined Model) ที่ฝึกตัวแปรผลลัพธ์ทั้งหมดพร้อมกัน ซึ่งแบบจำลองจะพยายามประมาณฟังก์ชันที่มีหลายตัวแปร $f: \mathbb{R}^p \rightarrow \mathbb{R}^n$ ในการสร้างแบบจำลองรวม ต้นไม้การถดถอยจะมีใบที่เป็น

เวกเตอร์แทนที่จะเป็นสเกลาร์ ซึ่งหมายความว่าแต่ละใบของต้นไม้ไม่สามารถให้ค่าทำนายสำหรับหลายผลลัพธ์ได้ในเวลาเดียวกัน มีหลักการดังนี้

1. กำหนดฟังก์ชันเริ่มต้น $\hat{f}_0(x) = \arg \min_{\beta} \sum_{i=1}^n L(y_i, \beta)$ ซึ่งเป็นฟังก์ชันที่ทำให้การสูญเสีย $L(y_i, \beta)$ ต่ำที่สุดเท่าที่จะเป็นไปได้จากข้อมูลทั้งหมด

ฟังก์ชันที่มีหลายตัวแปรได้ สำหรับฟังก์ชันการสูญเสียใด ๆ นั่นคือเมทริกซ์ที่มีหลายมิติ ใช้การคูณเทนเซอร์ (Tensor Product) เมื่อ U, V และ W เป็นเวกเตอร์สเปซ

$$\dim(W) = \dim(U) \cdot \dim(V)$$

ทำให้การขยายเทย์เลอร์ (Taylor expansion) สามารถการประมาณค่าการสูญเสียในบริบทของการทำนายหลายผลลัพธ์ได้

$$f(x) = f(x_0) + (Df(x_0))(x - x_0) + \cdots + \frac{D^k f(x_0)}{k!} (x - x_0)^{\otimes k} + R(x_0)$$

เมื่อ	$f(x)$	คือ ฟังก์ชันที่ต้องการประมาณค่า
	$f(x_0)$	คือ ค่าของฟังก์ชันที่จุดเริ่มต้น x_0
	$Df(x_0)$	คือ อนุพันธ์อันดับแรกของฟังก์ชันที่จุด x_0
	$x - x_0$	คือ ความแตกต่างระหว่างค่าปัจจุบันกับค่าที่จุดเริ่มต้น
	$\frac{D^k f(x_0)}{k!}$	คือ อนุพันธ์อันดับ k
	$\otimes k$	คือ การคูณแบบ Tensor Product โดยทำซ้ำ k ครั้ง

หลังจากนั้นทำให้เป็นมาตรฐานโดยการใช้สองพจน์แรก (2-jet) ในการประมาณค่าฟังก์ชันเริ่มต้น มีสมการดังนี้

$$L(\hat{f}_m) = L(\hat{f}_{m-1}) + h_m \cdot DL(\hat{f}_{m-1}) + \frac{1}{2} h_m^{\otimes 2} \cdot D^2 L(\hat{f}_{m-1})$$

เมื่อ	$L(\hat{f}_{m-1})$	คือ ฟังก์ชันการสูญเสียที่คำนวณจากค่าทำนายในรอบก่อนหน้า
	h_m	คือ ค่าการปรับปรุงที่แบบจำลองทำในรอบที่ m
	$DL(\hat{f}_{m-1})$	คือ อนุพันธ์แรกของฟังก์ชันความสูญเสียที่จุด \hat{f}_{m-1}
	$D^2 L(\hat{f}_{m-1})$	คือ อนุพันธ์อันดับสองของฟังก์ชันความสูญเสียที่จุด \hat{f}_{m-1}

2. เรียนรู้ซ้ำ สำหรับแต่ละรอบ $m = 1, \dots, M$:

ฟังก์ชันการสูญเสียในสำหรับของการทำนายหลายผลลัพธ์ โดยการรวมทุกมิติของใบ (leaves) ช่วยในการคำนวณฟังก์ชันสูญเสียของต้นไม้

$$\hat{L}(\hat{y}_m) = -\frac{1}{2} \sum_{i=1}^u \left(\frac{\sum_j ((G_m^i)_j)^2}{\sum_j ((\hat{H}_m^i)_j) + \lambda} \right) + \gamma T_m$$

เมื่อ $\hat{L}(\hat{y}_m)$ คือ ฟังก์ชันการสูญเสียที่คำนวณจากค่าทำนายในรอบที่ m
 u คือ จำนวนของมิติของตัวแปรผลลัพธ์ หรือ จำนวนตัวแปรผลลัพธ์
 T_m คือ จำนวนใบ (leaf nodes) ของต้นไม้ รอบที่ m

2.1 คำนวณค่าอนุพันธ์อันดับหนึ่งและอนุพันธ์อันดับสองของฟังก์ชันการสูญเสียสำหรับหลายผลลัพธ์

$$\text{Gradient} \quad G(\hat{y}) = \left(\frac{\partial L}{\partial (\hat{y})_1}(\hat{y}), \dots, \frac{\partial L}{\partial (\hat{y})_u}(\hat{y}) \right)$$

$$\text{Hessian} \quad H(\hat{y}) = \begin{pmatrix} \frac{\partial^2 L}{\partial (\hat{y})_1^2}(\hat{y}) & \cdots & \frac{\partial^2 L}{\partial (\hat{y})_1 \partial (\hat{y})_u}(\hat{y}) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 L}{\partial (\hat{y})_u \partial (\hat{y})_1}(\hat{y}) & \cdots & \frac{\partial^2 L}{\partial (\hat{y})_u^2}(\hat{y}) \end{pmatrix}$$

2.2 สร้างแบบจำลองการถดถอยใหม่ h_m โดยใช้ชุดข้อมูล x_i จาก Gradient และ Hessian ที่คำนวณได้ เพื่อให้ต้นไม้สามารถปรับปรุงฟังก์ชันการพยากรณ์ได้ดีขึ้น

$$\left\{ \left(x_i, -\frac{\sum_j (G_m^i)_j}{\sum_j ((\hat{H}_m^i)_j) + \lambda} \right) : i = 1, \dots, n \right\}$$

ในแต่ละโหนดของต้นไม้จะมีการเลือกจุดแบ่งที่เพิ่มความแม่นยำสูงสุด ซึ่ง Gain เป็นค่าที่ใช้วัดประสิทธิภาพของการแบ่งใบ หากจุดใดมีค่าสูงจะถือว่าความสามารถในการลดความสูญเสียมาก มีสมการดังนี้

$$\text{Gain} = \frac{1}{2} \left[\sum_{i=1}^u \left(\frac{\sum_{j \in l_l} ((G_m^i)_j)^2}{\sum_{j \in l_l} ((\hat{H}_m^i)_j) + \lambda} \right) + \sum_{i=1}^u \left(\frac{\sum_{j \in l_r} ((G_m^i)_j)^2}{\sum_{j \in l_r} ((\hat{H}_m^i)_j) + \lambda} \right) - \sum_{i=1}^u \left(\frac{\sum_{j \in l} ((G_m^i)_j)^2}{\sum_{j \in l} ((\hat{H}_m^i)_j) + \lambda} \right) \right] - \gamma$$

หลังจากนั้นจะทำการอัปเดตฟังก์ชันการทำนาย $\hat{f}_m = \hat{f}_{m-1} + \eta h_m$ เมื่อ η คือค่าพารามิเตอร์ควบคุมขนาดการเปลี่ยนแปลงที่เกิดจากต้นไม้ใหม่ หรือ learning rate เป็นค่าคงที่ที่ใช้ควบคุมอัตราการเรียนรู้ เพื่อไม่ให้เกิดการอัปเดตแต่ละรอบมีขนาดใหญ่เกินไป

3. เมื่อทำซ้ำหุ้รอบ M หรือจนกว่าจะถึงเงื่อนไขหยุด เช่น จนกว่าจะถึงเงื่อนไขการหยุดฟังก์ชันการสูญเสียลดลงหรือครบเงื่อนไขของพารามิเตอร์ที่กำหนด ดังนั้นผลลัพธ์คือฟังก์ชันการทำนาย \hat{f}_M ที่รวมการเรียนรู้จากทุกต้นไม้เข้าด้วยกัน

การถดถอยพหุผลลัพธ์แบบเอ็กตรีเมียร์เกรเดียนต์บูสติง มีความสามารถในการจับความสัมพันธ์ระหว่างตัวแปรผลลัพธ์ได้ดี เนื่องจากใช้โมเดลเดียวในการพยากรณ์หลายผลลัพธ์พร้อมกัน การใช้ฟังก์ชันการสูญเสียที่เหมาะสมร่วมกับ L2 Regularization ช่วยลดความซับซ้อนของต้นไม้และลดโอกาสในการเกิด overfitting นอกจากนี้ การคำนวณอนุพันธ์อันดับแรกและอันดับสองทำให้สามารถหาจุดแบ่งได้อย่างรวดเร็วและแม่นยำ ในขณะที่การอัปเดตพารามิเตอร์ในแต่ละรอบช่วยเพิ่มประสิทธิภาพในการเรียนรู้ ทำให้กระบวนการสร้างแบบจำลองรวดเร็วขึ้น (Schagen, 2023) โดยในงานวิจัยนี้ใช้ Python library 'xgboost' ร่วมกับ MultiOutputRegressor คำสั่ง MultiOutputRegressor(xgb) ของ scikit-learn ในการทำนายพหุผลลัพธ์ สำหรับพารามิเตอร์หลักในการสร้างแบบจำลองมีดังนี้

- n_estimators: จำนวนต้นไม้ที่ใช้ในการสร้างแบบจำลอง
- eta: ควบคุมความเร็วในการเรียนรู้ของแบบจำลอง
- max_depth: ความลึกสูงสุดของต้นไม้แต่ละต้น
- min_child_weight: ค่าขั้นต่ำของน้ำหนักที่ใช้ในการแบ่งโหนด
- subsample: อัตราการสุ่มข้อมูลตัวอย่างในแต่ละรอบการสร้างต้นไม้

2.1.5.3 การถดถอยพหุผลลัพธ์แบบไลต์เกรเดียนต์บูสติง (Multi-Output Light Gradient Boosting Regression: M-LightGBM)

แบบจำลอง LightGBM เป็นแบบจำลองที่มีโครงสร้างเป็นแบบต้นไม้หลาย ๆ ต้น โดยต้นไม้เหล่านี้จะถูกสร้างขึ้นจากข้อมูลฝึก คือ แบบจำลองจะทำการค้นหาตัวแปรอิสระที่ส่งผลต่อตัวแปรตามโดยจะเปรียบเทียบความสามารถในการแบ่งกลุ่มค่าผลลัพธ์ของตัวแปรตาม ตามค่าตัวแปรอิสระที่เปลี่ยนแปลงไป จากนั้นจะสร้างเงื่อนไขว่า ถ้าตัวแปรอิสระดังกล่าวมีค่าอยู่ในช่วงแต่ละช่วง ตัวแปรตามควรมีค่าเท่าไร โดยจะใช้ค่าเฉลี่ยของตัวแปรตามที่ถูกจัดกลุ่มอยู่ในกลุ่มเดียวกัน ต่อมาแบบจำลองจะตรวจสอบการแบ่งกลุ่มว่าจะสามารถแยกค่าตัวแปรตามที่มีค่าแตกต่างกันออกจากกันและจับกลุ่มค่าตัวแปรตามไว้ให้อยู่ใกล้เคียงกัน จะดีกว่าการไม่แบ่งกลุ่มหรือไม่ ถ้าหากไม่ดีกว่าจะไม่ทำการแบ่งกลุ่มนั้น เมื่อได้ตัวแปรอิสระที่ใช้ในการแบ่งกลุ่มแล้ว แบบจำลองจะเริ่มสร้างต้นไม้โดยใช้ตัวแปรอิสระดังกล่าวแตกกิ่งค่าความเป็นไปได้ของตัวแปรตามออกมา จากนั้นแบบจำลองจะทำซ้ำตามทีกล่าวมากับตัวแปรอิสระอื่น ๆ ที่ยังไม่ถูกเลือกจนได้ต้นไม้ที่ประกอบด้วยกิ่งจำนวนมาก ซึ่งจะมีความ

ละเอียดในการทำนายค่าตัวแปรตาม โดยแบบจำลอง LightGBM มีความสามารถในการจัดการข้อมูลขนาดใหญ่ และทำงานได้เร็วมากเมื่อเทียบกับ Gradient Boosting (FINNOMENA, 2565) และคุณสมบัติที่โดดเด่นของ LightGBM คือความสามารถในการฝึกโดยอิงจากฟังก์ชันการสูญเสีย (loss function) ที่แตกต่างกัน

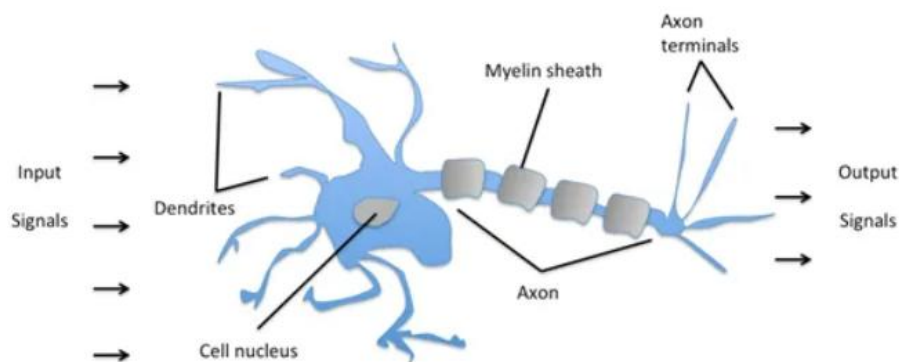
โดยหลักการหลักของ LightGBM คือใช้กลยุทธ์การเติบโตของต้นไม้แบบ Leaf-wise ซึ่งต่างจาก Gradient Boosting แบบเดิมที่ใช้ Level-wise โดย Leaf-wise จะเลือกขยายใบที่ลดค่าฟังก์ชันการสูญเสียได้มากที่สุดก่อน จะทำให้การเรียนรู้ที่รวดเร็วและมีประสิทธิภาพมากขึ้น และ LightGBM ใช้เทคนิค Histogram-based Algorithm ที่ลดการคำนวณที่จำเป็นในการหาจุดแบ่ง (split points)

โดยในงานวิจัยนี้ใช้ Python library ‘lightgbm’ ร่วมกับ MultiOutputRegressor คำสั่ง MultiOutputRegressor(lightgbm) ของ scikit-learn ในการทำนายพหุผลลัพธ์ สำหรับพารามิเตอร์หลักในการสร้างแบบจำลองมีดังนี้

- learning_rate: อัตราการเรียนรู้ที่ใช้ในการปรับค่าถ่วงน้ำหนักของแบบจำลอง
- n_estimators: จำนวนต้นไม้ที่ใช้ในการสร้างแบบจำลอง
- max_depth: ความลึกสูงสุดของต้นไม้แต่ละต้น
- subsample: อัตราการสุ่มข้อมูลตัวอย่างในแต่ละรอบการสร้างต้นไม้
- min_child_weight: ค่าขั้นต่ำของน้ำหนักที่ใช้ในการแบ่งโหนด
- num_leaves: จำนวนใบสูงสุดในแต่ละต้นไม้
- min_data_in_leaf: จำนวนข้อมูลขั้นต่ำในแต่ละใบ

2.1.5.4 การถดถอยพหุผลลัพธ์โครงข่ายประสาทเทียม (Multi-Output artificial neural network: M-ANN)

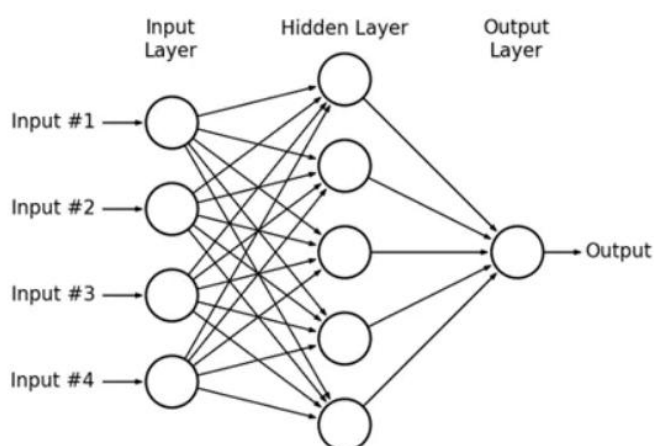
แบบจำลอง Artificial Neural Network (ANN) เป็นหนึ่งในวิธีการของการเรียนรู้เชิงลึกซึ่งอยู่ภายใต้แนวทางของการเรียนรู้แบบมีผู้สอนที่สามารถนำไปประยุกต์ใช้กับงานได้หลายด้าน เช่น การทำนาย การจัดกลุ่ม เป็นต้น โดยแนวคิดหลักของ ANN ได้รับแรงบันดาลใจจากโครงสร้างและการทำงานของสมองมนุษย์ ซึ่งพยายามเลียนแบบการทำงานของเซลล์ประสาทเพื่อประมวลผลข้อมูล แสดงดังภาพที่ 2.10



ภาพที่ 2.10 แผนผังเซลล์ประสาทมนุษย์

หมายเหตุ. จาก https://www.researchgate.net/publication/343921333_Modeling_and_Predicting_of_Motor_Insurance_Claim_Amount_using_Artificial_Neural_Network

ANN ประกอบไปด้วย โหนด หลายโหนด ซึ่งคล้ายกับเซลล์ประสาทของมนุษย์ที่เชื่อมต่อกันด้วยลิงค์และโต้ตอบกัน มีส่วนประกอบสำคัญ คือ ข้อมูลป้อนเข้า (Input layer) ชั้นซ่อน (Hidden Layer) และข้อมูลส่งออก (Output layer) ซึ่งมีหลักการทำงานดังภาพที่ 2.11



ภาพที่ 2.11 ตัวอย่างการทำงานแบบจำลอง ANN

หมายเหตุ. จาก https://www.researchgate.net/publication/343921333_Modeling_and_Predicting_of_Motor_Insurance_Claim_Amount_using_Artificial_Neural_Network

จากภาพที่ 2.11 โครงสร้างพื้นฐานของ ANN จะประกอบไปด้วย 3 ส่วนหลัก **ชั้นข้อมูลป้อนเข้า (Input layer)** เป็นชั้นที่จะทำการป้อนข้อมูลเข้า จำนวนโหนดในชั้นนี้จะขึ้นอยู่กับข้อมูลที่นำเข้าแบบจำลองหรือตัวแปรอิสระที่จะนำเข้าสู่แบบจำลอง

ชั้นซ่อน (Hidden Layer) คือชั้นที่อยู่ระหว่างชั้นข้อมูลนำเข้าและชั้นข้อมูลส่งออก ซึ่งมีบทบาทสำคัญในการเรียนรู้ของแบบจำลอง ในชั้นซ่อนสามารถมีหลายชั้นและแต่ละชั้นสามารถมีจำนวนโหนดเพอร์เซพตรอน (Perceptron) ได้ตามต้องการ เพอร์เซพตรอนเป็นหน่วยที่รับข้อมูลจากชั้นก่อนหน้า แล้วคำนวณและส่งผลลัพธ์ผ่านฟังก์ชันกระตุ้น (Activation Function)

ชั้นข้อมูลส่งออก (Output layer) เป็นชั้นที่จะนำข้อมูลจากการคำนวณชั้นที่ผ่านมาไปใช้ จำนวนโหนดในชั้นนี้จะขึ้นอยู่กับตัวแปรตาม ซึ่งสามารถมีมากกว่า 1 โหนดได้ (อัศพรพล พรหมพิริยะ พงษ์, 2566)

ซึ่ง ANN สามารถทำนายแบบพหุผลลัพธ์ได้โดยใช้ MLPRegressor (Multi-Layer Perceptron Regressor) ใน scikit-learn ซึ่งเป็นแบบจำลองปัญญาประดิษฐ์ประเภท ANN ออกแบบมาเพื่อใช้ในการแก้ปัญหาการถดถอย ซึ่งคาดการณ์ตัวแปรตามแบบต่อเนื่องและสามารถทำนายผลลัพธ์ได้หลายตัวแปรโดยไม่ต้องใช้ร่วมกับ Multi-Output Regression ดังนั้น MLPRegressor เป็นทางเลือกที่ดีสำหรับการแก้ปัญหาทำนายค่าตัวเลขที่ต้องการแบบจำลองที่มีความยืดหยุ่นสูงและสามารถจัดการกับความซับซ้อนของข้อมูลได้ซึ่งในงานวิจัยจะใช้พารามิเตอร์หลักดังนี้

- hidden_layer_sizes: การกำหนดจำนวนชั้นและจำนวนหน่วยประสาทในแต่ละ hidden layer เช่น (4,2) หมายความว่าในชั้นซ่อนจะมีทั้งหมด 2 ชั้น โดยชั้นแรกจะมี 4 โหนด และ ชั้นที่สองจะมี 2 โหนด
- solver: การปรับน้ำหนักของแบบจำลอง เช่น lbfgf ใช้สำหรับปัญหา optimization และ adam ใช้สำหรับการฝึกโครงข่ายประสาท
- activation: ฟังก์ชันที่ใช้ในการแปลงข้อมูลจากแต่ละชั้นของโครงข่ายประสาทเทียมก่อนส่งไปยังชั้นถัดไป
- learning_rate: อัตราการเรียนรู้ที่ใช้ในการปรับพารามิเตอร์

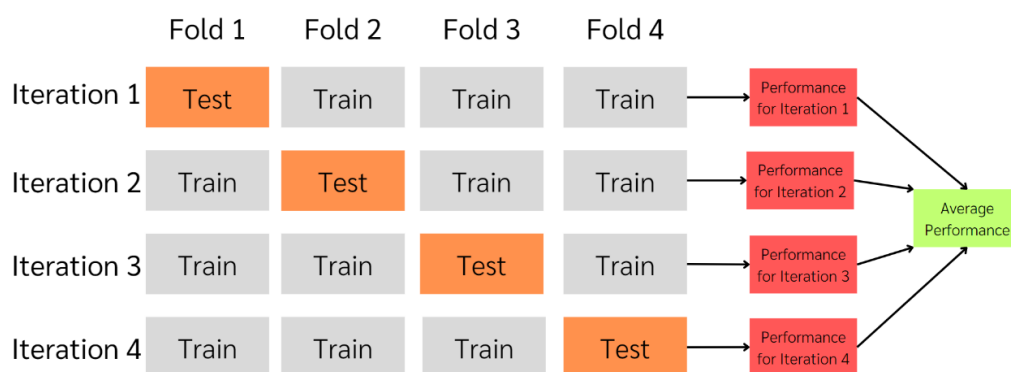
2.1.6 การฝึกแบบจำลอง (Model train)

การฝึกสอนแบบจำลองการเรียนรู้ของเครื่องหมายถึงกระบวนการที่ใช้ข้อมูลการฝึกเพื่อปรับพารามิเตอร์ต่าง ๆ ของแบบจำลองให้สามารถทำนายผลลัพธ์ได้อย่างแม่นยำ การฝึกแบบจำลองจะทำงานกับข้อมูลที่มีป้ายกำกับ ซึ่งประกอบด้วยคุณลักษณะและค่าผลลัพธ์ที่ต้องการ โดยผ่านขั้นตอนการเรียนรู้หลายรอบเพื่อค้นหาพารามิเตอร์ที่เหมาะสมที่สุด ซึ่งสามารถลดความผิดพลาดระหว่างผลลัพธ์ที่คาดการณ์และผลลัพธ์จริง ในงานวิจัยนี้จะใช้เทคนิค K-Fold Cross-

Validation โดยใช้ชุดข้อมูลฝึกอบรมที่แบ่งมา 80% สำหรับการฝึกและทดสอบแบบจำลอง เพื่อตรวจสอบความแม่นยำของแบบจำลองและปรับปรุงประสิทธิภาพ

K-fold Cross-validation เป็นเทคนิคการเรียนรู้ของเครื่องที่ใช้เพื่อประเมินประสิทธิภาพของแบบจำลอง โดยเฉพาะเมื่อข้อมูลมีจำกัด เทคนิคนี้จะทำการแบ่งข้อมูลออกเป็น K ส่วน ซึ่งเรียกว่า folds โดยในแต่ละรอบการประเมินหนึ่ง fold จะถูกใช้เป็นชุดทดสอบ (Test Set) หรือชุดตรวจสอบ (Validation Set) ขณะที่ folds ที่เหลือจะรวมกันเป็นชุดฝึก (Training Set) แบบจำลอง จะได้รับการฝึกและทดสอบ K ครั้ง โดยการเปลี่ยน fold ที่ใช้เป็นชุดทดสอบในแต่ละครั้ง วิธีการนี้ ช่วยให้การประเมินแบบจำลองมีความแม่นยำและสม่ำเสมอมากขึ้น แสดงดังภาพที่ 2.12

หลังจากเสร็จสิ้นกระบวนการ K-fold Cross-validation แบบจำลองจะถูกทดสอบเพิ่มเติมกับข้อมูลที่ไม่ได้ใช้ในการฝึกหรือในการ Cross-validation ซึ่งเรียกว่า test set เพื่อประเมินความแม่นยำที่แท้จริงของแบบจำลองในสภาพแวดล้อมที่ไม่เคยเห็นมาก่อน



ภาพที่ 2.12 กระบวนการทำงานของ 4-fold Cross-validation

2.1.7 การปรับไฮเปอร์พารามิเตอร์ (Hyperparameter Tuning)

การปรับไฮเปอร์พารามิเตอร์เป็นกระบวนการสำคัญในการพัฒนาแบบจำลองการเรียนรู้ของเครื่อง เนื่องจากไฮเปอร์พารามิเตอร์ที่เลือกใช้นั้นมีผลต่อประสิทธิภาพของแบบจำลองอย่างมาก ไม่น้อยไปกว่าพารามิเตอร์ที่เรียนรู้จากข้อมูล การกำหนดค่าไฮเปอร์พารามิเตอร์จะทำล่วงหน้าก่อนการฝึกแบบจำลอง ซึ่งหากไฮเปอร์พารามิเตอร์ถูกตั้งค่าอย่างเหมาะสม จะช่วยเพิ่มประสิทธิภาพการทำนายของแบบจำลองได้แม่นยำยิ่งขึ้น

ไฮเปอร์พารามิเตอร์ คือค่าที่ผู้ใช้กำหนดเองก่อนที่แบบจำลองจะเริ่มการเรียนรู้ เช่น สำหรับ Random Forest การตั้งค่า `n_estimators` เป็นจำนวนต้นไม้ที่ใช้ในการสร้างแบบจำลอง การกำหนดค่าไฮเปอร์พารามิเตอร์อย่างเหมาะสมมีความสำคัญเพื่อควบคุมการฝึกแบบจำลองและให้

ได้ผลลัพธ์ที่ดีที่สุด โดยวิธีค้นหาไฮเปอร์พารามิเตอร์ วิธีการค้นหาไฮเปอร์พารามิเตอร์ที่เหมาะสมมีหลายวิธี เช่น Manual Search Grid Search และ Random Search มีรายละเอียดดังนี้

1) Manual Search คือวิธีการที่เลือกค่าไฮเปอร์พารามิเตอร์จากประสบการณ์และความคิดเห็นส่วนตัว โดยสร้างแบบจำลองจากค่าที่เลือกและทำการวัดความแม่นยำของแบบจำลองนั้น ๆ เป็นระยะ ๆ จนกว่าจะได้ค่าความแม่นยำที่พึงพอใจ วิธีนี้ขึ้นอยู่กับความรู้และประสบการณ์ของผู้พัฒนาเพื่อกำหนดค่าไฮเปอร์พารามิเตอร์ที่ดีที่สุดสำหรับแบบจำลอง

2) Grid Search คือ เทคนิคการค้นหาแบบกริดที่ใช้ในการค้นหาค่าไฮเปอร์พารามิเตอร์ที่เหมาะสมอย่างเป็นระบบ โดยจะสร้างตารางของค่าที่กำหนดไว้ล่วงหน้าและทำการทดสอบแต่ละชุดพารามิเตอร์อย่างครบถ้วน จากนั้นจะประเมินประสิทธิภาพของแบบจำลองที่สร้างขึ้นจากชุดพารามิเตอร์เหล่านั้น เมื่อการทดสอบครบทุกชุดแล้ว แบบจำลองที่มีชุดพารามิเตอร์ที่ให้ผลลัพธ์แม่นยำที่สุดจะถือว่าดีที่สุด แม้ว่าจะเป็นวิธีที่เข้าใจง่ายและตรงไปตรงมา แต่การใช้ Grid Search อาจใช้เวลาค่อนข้างมาก โดยเฉพาะเมื่อจำนวนชุดพารามิเตอร์และความละเอียดของกริดสูง (BDI, 2024)

3) Random Search เป็นเทคนิคที่ใช้ในการค้นหาค่าพารามิเตอร์ที่ดีที่สุดสำหรับแบบจำลองการเรียนรู้ของเครื่อง โดยการสุ่มเลือกค่าพารามิเตอร์จากช่วงที่กำหนดไว้ ซึ่งแตกต่างจาก Grid Search ที่จะทดสอบทุกค่าที่เป็นไปได้ในช่วงที่กำหนดไว้ Random Search จะเลือกทดสอบค่าบางส่วนที่สุ่มมาเท่านั้น ทำให้ใช้เวลาน้อยกว่าและสามารถค้นหาพารามิเตอร์ที่เหมาะสมได้อย่างรวดเร็ว (Bergstra & Bengio, 2012)

2.1.8 การประเมินและเปรียบเทียบประสิทธิภาพของแบบจำลอง (Model Evaluation and Comparison)

2.1.8.1 การประเมินประสิทธิภาพของแบบจำลอง (Model Evaluation)

เพื่อประเมินและเปรียบเทียบประสิทธิภาพแบบจำลองการถดถอยพหุผลลัพธ์โดยทั่วไป จะคำนวณค่าเฉลี่ยของความคลาดเคลื่อนในการทำนายตัวแปรตามแต่ละตัวแปรแยกกัน โดยมักจะวัดจากค่าดังนี้ 1. ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean Square Error: MSE) 2. ค่ารากที่สองของค่าเฉลี่ยค่าความคลาดเคลื่อนกำลังสอง (Root Mean Squared Error: RMSE) เนื่องจากค่า MSE เป็นค่าวัดความถูกต้องของการทำนายที่วัดจากการยกกำลังสองของค่าความคลาดเคลื่อน โดยจะให้ความสำคัญกับค่าทำนายที่แตกต่างไปจากค่าจริงมาก จึงทำให้อ่อนไหวต่อค่าความคลาดเคลื่อนที่มีขนาดใหญ่ หากมีค่าความคลาดเคลื่อนสูงเมื่อยกกำลังจะทำให้มีค่ามาก จึงนิยมใช้ RMSE ในการเปรียบเทียบประสิทธิภาพของแบบจำลอง

1) ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean Square Error: MSE) มีสมการดังนี้

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

2) ค่ารากที่สองของค่าเฉลี่ยค่าความคลาดเคลื่อนกำลังสอง (Root Mean Squared Error: RMSE) มีสมการดังนี้

$$RMSE = \sqrt{MSE}$$

การเปรียบเทียบประสิทธิภาพของทั้งค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (MSE) และค่ารากที่สองของค่าเฉลี่ยค่าความคลาดเคลื่อนกำลังสอง (RMSE) หากค่าของการคำนวณมีค่าน้อยแสดงว่าแบบจำลองมีประสิทธิภาพมาก (Borchani et al., 2015)

3) ค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (Mean Absolute Error: MAE) เป็นการประเมินความแม่นยำของแบบจำลองการถดถอยโดยการวัดผลต่างสัมบูรณ์โดยเฉลี่ยระหว่างค่าที่คาดการณ์ไว้กับค่าเป้าหมาย ข้อดีของมาตรวัดนี้คือมีความทนทานต่อค่าผิดปกติมากกว่าวิธีอื่น โดยให้การแสดงข้อผิดพลาดที่สมดุล และตีความได้ง่ายเนื่องจากไม่สนใจทิศทางของความคลาดเคลื่อน (Ahmed, 2023) โดยมีสมการดังนี้

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

โดย n คือ จำนวนข้อมูลทดสอบ

y_i คือ ค่าจริงของข้อมูลลำดับที่ i เมื่อ $i = 1, 2, 3, \dots, n$

\hat{y}_i คือ ค่าทำนายของข้อมูลลำดับที่ i เมื่อ $i = 1, 2, 3, \dots, n$

4) ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (Mean Absolute Percentage Error: MAPE) ซึ่งใช้การรายงานผลเป็นร้อยละความคลาดเคลื่อนระหว่างค่าที่ทำนายกับค่าผลลัพธ์จริง ทำให้สามารถทำความเข้าใจถึงระดับความแม่นยำได้ แต่ยังมีข้อจำกัดอยู่คือตัวหารของสมการคือค่าจริง หากค่าจริงเป็น 0 จะไม่สามารถหาค่าได้ หรือหากใกล้เคียง 0 จะทำให้ค่าความคลาดเคลื่อนมีค่าผิดปกติ โดยมีสมการดังนี้

$$MAPE = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{|y_i|} \times 100$$

ค่าความคลาดเคลื่อนเฉลี่ยเปอร์เซ็นต์ (MAPE) ไม่สามารถคำนวณได้เมื่อค่าจริงเท่ากับศูนย์ เนื่องจากต้องมีการหารด้วยค่าจริง จึงได้มีการพัฒนาตัวชี้วัดประสิทธิภาพใหม่ที่เรียกว่า ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร

5) **ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (Symmetric Mean Absolute Percentage Error: SMAPE)** เพื่อแก้ไขปัญหาดังกล่าว SMAPE เป็นตัววัดความผิดพลาดที่ปรับปรุงจาก MAPE โดยเปลี่ยนตัวหารจากค่าจริงเพียงอย่างเดียว เป็นค่าเฉลี่ยระหว่างค่าจริงและค่าที่ทำนายในรูปแบบค่าสัมบูรณ์ ซึ่งช่วยลดผลกระทบของค่าจริงที่เป็นศูนย์หรือใกล้ศูนย์ในการคำนวณ (Hmong.in.th., n.d.)

$$SMAPE = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{\frac{|y_i| + |\hat{y}_i|}{2}} \times 100$$

y_i คือ ค่าจริงของข้อมูลลำดับที่ i เมื่อ $i = 1, 2, 3, \dots, n$

\hat{y}_i คือ ค่าทำนายของข้อมูลลำดับที่ i เมื่อ $i = 1, 2, 3, \dots, n$

2.1.8.2 แบบจำลองฐาน (Baseline Models)

ในการประเมินประสิทธิภาพของแบบจำลองการถดถอย (Regression Model) จำเป็นต้องมีจุดอ้างอิงที่เป็นมาตรฐานสำหรับใช้เปรียบเทียบผลลัพธ์และวัดศักยภาพของแบบจำลองที่พัฒนาขึ้น ดังนั้น การกำหนดแบบจำลองฐาน (Baseline Model) จึงมีบทบาทสำคัญในฐานะเกณฑ์ประเมินเบื้องต้นโดยทั่วไป Baseline Model มักมีโครงสร้างที่เรียบง่าย ไม่อาศัยกระบวนการเรียนรู้เชิงซับซ้อน (Learning Process) และไม่พยายามค้นหารูปแบบเชิงลึกจากข้อมูล แต่ใช้หลักการทางสถิติพื้นฐาน เช่น การทำนายค่ากลางจากข้อมูลฝึกสอน (Géron, 2019)

หนึ่งในวิธีที่นิยมสำหรับการสร้าง Baseline Model ในงานถดถอย คือ Dummy Regressor ซึ่งได้รับการออกแบบให้ทำนายค่าเป้าหมายโดยอิงจากค่าสถิติเบื้องต้นของข้อมูลฝึกสอน เช่น ค่าเฉลี่ย ค่ามัธยฐาน หรือค่าคงที่ที่กำหนดไว้ (Scikit-learn, 2024) โดย Dummy Regressor รองรับการทำนายที่หลากหลาย ได้แก่

Mean Strategy ทำนายค่าคงที่เท่ากับค่าเฉลี่ยของตัวแปรเป้าหมายในชุดฝึกสอน

Median Strategy ทำนายค่าคงที่เท่ากับค่ามัธยฐานของตัวแปรเป้าหมายในชุดฝึกสอน

Quantile Strategy ทำนายค่าคงที่ตามค่าเปอร์เซ็นต์ไทล์ที่ผู้ใช้กำหนด

Constant Strategy ทำนายค่าคงที่ที่ระบุโดยผู้ใช้งานเอง

สำหรับการประเมินประสิทธิภาพของ Dummy Regressor นั้น มักใช้ตัวชี้วัดเดียวกับแบบจำลองถดถอยอื่น ๆ เช่น Mean Absolute Error (MAE), Mean Squared Error (MSE) และ

Root Mean Squared Error (RMSE) ซึ่งการใช้ Baseline Model จะช่วยให้สามารถพิจารณาได้ว่าแบบจำลองที่พัฒนาขึ้นมีศักยภาพในการเรียนรู้จากข้อมูลและทำนายผลลัพธ์ได้ดีกว่าการทำนายเชิงพื้นฐานหรือไม่ หากแบบจำลองมีประสิทธิภาพต่ำกว่า Baseline Model บ่งชี้ถึงความจำเป็นในการปรับปรุงแบบจำลองหรือพิจารณาคุณภาพของข้อมูล แต่หากมีประสิทธิภาพสูงกว่าจึงจะถือว่าแบบจำลองดังกล่าวมีศักยภาพและเหมาะสมสำหรับการนำไปใช้งานจริง (Patel, n.d.)

2.1.9 Shapley Additive Explanation (SHAP)

Shapley Additive Explanation เป็นเครื่องมือแสดงภาพที่สามารถใช้ในการอธิบายการทำนายของแบบจำลองการเรียนรู้ของเครื่อง SHAP ใช้แนวคิดพื้นฐานจาก ค่า Shapley (Shapley Value) ที่มาจาก ทฤษฎีเกม (Game Theory) มีหลักการคือ ตัวแปรในแบบจำลองจะถูกมองว่าเป็น "ผู้เล่น" และการคำนวณค่า Shapley เพื่อค่าประเมินผลกระทบของตัวแปรแต่ละตัวในแบบจำลอง กล่าวคือ ผลกระทบจากการร่วมมือของผู้เล่นแต่ละคนเป็นเท่าไร โดยเฉลี่ย การเปลี่ยนแปลงของการคาดการณ์เมื่อมีตัวแปรเข้าร่วมชุดข้อมูลแต่ละชุด หากมีค่า Shapley Value สูงแสดงว่าทำให้ผลการทำนายมีความแม่นยำขึ้นมาก โดยมีสมการและแนวคิดของ Shapley Value ดังนี้

ขนาดของผลกระทบ ϕ_j ที่มีค่าของตัวแปรตัวที่ j ส่งผลให้การทำนายของจุดข้อมูลที่พิจารณาต่างจากค่าเฉลี่ยของทั้งชุดข้อมูล เช่น หากแบบจำลองที่ใช้เป็นแบบจำลองเชิงเส้น

$$\hat{f}(x) = \beta_0 + \beta_1 + \dots + \beta_p x_p$$

ค่าผลกระทบของ ϕ_j ที่ค่าของตัวแปรที่ j จะเขียนได้เป็น

$$\phi_j(\hat{f}) = \beta_j x_j - E(\beta_j X_j) = \beta_j x_j - \beta_j E(X_j)$$

โดย $\beta_j x_j$ คือการประเมินผลกระทบของตัวแปร x_j ณ จุดข้อมูลนั้น และ $E(\beta_j X_j)$ คือค่าเฉลี่ยของผลกระทบของตัวแปรนี้จากจุดข้อมูลทั้งหมด ซึ่งค่า ϕ_j คือค่าผลต่างของการทำนายเมื่อมีตัวแปรนั้นเทียบกับค่าเฉลี่ยของการทำนาย เมื่อนำค่าผลกระทบจากตัวแปรทุกตัวมารวมกัน จะได้ผลลัพธ์ของการทำนายเท่ากับผลต่างของค่าทำนายจุดนั้นกับค่าเฉลี่ยการทำนายของแบบจำลอง ดังสมการ

$$\begin{aligned} \sum_{j=1}^p \phi_j(\hat{f}) &= \sum_{j=1}^p (\beta_j x_j - E(\beta_j X_j)) \\ &= (\beta_0 + \sum_{j=1}^p \beta_j x_j) + (\beta_0 + \sum_{j=1}^p E(\beta_j X_j)) \\ &= \hat{f}(x) - E(\hat{f}(x)) \end{aligned}$$

(ปฏิภาณ ประเสริฐสม และ พีรต สามะศิริ, 2566)

2.2 งานวิจัยที่เกี่ยวข้อง

Poufinas et al.(2023) ได้ศึกษาการทำนายมูลค่าการเรียกร้อยค่าสินไหมทดแทนในประกันภัยรถยนต์ เนื่องจากมีผลต่อกระแสเงินสด การตั้งราคาเบี้ยประกันภัย และการจัดการความเสี่ยงของบริษัทประกันภัย โดยมีการนำเสนอตัวแปรที่ไม่ใหม่ในการเรียกร้อยค่าสินไหมทดแทน คือ สภาพอากาศและยอดขาย ทำนายด้วยแบบจำลองการเรียนรู้ของเครื่อง (Machine Learning) ดังนี้ 1. Support Vector Machines (SVM) 2. Decision Trees 3. Random Forests 4. Xgboost เพื่อทำนายมูลค่าการเรียกร้อยค่าสินไหมทดแทนเฉลี่ยต่อรถยนต์ที่ทำประกันในแต่ละไตรมาส ใช้ข้อมูลการเรียกร้อยค่าสินไหมทดแทนจากพอร์ตโฟลิโอประกันภัยรถยนต์ในกรุงเอเธนส์ ประเทศกรีซ ตั้งแต่ปี ค.ศ. 2008 ถึง ค.ศ. 2020 ผลการวิจัยแสดงให้เห็นว่าตัวแปรที่มีอิทธิพลมากที่สุด 3 ตัว ได้แก่ ยอดขายรถยนต์ใหม่ที่มีการชะลอข้อมูล 3 ไตรมาสและ 1 ไตรมาส และอุณหภูมิต่ำสุดของสถานีอากาศ Elefsina ที่ชะลอข้อมูล 3 ไตรมาส และผลการประเมินประสิทธิภาพของแบบจำลองการเรียนรู้ของเครื่อง พบว่าแบบจำลองที่ดีที่สุดคือ Random Forests และรองลงมาคือ XGBoost

Kumar et al. (2020) ได้ศึกษาการทำนายจำนวนมูลค่าการเรียกร้อยค่าสินไหมทดแทนจากการประกันภัยรถยนต์ในประเทศอินเดีย ซึ่งใช้ข้อมูลการเรียกร้อยค่าสินไหมทดแทนประกันภัยรถยนต์ต่อทรัพย์สิน ตั้งแต่ปี ค.ศ.1981 ถึงปี ค.ศ. 2016 สำหรับวิเคราะห์ โดยใช้แบบจำลองหลายประเภท คือ Generalized Linear Model (GLM) สำหรับการวิเคราะห์ข้อมูลที่มีการแจกแจงแบบไม่ปกติ ARIMA (Autoregressive Integrated Moving Average) ใช้ข้อมูลในอดีตเพื่อทำนายอนาคต และ Artificial Neural Network (ANN) ทำนายข้อมูลที่มีโครงสร้างไม่เชิงเส้น โดยมีความสามารถในการปรับตัวเองตามข้อมูลที่ได้รับ และได้ทำการเปรียบเทียบ ผลการทำนายกับข้อมูลจริงของมูลค่าการเรียกร้อยค่าสินไหมทดแทนในช่วง 36 ปี ผลลัพธ์การวิจัยพบว่าแบบจำลอง ANN มีประสิทธิภาพของการทำนายดีที่สุด โดยมีค่า RMSE เท่ากับ 0.17601 ในขณะที่แบบจำลอง GLM ได้ค่า RMSE เท่ากับ 1.183 และ ARIMA ได้ค่า RMSE เท่ากับ 1.3748 ดังนั้นแบบจำลอง ANN มีความแม่นยำสูงกว่าแบบจำลองอื่นๆในการทำนายมูลค่าการเรียกร้อยค่าสินไหมทดแทนประกันภัยรถยนต์ต่อทรัพย์สิน ซึ่งจะช่วยให้บริษัทประกันภัยสามารถคาดการณ์การเรียกร้อยค่าสินไหมทดแทนในอนาคตได้แม่นยำยิ่งขึ้น

Chen et al. (2020) ศึกษาการสร้างแบบจำลอง LightGBM (Light Gradient Boosting Machine) เพื่อทำนายความถี่ในการเรียกร้อยค่าสินไหมทดแทนในประกันภัยรถยนต์ ข้อมูลที่ใช้มาจากบริษัทประกันภัยรถยนต์ของสหรัฐอเมริกา มีทั้งหมด 10,305 แถว และ 24 ตัวแปร แบ่งออกเป็นสามกลุ่มได้แก่ ตัวแปรที่เกี่ยวกับเจ้าของรถ ตัวแปรที่เกี่ยวกับยานพาหนะ ตัวแปรที่เกี่ยวกับกรรมธรรม์

เพื่อใช้ในการทำนายความถี่ในการเรียกร้องค่าสินไหมทดแทนในประกันภัยรถยนต์ โดยมีการเปรียบเทียบวิธี LightGBM กับวิธีต่างๆ ได้แก่ 1. Gradient Decision Boosting Tree 2. Artificial Neural Network 3. Support Vector Machine 4. Generalized Linear Models ซึ่งพบว่า LightGBM สามารถจัดการกับข้อมูลที่ไม่สมดุลและมีความสามารถในการทำนายที่ดีที่สุดเมื่อเทียบกับแบบจำลองอื่น มีประสิทธิภาพดังนี้ Accuracy เท่ากับ 0.835 และค่า AUC เท่ากับ 0.907

Jin (2021) ศึกษาการทำนายความรุนแรงในการเรียกร้องค่าสินไหมทดแทนรายบุคคล โดยพิจารณามูลค่าการเรียกร้องค่าสินไหมทดแทนที่มีการรายงานแต่ยังไม่ได้รับการชำระ (Reported But Not Settled) สำหรับประกันวินาศภัยด้วยแบบจำลองการเรียนรู้ของเครื่องประเภท ensemble ML models ทั้งหมด 3 แบบจำลอง ได้แก่ 1. XGBoost 2. Random Forest 3. Extra Trees มีข้อมูลสำหรับการวิจัยทั้งหมด 500,914 ข้อมูล และ 32 ตัวแปร โดยใช้วิธีการ Grid Search เพื่อค้นหาค่าพารามิเตอร์ที่ดีที่สุดสำหรับแต่ละแบบจำลอง และนำมาทดสอบกับชุดข้อมูลทดสอบ ผลการวิจัยพบว่า XGBoost ให้ค่า RMSE ที่ต่ำที่สุดและมีเวลาในการคำนวณที่เร็วที่สุดเมื่อเปรียบเทียบกับวิธี ensemble ML models อื่น ๆ นอกจากนี้ XGBoost ยังช่วยให้สามารถวิเคราะห์และตีความปัจจัยที่มีผลต่อการทำนายมูลค่าการเรียกร้องค่าสินไหมทดแทน โดยใช้เทคนิค Tree SHAP พบว่าตัวแปรที่มีอิทธิพลต่อการทำนายมากที่สุดได้แก่ อายุของผู้เอาประกันภัย ไตรมาสที่เกิดอุบัติเหตุ จำนวนปีของการพัฒนาระหว่างชุดฝึกและชุดทดสอบ และส่วนของร่างกายที่ได้รับบาดเจ็บ ตามลำดับ

กิตติศักดิ์ และคณะ (2564) ได้ศึกษาการตรวจจบการเรียกร้องค่าสินไหมทดแทนประกันภัยรถยนต์ ซึ่งมีวัตถุประสงค์ในงานวิจัยคือ 1) หาปัจจัยที่มีอิทธิพลต่อการเรียกร้องค่าสินไหมทดแทน 2) เปรียบเทียบประสิทธิภาพการทำนายของวิธีการคัดเลือกตัวแปรอิสระ 3 วิธี คือ ไม่มีการคัดเลือกคุณลักษณะ วิธีการถดถอยลอจิสติกทีละขั้น และวิธีต้นไม้ตัดสินใจ 3) เปรียบเทียบประสิทธิภาพการทำนายของ 3 อัลกอริทึม คือ วิธีนาอูฟเบย์ วิธีสุ่มป่าไม้ และวิธีบูสต์ติงปรับได้ และ 4) เปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมร่วมกับวิธีการคัดเลือกคุณลักษณะ จากการวิจัยพบว่าคุณลักษณะที่มีอิทธิพลต่อการเรียกร้องค่าสินไหมทดแทนคือ ลักษณะการเกิดเหตุ ภูมิภาคที่เกิดเหตุ อายุผู้เอาประกันภัย ประเภทการซ่อม ทุนประกันภัย และจำนวนเงินที่จ่ายรวมทั้งการเรียกร้องค่าสินไหมทดแทน

Sun et al. (2024) ได้ศึกษาการประเมินความเสี่ยงการประกันภัยรถยนต์ซึ่งต้องการหาปัจจัยที่มีผลต่อประกันภัยรถยนต์และเปรียบเทียบประสิทธิภาพของแบบจำลองโดยใช้แบบจำลอง Actuarial Transformer (AT) ร่วมกับแบบจำลองแบบ Tree-Based คือ XGBoost, LightGBM,

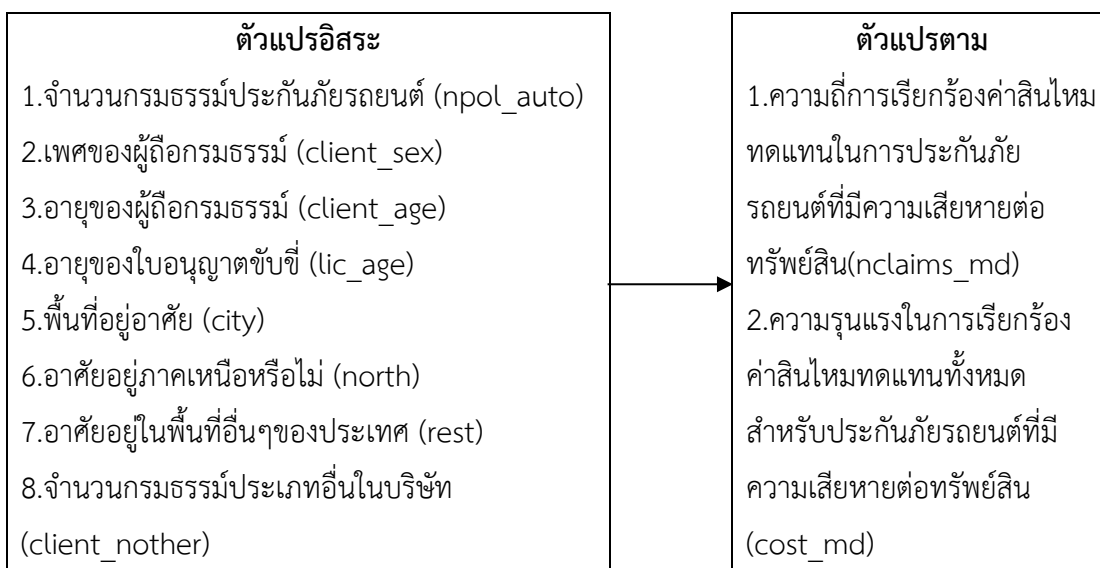
ตารางที่ 2.1 แบบจำลองและวิธีการที่ใช้ในการสร้างแบบจำลองจากการทบทวนวรรณกรรม

อ้างอิง	แบบจำลอง				วิธีการ	
	Random Forest	XGBoost	LightGBM	ANN	แบ่งข้อมูลอัตราส่วน 80:20	Hyperparameter Tuning
Poufinas et al. (2023)	✓	✓			✓	
Kumar et al. (2020)				✓		
Chen et al. (2020)			✓	✓		✓
Jin (2021)	✓	✓				✓
กิตติศักดิ์ และคณะ (2564)	✓				✓	
<div> <div>✓</div> <div>งานวิจัยที่เกี่ยวข้องของการศึกษาแบบจำลองหรือวิธีการดังกล่าว</div> </div> <div> <div>✓</div> <div>แบบจำลองที่ศึกษาและมีประสิทธิภาพดีที่สุดเมื่อเทียบกับแบบจำลองอื่นในงานวิจัยที่เกี่ยวข้อง</div> </div>						

ตารางที่ 2.2 ตัวแปรที่ส่งผลต่อการทำนายการเรียกร้องค่าสินไหมทดแทนจากการทบทวนวรรณกรรม

อ้างอิง	ตัวแปร							
	จำนวนกรมธรรม์ประกันภัยรถยนต์	เพศของผู้ถือกรมธรรม์	อายุของผู้ถือกรมธรรม์	อายุของใบอนุญาตขับขี่	พื้นที่อยู่อาศัย	อาศัยอยู่ภาคเหนือหรือไม่	อาศัยอยู่ในพื้นที่อื่นๆของประเทศ	จำนวนกรมธรรม์ประเภทอื่นในบริษัท
Chen et al. (2020)		✓	✓		✓			
Jin (2021)			✓					
กิตติศักดิ์ และอื่น ๆ (2564)		✓	✓			✓		
Sun (2024)			✓			✓		
✓ งานวิจัยที่เกี่ยวข้องของการศึกษาตัวแปรดังกล่าว ✓ ตัวแปรดังกล่าวมีผลต่อการทำนายหรือส่งผลต่อการเรียกร้องค่าสินไหมทดแทน								

2.3 กรอบแนวคิดการวิจัย



ภาพที่ 2.13 กรอบแนวคิดการวิจัย

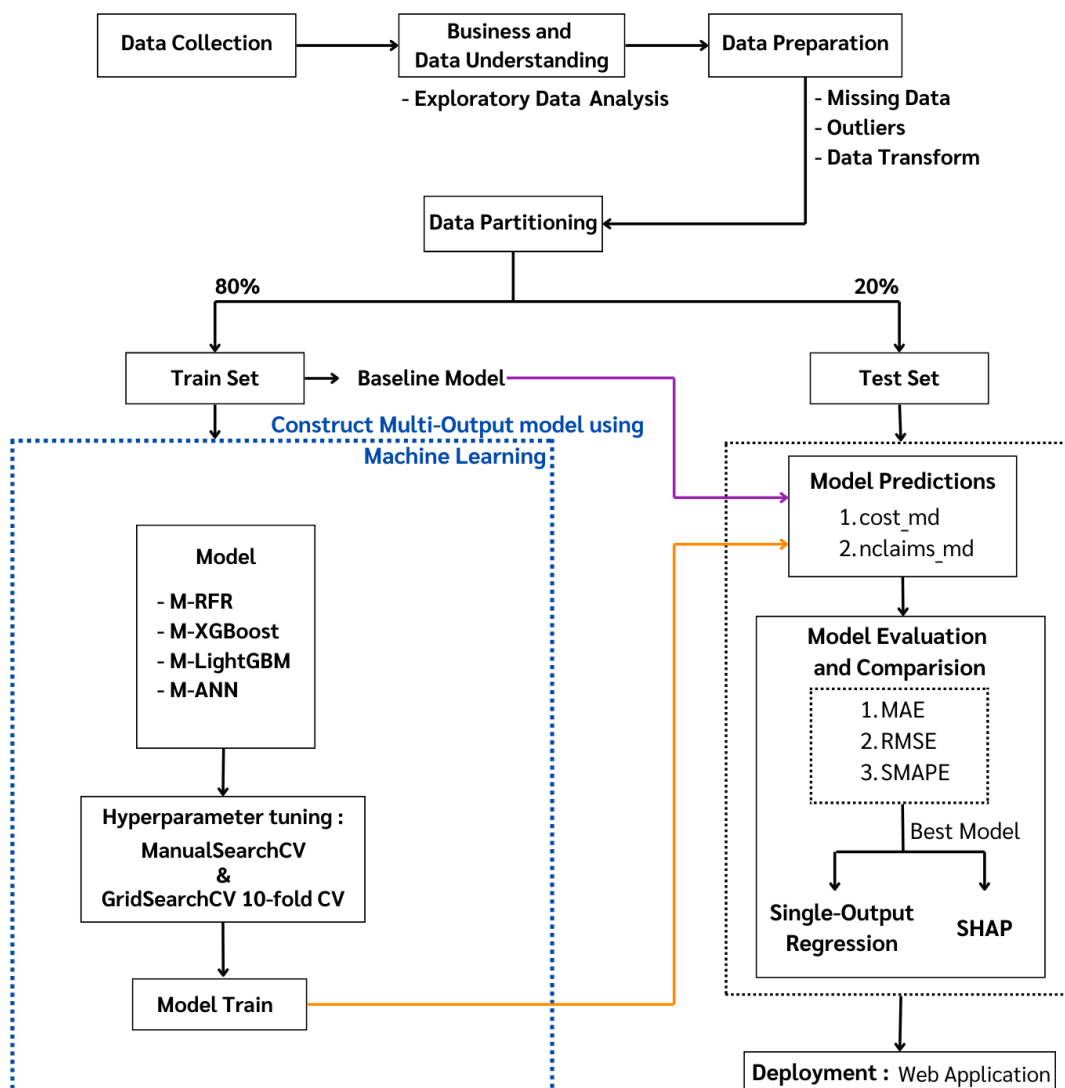
บทที่ 3

วิธีการดำเนินงานวิจัย

งานวิจัยเรื่อง การวิเคราะห์การถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ มีวัตถุประสงค์เพื่อสร้างและเปรียบเทียบประสิทธิภาพของแบบจำลองในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนและพัฒนาเว็บแอปพลิเคชันที่ช่วยในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน สำหรับการเตรียมและวิเคราะห์ข้อมูล ผู้วิจัยใช้ภาษา Python ผ่าน Google Colaboratory เวอร์ชัน 3.10.12 โดยมีขั้นตอนการดำเนินการวิจัย ดังนี้

- 3.1 การทำความเข้าใจธุรกิจ
- 3.2 การศึกษาและทำความเข้าใจข้อมูล
- 3.3 การเตรียมข้อมูล
 - 3.3.1 การทำความสะอาดข้อมูล
 - 3.3.2 การแปลงข้อมูล
- 3.4 การสร้างแบบจำลองการถดถอยพหุผลลัพธ์โดยใช้การเรียนรู้ของเครื่อง
 - 3.4.1 การแบ่งข้อมูล
 - 3.4.2 การสร้างแบบจำลองการเรียนรู้ของเครื่อง
- 3.5 การประเมินและเปรียบเทียบประสิทธิภาพแบบจำลอง
- 3.6 การนำแบบจำลองไปใช้งานจริง

โดยมีรายละเอียดขั้นตอนของการดำเนินงานวิจัย ดังภาพที่ 3.1



ภาพที่ 3.1 ขั้นตอนการดำเนินงานวิจัย

3.1 การทำความเข้าใจธุรกิจ (Business Understanding)

ในขั้นตอนแรก จะทำการศึกษาว่าธุรกิจประกันภัยรถยนต์คืออะไร มีความสำคัญต่อภาคเศรษฐกิจอย่างไร มีความเสี่ยงด้านใดบ้างที่ทำให้ธุรกิจประสบความล้มเหลว รวมถึงศึกษาว่ามีปัจจัยใดบ้างที่ส่งผลกระทบต่อความเสี่ยง

3.2 การศึกษาและทำความเข้าใจข้อมูล (Data Understanding)

ในขั้นตอนนี้จะทำความเข้าใจข้อมูล จัดทำ data dictionary และใช้ Exploratory Data Analysis: EDA ในการพรรณนาข้อมูล ได้แก่ การหาค่าเฉลี่ย ร้อยละ ส่วนเบี่ยงเบนมาตรฐาน ค่าสูงสุด ต่ำสุด การกระจายตัวของข้อมูล และ scatter plot ดูความสัมพันธ์ของข้อมูล

ข้อมูลที่ใช้ในการศึกษา คือข้อมูลการเรียกร่องค่าสินไหมทดแทนซึ่งเป็นข้อมูลทุติยภูมิ เป็นข้อมูลของบริษัทประกันภัยแห่งหนึ่งในประเทศสเปน ตั้งแต่วันที่ 1 มกราคม ค.ศ. 2006 ถึง 31 ธันวาคม ค.ศ. 2015 ระยะเวลารวม 10 ปี มีทั้งหมด 10 ตัวแปร จำนวน 80,924 แถว ได้มาจากงานวิจัยของ Catalina Bolance and Raluca Vernic ปี ค.ศ. 2019

3.3 การเตรียมข้อมูลก่อนการวิเคราะห์ (Data Preparation)

วัตถุประสงค์ของขั้นตอนนี้ คือการจัดการกับข้อมูลให้อยู่ในรูปแบบที่เหมาะสมก่อนการวิเคราะห์ โดยใช้โปรแกรม Google Collaboratory ในการเตรียมข้อมูล มีรายละเอียดการเตรียมข้อมูลดังนี้

3.3.1 การทำความสะอาดข้อมูล (Data Cleansing)

การทำความสะอาดข้อมูลเป็นการตรวจสอบข้อมูลและปรับปรุงข้อมูลให้มีคุณภาพมากขึ้น โดยหาจุดผิดของข้อมูลและหาวิธีการปรับปรุงข้อมูลนั้น ในงานวิจัยนี้ประกอบด้วย

1) การตรวจสอบความถูกต้องและความแนบเนียนของข้อมูล

เป็นการตรวจสอบข้อมูลว่า มีการบันทึกมาถูกต้องแนบเนียนหรือไม่ จะทำการตรวจสอบความแนบเนียนภายนอก (External consistency) ในแต่ละตัวแปร จะพิจารณาค่าที่เป็นไปได้ของข้อมูล หากพบว่าในแถวใด มีค่าของข้อมูลที่ไม่ถูกต้องแม้เพียงหนึ่งตัว จะทำการตัดข้อมูลทั้งแถวนั้นออก

2) การจัดการกับค่าสูญหาย (Missing Value)

ทำการตรวจสอบทุกแถว หากพบว่ามีข้อมูลสูญหายในแถวใด จะทำการตัดทั้งแถวนั้นออก

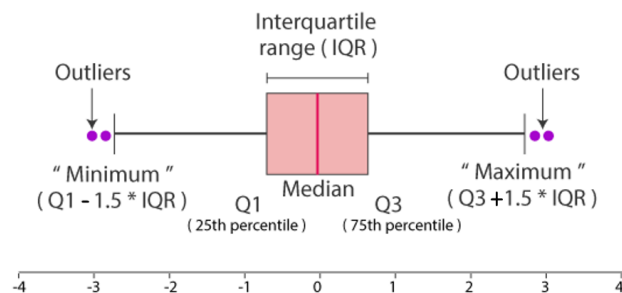
3) การจัดการกับค่านอกเกณฑ์ (Outliers)

โดยทั่วไปหากค่านอกเกณฑ์มีจำนวนน้อยกว่า 5% ของข้อมูลทั้งหมด จะคงค่านอกเกณฑ์ไว้เพื่อนำมาวิเคราะห์ต่อไป สำหรับข้อมูลเชิงคุณภาพที่มีค่าเป็น 0 และ 1 ไม่จำเป็นต้องจัดการกับค่านอกเกณฑ์ แต่สามารถใช้วิธีการถ่วงน้ำหนักเพื่อให้แบบจำลองรับรู้ถึงความไม่สมดุลของข้อมูล (Kutner et al., 2004) มีหลักการในการหาค่านอกเกณฑ์ดังนี้ ค่านอกเกณฑ์คือค่าที่มากกว่า $Q3 + 1.5 * IQR$ หรือ $Q1 - 1.5 * IQR$

เมื่อ Q1 คือ ค่าจากเปอร์เซ็นต์ไทล์ที่ 25 ของข้อมูล

Q3 คือ ค่าจากเปอร์เซ็นต์ไทล์ที่ 75 ของข้อมูล

IQR (Interquartile Range) คือ $Q3 - Q1$ แสดงรายละเอียดดังภาพที่ 3.3



ภาพที่ 3.2 ลักษณะของ Box plot

อย่างไรก็ตาม ข้อมูลทางการประกันภัยรถยนต์มีความเฉพาะเจาะจง โดยปกติในรอบปีของการประกันภัยรถยนต์ ผู้เอาประกันภัยส่วนใหญ่ไม่ได้มีการเรียกร้องค่าสินไหมทดแทน เมื่อเกิดการเรียกร้องค่าสินไหมทดแทนขึ้นค่าเหล่านี้จึงเป็นค่านอกเกณฑ์ และมูลค่าของการเรียกร้องค่าสินไหมทดแทนนั้นอาจสูงซึ่งผลต่อการดำเนินธุรกิจประกันภัย การตัดค่านอกเกณฑ์ออกอาจทำให้การวิเคราะห์ไม่สะท้อนความเป็นจริง นอกจากนี้ค่านอกเกณฑ์เป็นส่วนสำคัญของการวิเคราะห์ความเสี่ยงของบริษัทประกันภัย และบริษัทประกันภัยจำเป็นต้องใช้ข้อมูลที่ครอบคลุมทุกสถานการณ์เพื่อให้สามารถกำหนดเบี้ยประกันภัยที่เหมาะสม ดังนั้นผู้วิจัยจะทำการวิเคราะห์เบื้องต้นว่า ค่านอกเกณฑ์มีจำนวนร้อยละเท่าใด แต่จะไม่ทำการตัดค่านอกเกณฑ์ออก โดยจะนำข้อมูลทั้งหมดที่ผ่านการเตรียมข้อมูลแล้วไปสร้างแบบจำลอง

3.3.2 การแปลงข้อมูล (Data Transformation)

การแปลงข้อมูล เพื่อให้ข้อมูลพร้อมที่จะนำเข้าในการสร้างแบบจำลองทำนายการเรียกร้องค่าสินไหมทดแทน การแปลงข้อมูลนอกจากจะช่วยป้องกันปัญหาค่าทำนายที่ไม่เหมาะสมแล้ว ยังช่วยลดความแปรปรวนในข้อมูลและเพิ่มประสิทธิภาพการทำนายด้วยแบบจำลองการเรียนรู้ของเครื่อง ในงานวิจัยนี้ประกอบด้วย

- 1) แปลงข้อมูลเชิงคุณภาพให้อยู่ในรูปของตัวเลข
- 2) แปลงข้อมูลของตัวแปรตาม ได้แก่ ความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน (Y_1 และ Y_2 ตามลำดับ) เนื่องจากค่าตัวแปรตามทั้งสองไม่สามารถติดลบได้

ในการวิเคราะห์ข้อมูลที่เกี่ยวข้องกับตัวแปรตามหรือตัวแปรผลลัพธ์ด้วยแบบจำลองการเรียนรู้ของเครื่อง อาจเกิดปัญหาค่าทำนายของแบบจำลองที่มีค่าติดลบได้ ส่งผลให้ผลลัพธ์ไม่สอดคล้องกับบริบทตัวแปรผลลัพธ์ เพื่อแก้ไขปัญหานี้จึงมีการแปลงข้อมูลของตัวแปรโดยใช้ฟังก์ชันลอการิทึมธรรมชาติ (Logarithmic Transformation) ในการแปลงค่าตัวแปรผลลัพธ์ก่อนนำไปใช้ในการฝึกสอนแบบจำลอง (Model Training) มีการปรับค่า $\ln(y)$ โดยการเพิ่มค่าคงที่ 10^{-2}

เพื่อหลีกเลี่ยงปัญหาค่าศูนย์ $\ln(0)$ ที่ไม่สามารถหาค่าได้ การเพิ่มค่าคงที่ 10^{-2} ช่วยป้องกันปัญหาค่าทำนายที่เป็นลบ โดยสมการที่ใช้ในการแปลงข้อมูลมีดังนี้

$$y' = \ln(y + 10^{-2})$$

หลังจากการทำนายโดยแบบจำลองเสร็จสิ้นแล้ว ค่าทำนายจะถูกแปลงกลับมาเป็นค่าตัวแปรเดิม (y) เพื่อให้ผลลัพธ์สอดคล้องกับข้อมูลจริงและเหมาะสมต่อการนำไปใช้ในเชิงปฏิบัติ โดยสมการที่ใช้ในการแปลงกลับคือ

$$y = \exp(y') - 10^{-2}$$

3.4 การสร้างแบบจำลองการถดถอยพหุผลลัพธ์โดยใช้การเรียนรู้ของเครื่อง (Modeling)

ในงานวิจัยนี้ทำการการสร้างแบบจำลองการเรียนรู้ของเครื่องเพื่อทำนายจำนวน (ความถี่) และความรุนแรงในการเรียกร้องค่าสินไหมทดแทน จำนวน 4 แบบจำลอง ได้แก่

- 1) การถดถอยพหุผลลัพธ์แบบการสุ่มป่าไม้ (Multi-Output Random Forest Regression: M-RFR)
- 2) การถดถอยพหุผลลัพธ์แบบเอ็กซ์ตรีมเกรเดียนต์บูสติง (Multi-Output Extreme Gradient Boosting Regression: M-XGBoost)
- 3) การถดถอยพหุผลลัพธ์แบบไลต์เกรเดียนต์บูสติง (Multi-Output Light Gradient Boosting Regression: M-LightGBM)
- 4) การถดถอยพหุผลลัพธ์แบบโครงข่ายประสาทเทียม (Multi-Output Artificial Neural Network: M-ANN)

ก่อนที่จะทำการสร้างแบบจำลองจะทำการแบ่งข้อมูลที่ผ่านมาขั้นตอนการเตรียมข้อมูลมาแล้ว จากนั้นทำการปรับแต่งพารามิเตอร์เพื่อใช้ในการสร้างแบบจำลอง สำหรับการคัดเลือกตัวแปรอิสระที่จะนำมาใช้สร้างแบบจำลอง พบว่าข้อมูลที่ได้มีจำนวนตัวแปรอิสระไม่มาก กล่าวคือ มีจำนวน 8 ตัวแปร ดังนั้นผู้วิจัยจะนำตัวแปรอิสระทั้งหมดมาใช้ในการสร้างแบบจำลอง

3.4.1 การแบ่งข้อมูล (Data partitioning)

ข้อมูลที่ผ่านมาการเตรียมมาแล้วจะถูกแบ่งออกเป็นสองส่วน คือชุดข้อมูลฝึก (Train Set) คิดเป็น 80% จากข้อมูลทั้งหมด และชุดข้อมูลทดสอบ (Test Set) คิดเป็น 20% จากข้อมูลทั้งหมด หากข้อมูลของตัวแปรตามเป็นข้อมูลที่ไม่สมดุล (Imbalance data) จะทำการแบ่งข้อมูลด้วยวิธี stratify เพื่อให้สัดส่วนการแบ่งกลุ่มของข้อมูลทั้งสองชุดยังคงใกล้เคียงกัน

การฝึกแบบจำลอง งานวิจัยนี้ใช้วิธี K-fold Cross Validation โดยกำหนดค่า K เท่ากับ 10 หมายความว่าชุดข้อมูลฝึกจะถูกแบ่งเป็น 10 ส่วนย่อยที่มีขนาดเท่ากัน ในแต่ละรอบของการฝึกแบบจำลอง จะเลือกใช้ข้อมูล 9 ส่วนสำหรับฝึก และอีก 1 ส่วนที่เหลือใช้สำหรับทดสอบ แบบจำลองจะฝึกและทดสอบซ้ำกันทั้งหมด 10 รอบ โดยเปลี่ยนส่วนข้อมูลที่ใช้ทดสอบในแต่ละรอบ ทำให้ข้อมูลทั้งหมดถูกใช้ทั้งในการฝึกและทดสอบแบบสลับกัน กระบวนการนี้ช่วยเพิ่มความแม่นยำและความเสถียรของแบบจำลอง

3.4.2 การสร้างแบบจำลองการเรียนรู้ของเครื่อง (Model)

3.4.2.1 การสร้างแบบจำลองฐาน (Baseline Model)

แบบจำลองฐานใช้เป็นจุดอ้างอิงในการประเมินประสิทธิภาพของแบบจำลองการเรียนรู้ของเครื่อง โดยแบบจำลองการเรียนรู้ของเครื่องที่พัฒนาควรมีประสิทธิภาพการทำนายที่เหนือกว่าแบบจำลองฐาน เพื่อยืนยันความสามารถในการเรียนรู้และจับรูปแบบจากข้อมูลได้อย่างมีประสิทธิภาพกว่าการทำนายแบบง่าย โดยในการสร้างแบบจำลองฐาน จะใช้วิธี Dummy Regression ซึ่งดำเนินการตามขั้นตอนต่อไปนี้

1. เลือกรูปแบบการทำนาย (Strategy Selection) จากการพิจารณาลักษณะของข้อมูลในชุดฝึกสอน (Training Set) เพื่อเลือกกลยุทธ์ที่เหมาะสมสำหรับการทำนายค่าเป้าหมาย โดยอาจเลือกใช้

Mean Strategy: ทำนายค่าคงที่เป็นค่าเฉลี่ยของตัวแปรเป้าหมาย

Median Strategy: ทำนายค่าคงที่เป็นค่ามัธยฐานของตัวแปรเป้าหมาย

Quantile Strategy: ทำนายค่าคงที่ตามเปอร์เซ็นต์ไทล์ที่กำหนด

Constant Strategy: ทำนายค่าคงที่ที่กำหนดโดยผู้ใช้

2. สร้างแบบจำลอง Dummy Regressor โดยใช้ไลบรารี scikit-learn เพื่อสร้างแบบจำลองฐานตาม Strategy ที่เลือก โดยทำการฝึกแบบจำลองกับชุดฝึกสอนและสร้างค่าพยากรณ์ (predictor) สำหรับนำไปใช้ทดสอบกับชุดทดสอบ

3.4.2.2 การสร้างแบบจำลองการเรียนรู้ของเครื่อง (Machine Learning Model)

การสร้างแบบจำลองการเรียนรู้ของเครื่องเพื่อทำนายความถี่และความรุนแรงในการเรียกร้องค่าสินไหมทดแทน ใช้ข้อมูลชุดฝึกในการพัฒนาและปรับแต่งไฮเปอร์พารามิเตอร์ของแบบจำลอง ซึ่งจำเป็นต้องกำหนดค่าพารามิเตอร์ที่เหมาะสมสำหรับแต่ละแบบจำลอง ในงานวิจัยนี้มีวิธีการปรับแต่งค่าพารามิเตอร์ดังนี้

ขั้นที่ 1 ปรับแต่งไฮเปอร์พารามิเตอร์ด้วยวิธี Manual Search เพื่อสำรวจค่าพารามิเตอร์เบื้องต้น

ขั้นที่ 2 นำค่าไฮเปอร์พารามิเตอร์จากขั้นตอนที่ 1. ไปกำหนดขอบเขตสำหรับการปรับแต่งไฮเปอร์พารามิเตอร์เพิ่มเติมด้วยวิธี GridSearchCV เพื่อทดสอบทุกค่าที่เป็นไปได้ในช่วงพารามิเตอร์ที่กำหนด

ผู้วิจัยใช้ Google Colaboratory ในการสร้างแบบจำลอง มี library ที่เกี่ยวข้องกับงานวิจัยแสดงในตารางที่ 3.1

ตารางที่ 3.1 library ที่ใช้สำหรับวิเคราะห์ข้อมูล

library	Description
google.colab	ใช้สำหรับการเข้าถึงและทำงานร่วมกับ Google Drive
pandas	ใช้สำหรับการจัดการและวิเคราะห์ข้อมูลในรูปแบบ DataFrame เช่น CSV, Excel เป็นต้น
numpy	ใช้สำหรับการคำนวณเชิงตัวเลข เช่น คำนวณเกณฑ์
pyarrow.parquet	ใช้สำหรับการอ่านและเขียนไฟล์ Parquet ซึ่งเป็นไฟล์ที่มีการบีบอัดและจัดเก็บข้อมูลแบบ Column
scipy.stats	ใช้สำหรับการทดสอบทางสถิติต่าง ๆ เช่น Pearson correlation
sklearn.datasets	ใช้สำหรับการสร้างชุดข้อมูลตัวอย่างสำหรับการทดสอบหรือฝึกอบรมแบบจำลอง
sklearn.ensemble	ใช้สำหรับการสร้างแบบจำลองการเรียนรู้แบบรวม (ensemble learning) เช่น ExtraTreesRegressor
sklearn.multioutput	ใช้สำหรับการสร้างแบบจำลองการทำนายหลายผลลัพธ์
sklearn.model_selection	ใช้สำหรับการแบ่งชุดข้อมูล การทำ cross-validation (train_test_split) และการค้นหาพารามิเตอร์ที่ดีที่สุดด้วย GridSearchCV
sklearn.metrics	ใช้ในการวัดผลของแบบจำลอง เช่น MAE, RMSE
matplotlib.pyplot	ใช้สำหรับการสร้างกราฟและแผนภูมิเพื่อการแสดงผลข้อมูลที่เป็นภาพ
seaborn	ใช้สำหรับการสร้างกราฟที่มีความซับซ้อนและสวยงามมากขึ้น
time	ใช้สำหรับการจัดการและวัดเวลาในการประมวลผลของโค้ด
sklearn.neural_network	ใช้สำหรับการสร้างแบบจำลอง Artificial Neural Network (ANN)

library	Description
lightgbm	ใช้สำหรับการสร้างแบบจำลอง LightGBM
xgboost	สำหรับการสร้างแบบจำลอง Extreme Gradient Boosting
shap	สำหรับสร้างภาพเพื่อหาความสำคัญของตัวแปรอิสระในการทำนาย
flask	สำหรับเชื่อม backend ไปยัง frontend
dummy.Dummy	สำหรับสร้างแบบจำลองที่ใช้เป็น Baseline Model โดยไม่อาศัยการเรียนรู้จากข้อมูล แต่ใช้ค่าทางสถิติเช่น ค่าเฉลี่ย ค่ามัธยฐาน

ขั้นที่ 1 การปรับแต่งไฮเปอร์พารามิเตอร์ด้วยวิธี Manual Search มีขั้นตอน ดังนี้

1. ทดสอบด้วยค่าเริ่มต้น (default) ของแบบจำลอง เพื่อให้ทราบประสิทธิภาพเบื้องต้นของแบบจำลอง

2. เลือกพารามิเตอร์ที่ต้องการทดสอบ โดยทดสอบทีละพารามิเตอร์เพื่อตรวจสอบว่าการปรับค่าแต่ละค่าให้ประสิทธิภาพที่ดีขึ้นหรือไม่ เช่น Random Forest Regression ที่ทำการทดสอบ n_estimators ทดสอบที่ค่า 100, 200, 300, ...

3. ทดสอบแต่ละพารามิเตอร์ในแบบจำลอง จากนั้นพิจารณาประสิทธิภาพของแบบจำลองที่ได้ เมื่อพบว่าการปรับค่าพารามิเตอร์ไปในทิศทางใด ๆ ทำให้ประสิทธิภาพของแบบจำลองลดลงหรือไม่ดีขึ้น จะถือว่าเป็นจุดที่เหมาะสมในการหยุดการทดสอบและนำค่าไฮเปอร์พารามิเตอร์ไปกำหนดขอบเขตของค่าไฮเปอร์พารามิเตอร์ สำหรับทำ GridSearchCV

ค่าเริ่มต้นของพารามิเตอร์ในแต่ละแบบจำลอง แสดงดังตารางที่ 3.2 - 3.5 ตามลำดับ

ตารางที่ 3.2 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-RFR

ไฮเปอร์พารามิเตอร์	ค่าเริ่มต้น
n_estimators	100
max_depth	None
min_samples_split	2
min_sample_leaf	1

max_depth = None หมายถึงไม่มีการกำหนดความลึกสูงสุดของต้นไม้ โดยโหนดจะขยายออกจนกว่าโหนดสุดท้ายของต้นไม้จะไม่สามารถแบ่งได้อีก หรือจำนวนตัวอย่างในแต่ละใบจะน้อยกว่า min_samples_split

ตารางที่ 3.3 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-XGBoost

ไฮเปอร์พารามิเตอร์	ค่าเริ่มต้น
n_estimators	100
Eta	0.3
max_depth	6
min_child_weight	1
Subsample	1

ตารางที่ 3.4 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-LightGBM

ไฮเปอร์พารามิเตอร์	ค่าเริ่มต้น
learning_rate	0.1
n_estimators	100
max_depth	-1
Subsample	1
min_child_weight	0.001
num_leaves	31
min_data_in_leaf	20

Max_depth = -1 หมายความว่าไม่มีจำกัดความลึกของต้นไม้

ตารางที่ 3.5 ค่าเริ่มต้นของพารามิเตอร์สำหรับแบบจำลอง M-ANN

ไฮเปอร์พารามิเตอร์	ค่าเริ่มต้น
Hidden_layer_sizes	(100,)
Solver	adam
Activation	relu
Learning_rate	constant

Solver = adam หมายถึง แบบจำลองจะปรับค่าอัตราการเรียนรู้และน้ำหนักในโมเดลให้มีประสิทธิภาพสูงสุด Activation = relu หมายถึง ฟังก์ชันในโหนดซ่อน ให้ค่าเป็น 0 ถ้าอินพุตติดลบ และคงค่าเดิมถ้าอินพุตเป็นบวก

learning_rate = constant หมายถึง ใช้อัตราการเรียนรู้คงที่ตลอดการฝึกแบบจำลอง

ขั้นที่ 2 การปรับแต่งไฮเปอร์พารามิเตอร์เพิ่มเติมด้วยวิธี GridSearchCV

หลังจากได้ขอบเขตไฮเปอร์พารามิเตอร์สำหรับแต่ละแบบจำลองแล้ว จะทำการทดสอบทุกพารามิเตอร์ที่เป็นไปได้ในช่วงพารามิเตอร์ที่กำหนดด้วยวิธี GridSearchCV โดย Google Colaboratory มี library GridSearchCV ซึ่งมีค่าพารามิเตอร์ แสดงดังตารางที่ 3.6

ตารางที่ 3.6 พารามิเตอร์สำหรับ GridSearchCV

พารามิเตอร์	ความหมายและค่าที่กำหนด
estimator	แบบจำลองที่ต้องการ เช่น RandomForest, SVM ฯลฯ
param_grid	กำหนดค่าพารามิเตอร์ต่าง ๆ ที่ทดสอบจากแบบจำลอง เช่น แบบจำลอง Random Forest จะกำหนดพารามิเตอร์ดังนี้ 'estimator__n_estimators': [100, 200, 300], 'estimator__max_depth': [10, 20, 30] เป็นต้น
Scoring	ค่าที่ใช้ในการประเมินผลการทำงานของแบบจำลอง ในงานวิจัยนี้กำหนดเป็น MAE
n_jobs	กำหนดเป็น -1 หมายความว่าใช้ทุก CPU core ที่มีอยู่เพื่อเร่งกระบวนการคำนวณ
refit	กำหนดเป็น True แบบจำลองจะถูกฝึกซ้ำด้วยค่าพารามิเตอร์ที่ดีที่สุด
cv	Cross Validation ในงานวิจัยนี้จะใช้ K-fold Cross validation ซึ่งกำหนด K=10
pre_dispatch	จำนวนงานที่ถูกเตรียมไว้ล่วงหน้าก่อนที่การทำงานจะเริ่ม โดยกำหนดให้เท่ากับ 2*n_jobs (ค่าเริ่มต้น) นั่นคือเตรียมงานล่วงหน้าเป็นสองเท่าของจำนวนการทำงาน
error_score	ค่าที่จะใช้แทนค่าเมื่อแบบจำลองไม่สามารถฝึกได้หรือเกิดข้อผิดพลาดในการประเมิน กำหนดเป็น nan
return_train_score	คะแนนของการฝึกซ้อม (training scores) กำหนดเป็น False หมายความว่าคะแนนเฉพาะคะแนนในการทดสอบ (test scores)
verbose	ความละเอียดของข้อความที่แสดงเมื่อกำลังประมวลผล กำหนดค่าเป็น 0 หมายความว่าจะไม่มีการพิมพ์ข้อความ

การสร้างตัวแบบในงานวิจัยนี้จะคงค่าไฮเปอร์พารามิเตอร์ของ GridSearchCV ไว้ในทุก ๆ แบบจำลอง แต่จะมีการปรับเปลี่ยน estimator และ ค่า param_grid ของแต่ละแบบจำลอง เพื่อให้

มีความเหมาะสมกับแบบจำลอง โดยแต่ละแบบจำลองมีขอบเขตของไฮเปอร์พารามิเตอร์และหลังจากที่ได้พารามิเตอร์ที่ดีที่สุดสำหรับแต่ละแบบจำลองจากวิธี GridSearchCV แล้ว จะนำพารามิเตอร์เหล่านั้นมาทดสอบกับชุดข้อมูลทดสอบที่แบบจำลองไม่เคยพบมาก่อน เพื่อประเมินและเปรียบเทียบประสิทธิภาพของแบบจำลองว่าแบบจำลองใดมีประสิทธิภาพดีที่สุด

3.5 การประเมินและเปรียบเทียบประสิทธิภาพแบบจำลอง (Evaluation)

3.5.1 การประเมินและเปรียบเทียบประสิทธิภาพแบบจำลอง

การประเมินประสิทธิภาพของแบบจำลองในงานวิจัยนี้จะดำเนินการโดยนำแบบจำลองที่ผ่านการพัฒนาแล้ว มาทดสอบกับข้อมูลชุดทดสอบเพื่อวัดความสามารถในการทำนาย ซึ่งในงานวิจัยมีการทำนายสองผลลัพธ์ แต่ละผลลัพธ์จะถูกประเมินด้วยมาตรวัดดังนี้

1. ค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (Mean Absolute Error: MAE) โดยยิ่งค่า MAE ต่ำ คือเข้าใกล้ศูนย์ยิ่งแสดงถึงความแม่นยำในการทำนาย
2. ค่ารากที่สองของค่าคลาดเคลื่อนกำลังสองเฉลี่ย (Root Mean Squared Error: RMSE) โดยยิ่งค่า RMSE ต่ำ คือเข้าใกล้ศูนย์ยิ่งแสดงถึงความแม่นยำในการทำนาย
3. ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (Symmetric Mean Absolute Percentage Error: SMAPE) เป็นค่าที่ใช้วัดขนาดของความคลาดเคลื่อนของการทำนาย กับค่าจริงโดยไม่คำนึงถึงทิศทางและถูกแปลงให้อยู่ในรูปแบบของร้อยละ ยิ่งค่า SMAPE ต่ำ คือเข้าใกล้ศูนย์ยิ่งแสดงถึงความแม่นยำในการทำนาย

ในงานวิจัยนี้ การเปรียบเทียบประสิทธิภาพของแบบจำลองจะให้น้ำหนักกับค่า MAE มากที่สุดโดยพิจารณาแบบจำลองที่มีค่า MAE ต่ำที่สุดว่าเป็นแบบจำลองที่มีความสามารถในการทำนายมากที่สุด เนื่องจาก MAE เหมาะสมกับการพิจารณาสำหรับข้อมูลมีค่านอกเกณฑ์ (outliers) และจะนำค่าประสิทธิภาพของแต่ละแบบจำลองการเรียนรู้ของเครื่องมาเทียบกับแบบจำลองฐาน หากแบบจำลองการเรียนรู้ของเครื่องที่พัฒนามีประสิทธิภาพแย่กว่าแบบจำลองฐาน อาจจำเป็นต้องปรับปรุงแบบจำลองเพิ่มเติม เช่น การเลือกคุณลักษณะ (Feature Selection) หรือการปรับแต่งพารามิเตอร์ (Hyperparameter Tuning) อีกครั้ง เพื่อเพิ่มประสิทธิภาพของแบบจำลอง

3.5.2 การเปรียบเทียบประสิทธิภาพของแบบจำลองพหุผลลัพธ์ที่ดีที่สุดกับแบบจำลองผลลัพธ์เดียว

เพื่อประเมินว่าแบบจำลองประเภทการถดถอยพหุผลลัพธ์ที่ได้จากขั้นตอน 3.5.1 ว่ามีประสิทธิภาพเหนือกว่าแบบจำลองเดียวกันที่เป็นประเภทการถดถอยผลลัพธ์เดียว (Single-

Output Regression) หรือไม่ จะทำการสร้างแบบจำลองการถดถอยผลลัพธ์เดียว โดยทำนายเป็น แยกผลลัพธ์ ได้แก่ 1. ทำนายความถี่ของการเรียกร้องค่าสินไหมทดแทน และ 2. ทำนายความรุนแรงของการเรียกร้องค่าสินไหมทดแทน โดยจะทำการปรับแต่งพารามิเตอร์ให้เหมาะสมในการสร้างแบบจำลอง ใช้ข้อมูลชุดสร้างและทดสอบแบบจำลองเช่นเดียวกันกับแบบจำลองพหุผลลัพธ์ ทำการประเมินประสิทธิภาพด้วยค่า MAE RMSE และ SMAPE

3.5.3 Shapley Additive Explanation (SHAP)

หลังจากได้แบบจำลองพหุผลลัพธ์ที่ดีที่สุดในขั้นตอน 3.5.1 จะทำการวิเคราะห์ SHAP ซึ่งเป็นเครื่องมือที่ช่วยในการอธิบายผลของการทำนายจากแบบจำลอง เพื่อทำความเข้าใจคุณลักษณะที่มีผลกระทบต่อการตัดสินใจของแบบจำลองการทำนาย หรือเพื่อประเมินผลกระทบของแต่ละตัวแปรต่อผลลัพธ์การทำนาย ค่าของ Shapley value เป็นค่าที่บ่งบอกว่าตัวแปรแต่ละตัวส่งผลต่อการตัดสินใจของแบบจำลองในลักษณะใด ถ้าค่าเป็นบวกหมายความว่าตัวแปรนั้นช่วยผลักดันค่าการทำนายให้สูงขึ้น และถ้าค่าเป็นลบหมายความว่าตัวแปรนั้นทำให้ค่าการทำนายลดลง ตัวแปรที่มีค่า Shapley Value สูงสุดจะถูกพิจารณาว่าเป็นปัจจัยสำคัญที่สุดในการทำนายผลลัพธ์

ผู้วิจัยจะนำเสนอแบบรูปภาพด้วย SHAP summary plot เป็นกราฟที่แสดงค่าผลกระทบของตัวแปรต่าง ๆ ในแบบจำลอง โดยจุดสีแดงจะแสดงผลในเชิงบวก (ค่าการทำนายเพิ่ม) และจุดสีฟ้าในเชิงลบ (ค่าการทำนายลด) การนำเสนอผลลัพธ์ในรูปแบบนี้ช่วยเพิ่มความโปร่งใสและความน่าเชื่อถือของแบบจำลอง ทำให้สามารถตรวจสอบความสมเหตุสมผลของผลการทำนายได้อย่างชัดเจนและเข้าใจง่าย อีกทั้งยังสนับสนุนการตัดสินใจที่แม่นยำมากขึ้นในการนำแบบจำลองไปประยุกต์ใช้ในสถานการณ์จริง

3.6 การนำแบบจำลองไปใช้งานจริง (Deployment)

การนำแบบจำลองไปใช้งานจริง ผู้วิจัยจะนำแบบจำลองที่มีประสิทธิภาพดีที่สุดจากการทดสอบมาพัฒนาเป็นเว็บแอปพลิเคชัน เพื่อให้บริษัทประกันภัยรถยนต์สามารถนำไปใช้งานทำนายผลลัพธ์ได้สะดวกและรวดเร็ว โดยใช้โปรแกรม Visual Studio Code (VS Code) และใช้ Flask ซึ่งเป็น Web API เชื่อมต่อแบบจำลอง Machine Learning เข้ากับระบบเว็บ โดยหน้าต่างของเว็บแอปพลิเคชัน (User Interface) พัฒนาโดยใช้ภาษา HTML, CSS, JavaScript มีขั้นตอนการทำงานของเว็บแอปพลิเคชัน ดังนี้

1. ผู้ใช้ป้อนข้อมูลผ่านเว็บแอปพลิเคชัน ดังนี้
 - 1) จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี

- 2) เพศของผู้ถือกรรมธรรม์
- 3) อายุของผู้ถือกรรมธรรม์
- 4) อายุของใบอนุญาตขับขี่
- 5) พื้นที่อยู่อาศัย
- 6) การอาศัยอยู่ภาคเหนือหรือไม่
- 7) อาศัยอยู่พื้นที่อื่น ๆ ของประเทศ
- 8) จำนวนกรรมธรรม์ประเภทอื่นในบริษัทเดียวกัน

2. ประมวลผลข้อมูลและทำนายผลลัพธ์ เว็บแอปพลิเคชันใช้ Flask เป็นตัวกลางในการรับข้อมูลที่ใช้ป้อนเข้ามา จากนั้นทำการประมวลผลข้อมูล เช่น การแปลงให้อยู่ในรูปแบบที่สามารถนำไปใช้งานกับแบบจำลองได้ ข้อมูลที่ผ่านการประมวลผลแล้วจะถูกส่งเข้าสู่แบบจำลองเพื่อทำนายผลลัพธ์ และ Flask จะรับค่าทำนายกลับมาเพื่อนำไปแสดงผลแก่ผู้ใช้

3. แสดงผลลัพธ์ ผลลัพธ์จากการทำนายจะถูกส่งกลับไปยังผู้ใช้ในรูปแบบตัวเลข ดังนี้

- 1) ความรุนแรงในการเรียกร้องค่าสินไหมทดแทนทั้งหมดสำหรับประกันภัยรถยนต์
- 2) จำนวนการเรียกร้องค่าสินไหมทดแทนสำหรับประกันภัยรถยนต์

สำหรับการนำเว็บแอปพลิเคชันไปใช้งานจริง ผู้วิจัยจะเผยแพร่โค้ดผ่าน GitHub เพื่อให้ผู้ใช้ที่สนใจสามารถดาวน์โหลดและติดตั้งใช้งาน โดยสามารถปรับแต่งหรือพัฒนาเพิ่มเติมให้เหมาะสมกับความต้องการ

บทที่ 4

ผลการวิจัย

ผลกาวิจัยเรื่องการวิเคราะห์การถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ มีวัตถุประสงค์เพื่อสร้างและเปรียบเทียบประสิทธิภาพของแบบจำลองในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนและเพื่อพัฒนาเว็บแอปพลิเคชันที่ช่วยในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน ซึ่งมีผลการวิจัยดังต่อไปนี้

- 4.1 ผลการศึกษาธุรกิจประกันภัยรถยนต์
- 4.2 ผลการศึกษาข้อมูลการเรียกร้องค่าสินไหมทดแทน
 - 4.2.1 รายละเอียดของข้อมูล
 - 4.2.2 ผลการวิเคราะห์ข้อมูลเชิงสำรวจ
- 4.3 ผลการเตรียมข้อมูลก่อนการวิเคราะห์
 - 4.3.1 ผลการทำความสะอาดข้อมูล
 - 4.3.2 ผลการแปลงข้อมูล
- 4.4 ผลการสร้างแบบจำลอง
 - 4.4.1 ผลการแบ่งข้อมูล
 - 4.4.2 ผลการสร้างแบบจำลองฐาน
 - 4.4.3 ผลการปรับแต่งพารามิเตอร์ของแต่ละแบบจำลอง
- 4.5 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลอง
 - 4.5.1 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการถดถอยพหุผลลัพธ์
 - 4.5.2 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการถดถอยพหุผลลัพธ์กับการถดถอยผลลัพธ์เดียว
 - 4.5.3 ผลการอธิบายแบบจำลองด้วยเทคนิค SHAP
- 4.6 ผลการพัฒนาเว็บแอปพลิเคชัน

4.1 ผลการศึกษาธุรกิจประกันภัยรถยนต์

ธุรกิจประกันภัยรถยนต์เป็นส่วนหนึ่งของประกันภัยเบ็ดเตล็ด (Casualty Insurance) ที่มีบทบาทสำคัญต่อเศรษฐกิจโลก ปัจจุบัน ตลาดประกันภัยมีมูลค่าประมาณ 5.8 ล้านล้านดอลลาร์สหรัฐ ขณะที่การลงทุนในเทคโนโลยีด้านประกันภัย (InsurTech) มีมูลค่าราว 7.2 พันล้านดอลลาร์สหรัฐ นอกจากนี้ ประกันภัยรถยนต์ยังเป็นหนึ่งในประเภทประกันภัยที่ได้รับความนิยมสูงสุดในอุตสาหกรรม

แนวโน้มของอุตสาหกรรมประกันภัยรถยนต์ คาดการณ์ว่าจะมีการเปลี่ยนแปลงในพฤติกรรม การซื้อและขายกรมธรรม์ โดย Priceza Money ได้สรุป 3 แนวโน้มที่สำคัญในอุตสาหกรรมประกันภัยรถยนต์ ดังนี้

1. Personalize insurance คือ รูปแบบประกันภัยที่สามารถปรับแต่งความคุ้มครองให้เหมาะสมกับความต้องการและพฤติกรรมการใช้รถของแต่ละบุคคล เช่น ประกันตามระยะทางการขับขี่ คือ คิดค่าเบี้ยประกันตามระยะทางที่ขับขี่จริง ประกันแบบเปิด-ปิด คือ ผู้ขับขี่สามารถเปิดหรือปิดการคุ้มครองได้ตามต้องการ และจ่ายเบี้ยประกันตามระยะเวลาการใช้งาน และประกันระยะสั้น เหมาะสำหรับผู้ที่ต้องการความคุ้มครองในช่วงเวลาจำกัดโดยไม่ต้องจ่ายเบี้ยประกันรายปี
2. Direct-to-Customer (D2C) Insurance คือ การขายกรมธรรม์โดยตรงจากบริษัทประกันภัยผ่านแพลตฟอร์มออนไลน์ เช่น เว็บไซต์หรือแอปพลิเคชันของบริษัท ช่วยให้ลูกค้าสามารถเลือกซื้อประกันได้สะดวกและรวดเร็วขึ้น โดยแนวโน้มนี้มีโอกาสเติบโตสูง โดยเฉพาะในกลุ่มผู้บริโภคที่ต้องการกรมธรรม์ที่ตอบโจทย์ความต้องการเฉพาะบุคคล
3. ผู้บริโภคให้ความสำคัญกับ ความน่าเชื่อถือและการให้บริการหลังการขาย มากกว่าการเลือกซื้อประกันจากราคาที่ถูกที่สุด ส่งผลให้บริษัทประกันภัยต้องปรับกลยุทธ์ โดยมุ่งเน้นการให้บริการที่ยืดหยุ่น สะดวก และสามารถตอบโจทย์ความต้องการเฉพาะบุคคลได้ดียิ่งขึ้น (Marketingoops, 2022)

ความก้าวหน้าทางเทคโนโลยีทำให้การประยุกต์ใช้การเรียนรู้ของเครื่องในภาคประกันภัยมีบทบาทสำคัญในการเพิ่มศักยภาพของบริษัทประกันภัย โดยช่วยให้สามารถวิเคราะห์ข้อมูลจำนวนมากได้อย่างมีประสิทธิภาพ การพัฒนาแบบจำลองคาดการณ์ที่ขับเคลื่อนด้วยการเรียนรู้ของเครื่องช่วยเพิ่มความแม่นยำในการคาดการณ์พฤติกรรมของลูกค้าและการประเมินความเสี่ยง ตลอดจนสนับสนุนมาตรการเชิงรุก เช่น การพยากรณ์ความต้องการของลูกค้า การให้คำแนะนำเฉพาะ

บุคคล และการปรับปรุงกระบวนการรับประกันภัย ซึ่งส่งผลให้การกำหนดราคาและการบริหารต้นทุนค่าสินไหมทดแทนมีประสิทธิภาพมากขึ้น (Arthur D. Little, n.d.)

ในปัจจุบันความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนจากผู้เอาประกันภัยมีความผันผวนสูง ส่งผลกระทบโดยตรงต่อการกำหนดเบี้ยประกันภัย ความมั่นคงของธุรกิจ และการรักษาความยั่งยืนในระยะยาว บริษัทประกันภัยจึงจำเป็นต้องมีการคาดการณ์ที่แม่นยำเกี่ยวกับความถี่ (Claim Frequency) และความรุนแรง (Claim Severity) ของการเรียกร้องค่าสินไหมทดแทนเพื่อลดความเสี่ยงในการตั้งราคาที่ไม่เหมาะสม อาจทำให้เกิดการสูญเสียรายได้หรือความไม่พึงพอใจของลูกค้า

ด้วยเหตุนี้งานวิจัยนี้จึงมีวัตถุประสงค์ที่จะพัฒนาแบบจำลองการถดถอยพหุผลลัพธ์ (Multi-output Regression) สำหรับการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนพร้อมกันโดยใช้เทคนิคการเรียนรู้ของเครื่อง เพื่อให้บริษัทประกันภัยสามารถคาดการณ์แนวโน้มของการเรียกร้องค่าสินไหมทดแทนได้แม่นยำ และสามารถนำไปกำหนดเบี้ยประกันภัยได้อย่างเหมาะสม ช่วยรักษาความมั่นคงของบริษัทประกัน ซึ่งมีปัจจัยหลักมาจากลักษณะเฉพาะของผู้เอาประกันภัย เช่น อายุของผู้ถือกรมธรรม์ ซึ่งเป็นตัวแปรที่ได้รับความสนใจและมีการศึกษาในหลายประเทศ (Chen et al., 2020) นอกจากนี้ การเลือกแบบจำลองที่เหมาะสมมีผลโดยตรงต่อความแม่นยำของทำนาย เนื่องจากการเลือกแบบจำลองที่ไม่เหมาะสมอาจทำให้ค่าทำนายคลาดเคลื่อนหรือมีความแปรปรวนต่ำเกินจริง (Selection Effect) ดังนั้น การใช้แบบจำลองแบบรวมตัวแปรทั้งหมดที่มีอยู่ในข้อมูลโดยไม่ทำการคัดเลือกตัวแปร อาจให้ผลลัพธ์ที่น่าเชื่อถือมากกว่า (Hong, Kuffner, & Martin, 2018)

4.2 ผลการศึกษาข้อมูลการเรียกร้องค่าสินไหมทดแทน

4.2.1 รายละเอียดของข้อมูล

จากการศึกษาข้อมูลการเรียกร้องค่าสินไหมทดแทนของบริษัทประกันภัยแห่งหนึ่งในประเทศสเปน พบว่า มีทั้งหมด 10 ตัวแปร จำนวน 80,924 แถว มีข้อมูลเกี่ยวกับลักษณะของผู้เอาประกันภัย ประกอบด้วย 1) ความรุนแรงในการเรียกร้องค่าสินไหมทดแทนทั้งหมดสำหรับประกันภัยรถยนต์ 2) จำนวนการเรียกร้องค่าสินไหมทดแทนสำหรับประกันภัยรถยนต์ 3) จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี 4) เพศของผู้ถือกรมธรรม์ 5) อายุของผู้ถือกรมธรรม์ 6) อายุของใบอนุญาตขับขี่ 7) พื้นที่อยู่อาศัย 8) การอาศัยอยู่ภาคเหนือหรือไม่ 9) อาศัยอยู่พื้นที่อื่นๆ ของประเทศ 10) จำนวนกรมธรรม์ประเภทอื่นในบริษัทเดียวกัน สำหรับประกันภัยรถยนต์ มีรายละเอียดข้อมูลแสดงดังตารางที่ 4.1

ตารางที่ 4.1 รายละเอียดข้อมูลการประกันภัยรถยนต์

ตัวแปร	ชื่อตัวแปร	ความหมายของตัวแปร	ประเภทของข้อมูล (มาตรวัด)	ช่วงของข้อมูล/ ค่าที่เป็นไปได้
Y_1	nclaims_md	ความถี่ในการเรียกร้องค่าสินไหมทดแทนประกันภัยรถยนต์ที่มีความเสียหายต่อทรัพย์สิน	ปริมาณ (อัตราส่วน)	0 – 40
Y_2	cost_md	ความรุนแรงในการเรียกร้องค่าสินไหมทดแทนทั้งหมดสำหรับประกันภัยรถยนต์ที่มีความเสียหายต่อทรัพย์สิน (หน่วย: ยูโร)	ปริมาณ (อัตราส่วน)	0 – 65875.22
X_1	npol_auto	จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี	ปริมาณ (อัตราส่วน)	1 – 35
X_2	client_sex	เพศของผู้ถือกรมธรรม์	คุณภาพ (นามบัญญัติ)	Man Woman
X_3	client_age	อายุของผู้ถือกรมธรรม์	ปริมาณ (อันตรภาค)	20 – 90
X_4	lic_age	อายุของใบอนุญาตขับขี่	ปริมาณ (อันตรภาค)	1 – 70
X_5	city	พื้นที่อยู่อาศัย	คุณภาพ (นามบัญญัติ)	0 = อื่นๆ 1 = อยู่ในเมืองใหญ่
X_6	north	อาศัยอยู่ภาคเหนือหรือไม่	คุณภาพ (นามบัญญัติ)	0 = ไม่ 1 = ใช่
X_7	rest	อาศัยอยู่พื้นที่อื่นๆของประเทศ	คุณภาพ (นามบัญญัติ)	0 = ไม่ใช่ 1 = ใช่
X_8	client_nother	จำนวนกรมธรรม์ประเภทอื่นในบริษัทเดียวกัน	ปริมาณ (อัตราส่วน)	0 – 23

ในกรณีตัวแปร lic_age (X_4) คือ อายุของใบอนุญาตขับขี่ เป็นระยะเวลาที่ผู้ขับขี่ได้รับใบอนุญาตขับขี่ โดยคำนวณจากปีที่ได้รับใบอนุญาตครั้งแรกจนถึงปีปัจจุบัน เช่น ถ้าลูกค้าได้รับใบอนุญาตขับขี่ตอนอายุ 20 ปี และปัจจุบันลูกค้าอายุ 30 ปี " lic_age " จะเท่ากับ 10 ปี

$rest$ (X_7) คืออาศัยอยู่พื้นที่อื่นๆของประเทศ มีค่าที่เป็นไปได้สองค่า ได้แก่ 1 = ใช่ ($X_5 = 0$ และ $X_6 = 0$) และมีค่าเป็น 0 = ไม่ใช่ (อยู่ในพื้นที่ใดพื้นที่หนึ่ง ยกเว้น $X_5 = 1$ และ $X_6 = 1$)

$client_nother$ (X_8) คือ จำนวนกรรมธรรม์ประเภทอื่นในบริษัท ที่ไม่ใช่ประกันภัยรถยนต์และประกันภัยบ้าน (เช่น ประกันภัยอุบัติเหตุ ประกันชีวิตแบบตลอดชีพ ประกันชีวิตแบบบำนาญ)

4.2.2 ผลการวิเคราะห์ข้อมูลเชิงสำรวจ

ผลการวิเคราะห์สถิติเบื้องต้นของข้อมูลผู้ถือกรรมธรรม์ แสดงดังภาพที่ 4.1

1) ผลการวิเคราะห์ข้อมูลเชิงสำรวจ

	nclaims_md	cost_md	npol_auto	client_sex	client_age	lic_age	client_nother	city	north	rest
count	80924.000000	80924.000000	80924.000000	80924	80924.000000	80924.000000	80924.000000	80924.0	80924.0	80924.0
unique	NaN	NaN	NaN	2	NaN	NaN	NaN	2.0	2.0	2.0
top	NaN	NaN	NaN	Man	NaN	NaN	NaN	0.0	0.0	1.0
freq	NaN	NaN	NaN	61731	NaN	NaN	NaN	64994.0	57566.0	41636.0
mean	0.255511	234.202818	2.44679	NaN	53.24177	29.937682	0.219020	NaN	NaN	NaN
std	0.822942	916.829434	2.10436	NaN	13.11958	11.359151	0.655378	NaN	NaN	NaN
min	0.000000	0.000000	1.00000	NaN	18.00000	1.000000	0.000000	NaN	NaN	NaN
25%	0.000000	0.000000	1.00000	NaN	43.00000	21.000000	0.000000	NaN	NaN	NaN
50%	0.000000	0.000000	2.00000	NaN	53.00000	30.000000	0.000000	NaN	NaN	NaN
75%	0.000000	0.000000	3.00000	NaN	63.00000	38.000000	0.000000	NaN	NaN	NaN
max	40.000000	65875.220000	35.00000	NaN	90.00000	70.000000	23.000000	NaN	NaN	NaN

ภาพที่ 4.1 ผลการวิเคราะห์สถิติเบื้องต้นของข้อมูล

จากภาพที่ 4.1 ข้อมูลผู้ถือกรรมธรรม์ประกอบด้วยตัวแปรเชิงปริมาณและเชิงคุณภาพ สำหรับตัวแปรเชิงปริมาณจะปรากฏค่า count mean std min 25% 50% 75% max และสำหรับตัวแปรเชิงคุณภาพจะปรากฏค่า count unique top freq ในกรณีที่เห็นค่า NaN หมายความว่าไม่สามารถหาสำหรับตัวแปรประเภทนั้นได้ โดยพบว่า ทุกตัวแปรมีจำนวนข้อมูลทั้งหมด 80,924 แถว

จำนวนการเรียกร้องค่าสินไหมทดแทน ($nclaims_md$) มีค่าเฉลี่ยเท่ากับ 0.2555 ส่วนเบี่ยงเบนมาตรฐานเท่ากับ 0.8229 และค่ามัธยฐานเท่ากับ 0 ซึ่งพบว่าข้อมูลส่วนใหญ่ไม่มีการเรียกร้องค่าสินไหมทดแทน และค่าสูงสุดของข้อมูลเท่ากับ 40 แสดงว่ามีค่าห่างจากข้อมูลส่วนใหญ่หรือค่านอกเกณฑ์อยู่

ความรุนแรงในการเรียกร้องค่าสินไหมทดแทน (cost_md) มีค่าเฉลี่ยเท่ากับ 234.2028 ยูโร ส่วนเบี่ยงเบนมาตรฐานเท่ากับ 916.8294 ค่าต่ำสุดของข้อมูลเท่ากับ 0 ค่าที่อยู่ตำแหน่ง 25% ของข้อมูลเท่ากับ 0 ค่ามัธยฐานเท่ากับ 0 ค่าที่อยู่ตำแหน่ง 75% ของข้อมูลเท่ากับ 0 และค่าสูงสุดของข้อมูลเท่ากับ 65875.2200 แสดงว่ามีค่าห่างจากข้อมูลส่วนใหญ่หรือค่านอกเกณฑ์อยู่

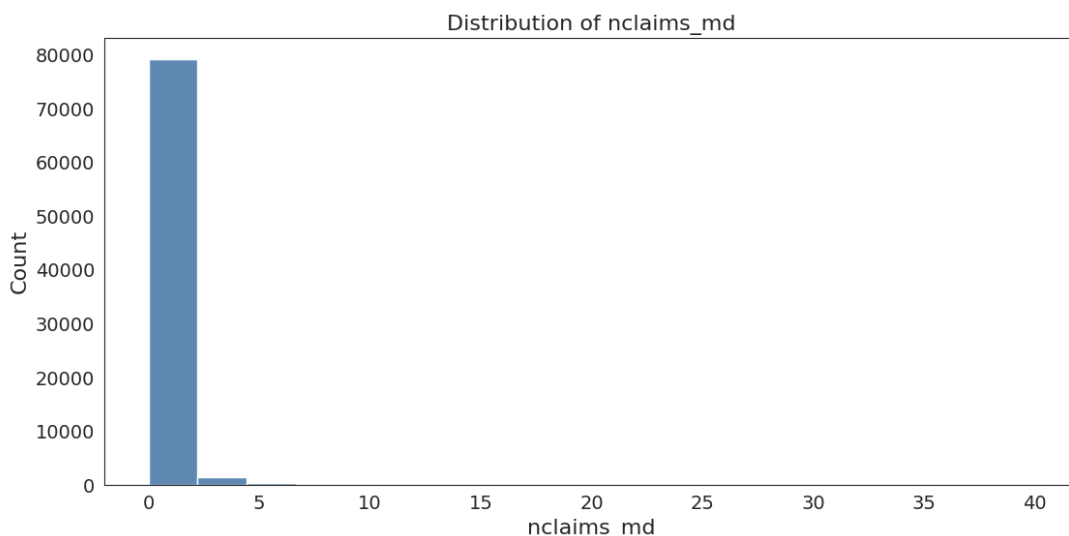
จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี (npol_auto) มีค่าเฉลี่ยเท่ากับ 2.4467 ส่วนเบี่ยงเบนมาตรฐานเท่ากับ 2.1043 ค่าต่ำสุดของข้อมูลเท่ากับ 1 ค่าที่อยู่ตำแหน่ง 25% ของข้อมูลเท่ากับ 1 ค่ามัธยฐานข้อมูลเท่ากับ 2 ค่าที่อยู่ตำแหน่ง 75% ของข้อมูลเท่ากับ 3 และค่าสูงสุดของข้อมูลเท่ากับ 35 ซึ่งมีค่าห่างจากค่ากลางค่อนข้างมาก

อายุของผู้ถือกรมธรรม์ (client_age) มีค่าเฉลี่ยเท่ากับ 53.2417 ส่วนเบี่ยงเบนมาตรฐานเท่ากับ 13.1195 นั่นคือผู้ถือกรมธรรม์ส่วนใหญ่อยู่ในช่วงวัยกลางคนถึงผู้สูงอายุ ค่าต่ำสุดของข้อมูลเท่ากับ 18 ค่าที่อยู่ตำแหน่ง 25% ของข้อมูลเท่ากับ 43 ค่ามัธยฐานเท่ากับ 53 ค่าที่อยู่ตำแหน่ง 75% ของข้อมูลเท่ากับ 63 และค่าสูงสุดของข้อมูลเท่ากับ 90

อายุของใบอนุญาตขับขี่ (lic_age) มีค่าเฉลี่ยเท่ากับ 29.9376 ส่วนเบี่ยงเบนมาตรฐานเท่ากับ 11.3591 ค่าต่ำสุดของข้อมูลเท่ากับ 1 ค่าที่อยู่ตำแหน่ง 25% ของข้อมูลเท่ากับ 21 ค่ามัธยฐานเท่ากับ 30 ค่าที่อยู่ตำแหน่ง 75% ของข้อมูลเท่ากับ 38 และค่าสูงสุดของข้อมูลเท่ากับ 70

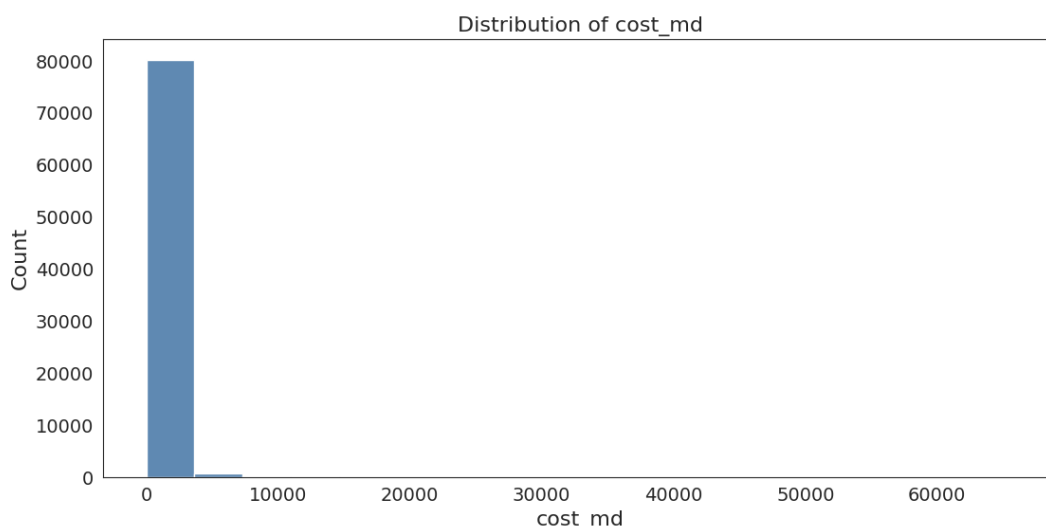
จำนวนกรมธรรม์ประเภทอื่นในบริษัท (client_nother) มีค่าเฉลี่ยเท่ากับ 0.2190 ส่วนเบี่ยงเบนมาตรฐานเท่ากับ 0.6553 ค่าต่ำสุดของข้อมูลเท่ากับ 0 ค่าที่อยู่ตำแหน่ง 25% ของข้อมูลเท่ากับ 0 ค่ามัธยฐานเท่ากับ 0 ค่าที่อยู่ตำแหน่ง 75% ของข้อมูลเท่ากับ 0 พบว่าส่วนใหญ่ไม่มีกรมธรรม์ประเภทอื่นในบริษัทเดียวกันและค่าสูงสุดของข้อมูลเท่ากับ 23

สำหรับตัวแปรที่เป็นข้อมูลเชิงคุณภาพ มีจำนวนค่าที่ไม่ซ้ำกันในตัวแปร (unique) เท่ากับ 2 และ ค่าอื่น ๆ ดังนี้ 1. เพศของผู้ถือกรมธรรม์ (client_sex) ส่วนใหญ่เป็นเพศชาย 2. พื้นที่อยู่อาศัย (city) ผู้ถือกรมธรรม์ส่วนใหญ่ไม่ได้อยู่ในเมืองใหญ่ โดยมีจำนวน 64,994 ราย 3. อาศัยอยู่ภาคเหนือหรือไม่ (north) 4. ผู้ถือกรมธรรม์ส่วนใหญ่ไม่ได้อาศัยพื้นที่ภาคเหนือ อาศัยอยู่พื้นที่อื่น ๆ ของประเทศ (rest) พบว่าส่วนใหญ่อาศัยอยู่ในพื้นที่อื่นของประเทศ คือไม่ได้อยู่ทั้งภาคเหนือและในเมืองใหญ่ จำนวน 41,636 ราย



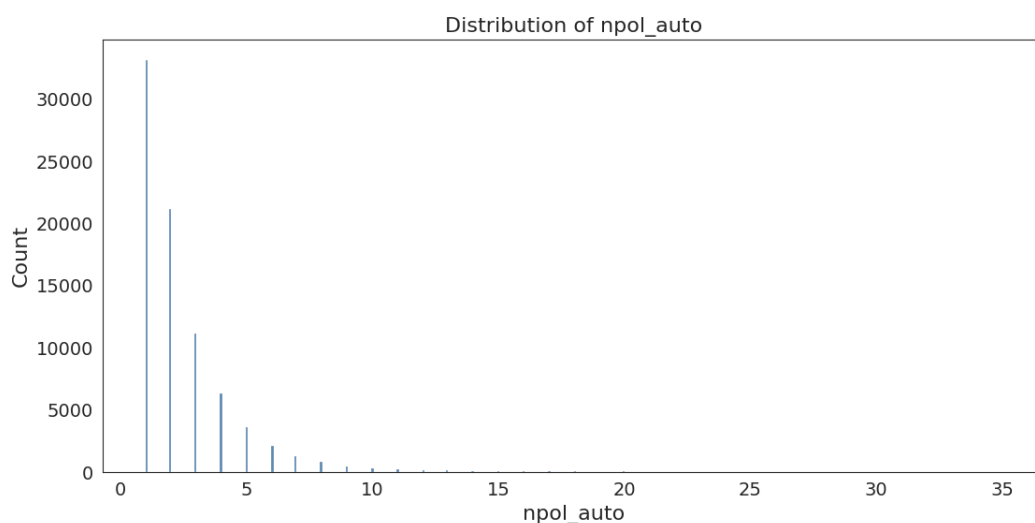
ภาพที่ 4.2 การกระจายตัวของข้อมูลความถี่ของการเรียกร้องค่าสินไหมทดแทน

จากภาพที่ 4.2 พบว่าจำนวนการเรียกร้องค่าสินไหมทดแทนประกันภัยรถยนต์ที่มีความเสียหายต่อทรัพย์สินมีการกระจายแบบเบ้ขวา ส่วนใหญ่มีจำนวนการเรียกร้องค่าสินไหมทดแทนต่ำมาก อยู่ในช่วง 0 ถึง 2 ครั้ง ส่งผลให้เกิดความไม่สมดุลในข้อมูล ซึ่งต้องพิจารณาปรับสมดุลข้อมูลเพื่อให้การสร้างแบบจำลองมีความแม่นยำและเสถียรมากขึ้น



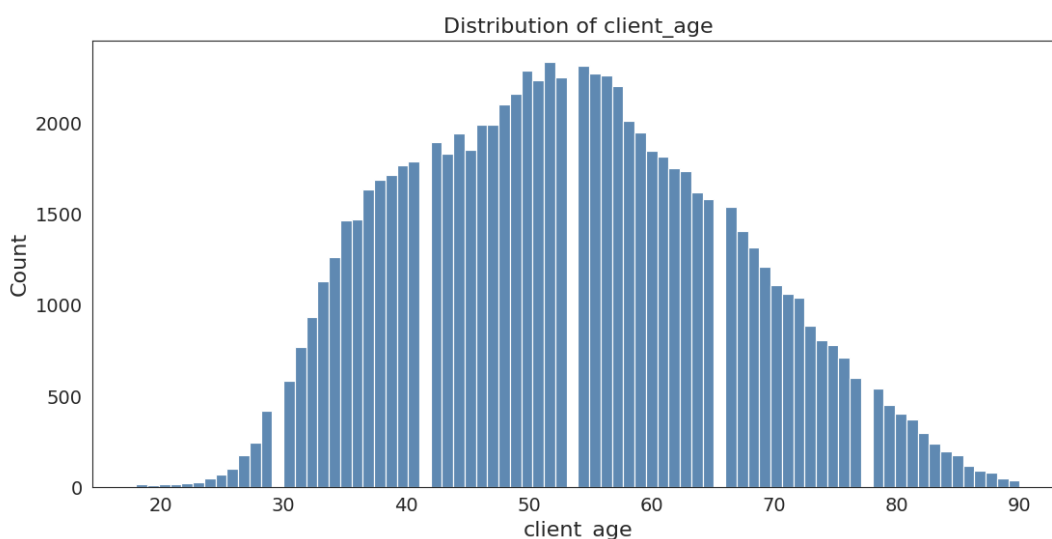
ภาพที่ 4.3 การกระจายตัวของข้อมูลความรุนแรงในการเรียกร้องค่าสินไหมทดแทน

จากภาพที่ 4.3 พบว่าความรุนแรงในการเรียกร้องค่าสินไหมทดแทนมีการกระจายแบบเบ้ขวา ส่วนใหญ่มีค่าต่ำกว่า 10000 การที่เป็น 0 มากหมายความว่าไม่เคยมีการเรียกร้องค่าสินไหมทดแทน ทำให้ไม่มีมูลค่าการเรียกร้องค่าสินไหมทดแทนทั้งหมดเช่นกัน



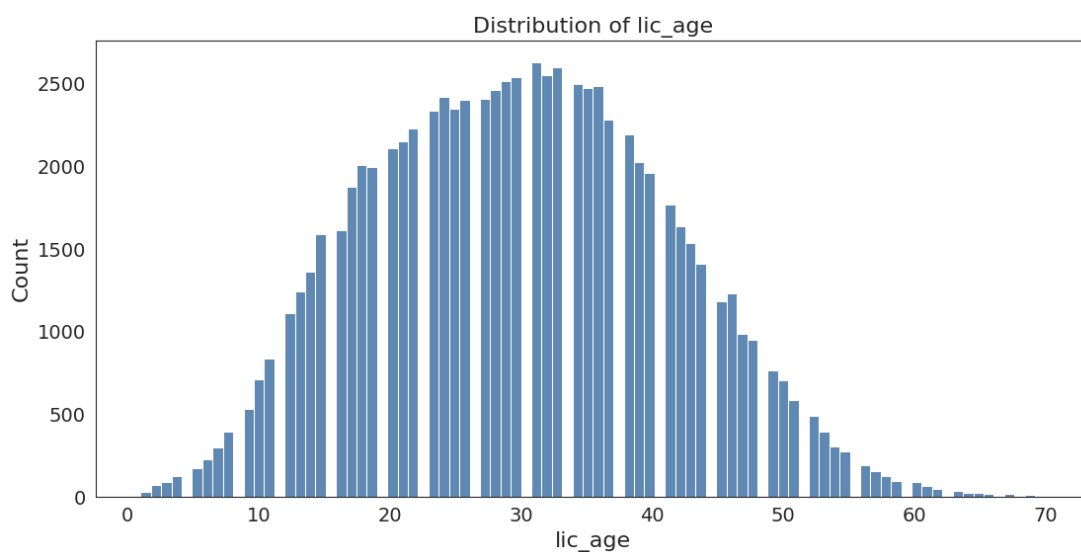
ภาพที่ 4.4 การกระจายตัวของข้อมูลจำนวนกรรมธรรม์ประกันภัยรถยนต์

จากภาพที่ 4.4 พบว่าจำนวนกรรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรรมธรรม์มีการกระจายแบบเบ้ขวา จำนวนกรรมธรรม์ที่รถยนต์ที่มีต่อบุคคลส่วนใหญ่อยู่ที่ 1 ถึง 3 เล่มกรรมธรรม์



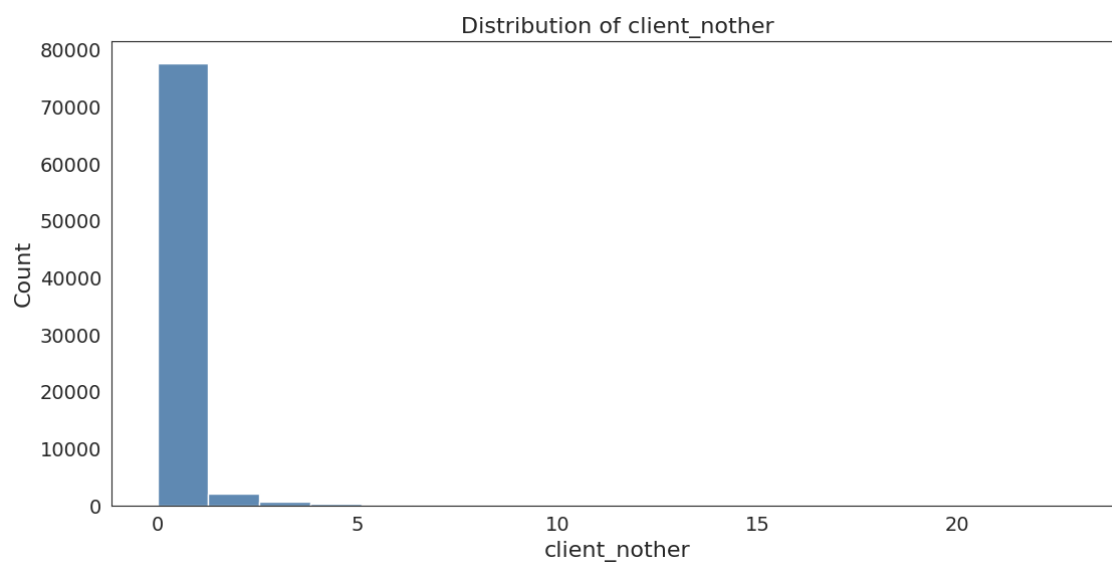
ภาพที่ 4.5 การกระจายตัวของข้อมูลอายุผู้ถือกรรมธรรม์

จากภาพที่ 4.5 พบว่าอายุของผู้ถือกรรมธรรม์ มีการกระจายตัวค่อนข้างสมมาตร อายุของผู้เอาประกันภัยส่วนใหญ่จะอยู่ในช่วงวัยกลางคนถึงวัยผู้ใหญ่ตอนปลาย คือช่วง 30 ถึง 70 ปี



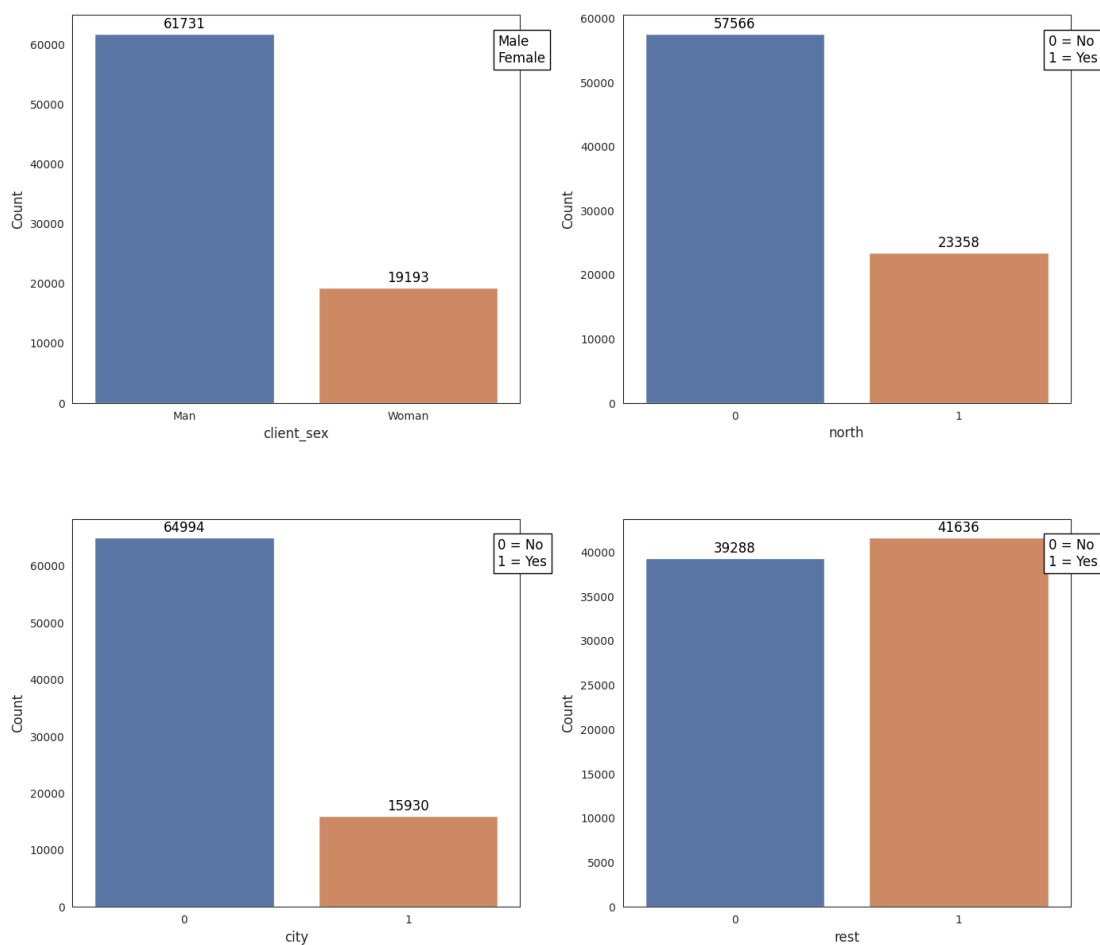
ภาพที่ 4.6 การกระจายตัวของข้อมูลอายุใบขับขี่

จากภาพที่ 4.6 พบว่าอายุของใบอนุญาตขับขี่ มีการกระจายตัวค่อนข้างสมมาตร ส่วนใหญ่มีใบอนุญาตขับขี่มานานประมาณ 20 ถึง 40 ปี



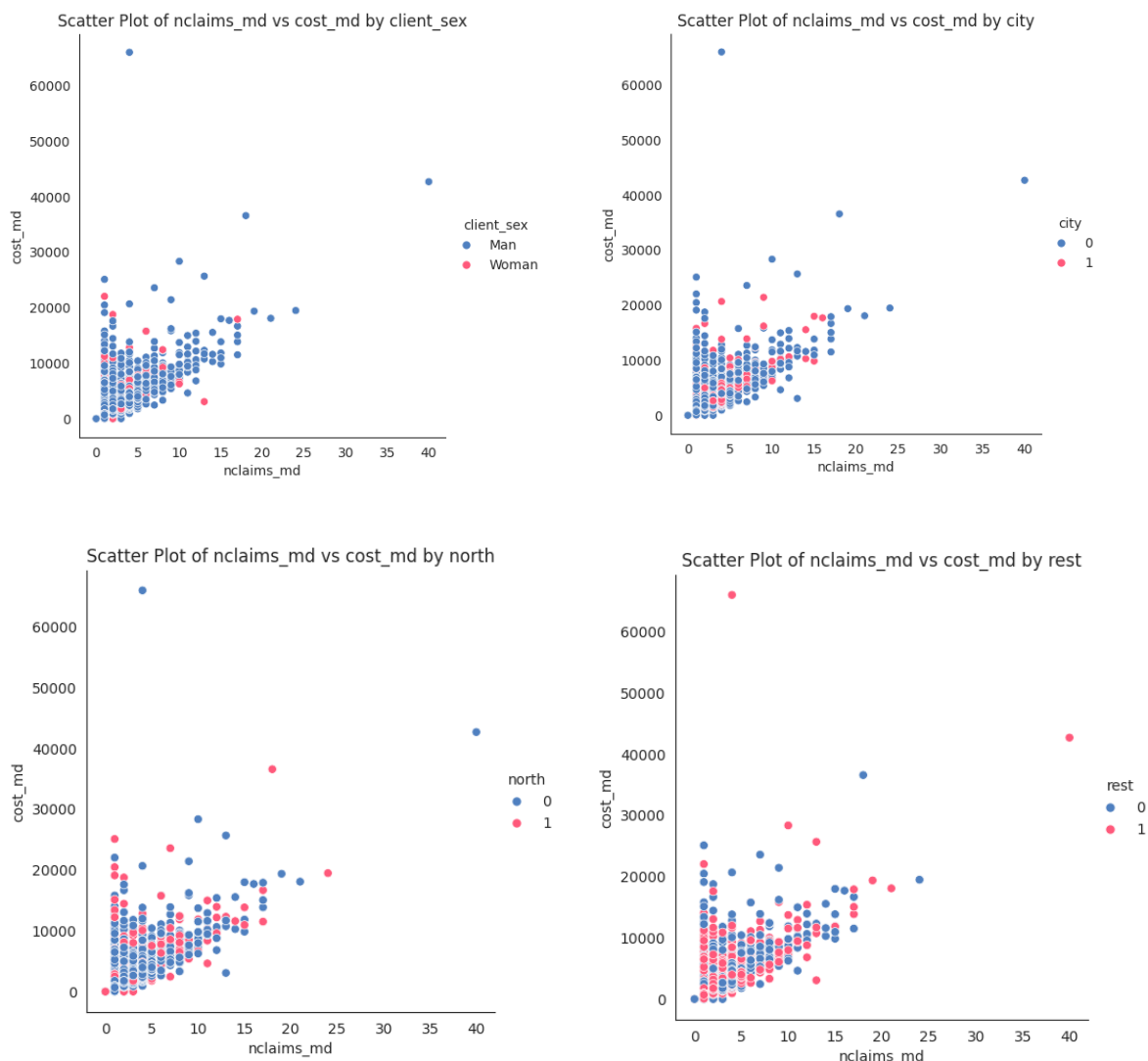
ภาพที่ 4.7 การกระจายตัวของข้อมูลจำนวนกรรมสิทธิ์ประกันภัยอื่น

จากภาพที่ 4.7 พบว่าจำนวนกรรมสิทธิ์ประเภทอื่นในบริษัท มีการกระจายแบบเบ้ขวา มีการถือกรรมสิทธิ์อื่นที่ไม่ใช่รถยนต์หรือบ้านค่อนข้างน้อยหรือไม่มีเลยคือ 0 ถึง 1 เล่มกรรมสิทธิ์



ภาพที่ 4.8 การตรวจสอบการกระจายของข้อมูลเชิงคุณภาพ

จากภาพที่ 4.8 สามารถอธิบายได้ดังนี้ เพศของผู้ถือกรรมธรรม์ (client_sex) พบว่าส่วนใหญ่ผู้ถือกรรมธรรม์เป็นเพศชาย พื้นที่อยู่อาศัย (city) พบว่าผู้ถือกรรมธรรม์ส่วนใหญ่ไม่ได้อาศัยอยู่ในเมืองใหญ่ (บาเซโลน่า หรือ มาดริด) อาศัยอยู่ภาคเหนือหรือไม่ (north) พบว่าผู้ถือกรรมธรรม์ส่วนใหญ่อยู่พื้นที่ที่ไม่ใช่ภาคเหนือ และการอาศัยอยู่พื้นที่อื่นๆของประเทศ (rest) พบว่าจำนวนคนที่อาศัยในเขตพื้นที่ทั้งในเมืองและภาคเหนือ มีความใกล้เคียงกับคนที่อาศัยในพื้นที่อื่น ๆ



ภาพที่ 4.9 Scatter plot แสดงความสัมพันธ์ระหว่างตัวแปรเชิงปริมาณกับเชิงคุณภาพ

จากภาพที่ 4.9 สามารถอธิบายได้ว่าส่วนใหญ่ข้อมูลมีการกระจายตัวอยู่ของจำนวนการเรียกร้องค่าสินไหมทดแทนมีค่าน้อยกว่า 10 และความรุนแรงในการเรียกร้องค่าสินไหมทดแทนอยู่ในช่วงไม่เกิน 20,000 ยูโร นอกจากนี้ยังชี้ให้เห็นถึงความสัมพันธ์ระหว่างจำนวนและความรุนแรงในการเรียกร้องค่าสินไหมทดแทนที่เป็นไปในเชิงบวก กล่าวคือเมื่อจำนวนการเรียกร้องค่าสินไหมทดแทนเพิ่มขึ้น ความรุนแรงในการเรียกร้องค่าสินไหมทดแทนก็มีแนวโน้มที่จะเพิ่มขึ้นตามไปด้วยและสามารถอธิบายแต่ละกราฟได้ดังนี้

Scatter Plot of nclaims_md vs cost_md by Client Sex ผู้ถือกรรมธรรม์ผู้ชายมีการกระจายตัวของจำนวนการเรียกร้องค่าสินไหมทดแทนในระดับสูงกว่า และบางกรณีก็มีความรุนแรงในการเรียกร้องค่าสินไหมทดแทนสูงกว่าผู้หญิง

Scatter Plot of nclaims_md vs cost_md by city ผู้ถือกรรมธรรม์ที่อาศัยอยู่ในเมืองใหญ่และผู้ที่ไม่ได้อาศัยอยู่ในเมืองใหญ่มีรูปแบบการกระจายตัวของการจำนวนและความรุนแรงในการเรียกร้องค่าสินไหมทดแทนคล้ายคลึงกัน แต่ในบางกรณีกลุ่มที่อาศัยในเมืองใหญ่มีแนวโน้มมากกว่า

Scatter Plot of nclaims_md vs cost_md by north ผู้ถือกรรมธรรม์ที่อาศัยอยู่ในภาคเหนือและผู้ที่ไม่ได้อาศัยอยู่ในภาคเหนือมีรูปแบบการกระจายตัวของการจำนวนและความรุนแรงในการเรียกร้องค่าสินไหมทดแทนคล้ายคลึงกัน

Scatter Plot of nclaims_md vs cost_md by rest ผู้ถือกรรมธรรม์อาศัยพื้นที่อื่น ๆ ของประเทศการกระจายตัวของความรุนแรงในการเรียกร้องค่าสินไหมทดแทนที่สูงกว่า หากเทียบกับจำนวนการเรียกร้องค่าสินไหมที่เท่ากัน

4.3 ผลการเตรียมข้อมูลก่อนการวิเคราะห์

4.3.1 ผลการทำความสะอาดข้อมูล

1) ผลการตรวจสอบความถูกต้องและความแนบเนียนของข้อมูล

เมื่อพิจารณาตัวแปร $client_age$ (X_3) คืออายุของผู้ถือกรรมธรรม์ และ lic_age (X_4) คืออายุของใบอนุญาตขับขี่ ดังนั้น " lic_age " จะต้องมีย่านน้อยกว่า " $client_age$ " เสมอ เพราะลูกค้าจะได้รับใบอนุญาตขับขี่หลังจากที่เกิดมาแล้ว และสำหรับอายุขั้นต่ำที่สามารถทำใบขับขี่ได้ในประเทศสเปน (ช่วงปี ค.ศ. 2006 ถึง ค.ศ. 2015) คืออายุ 18

ผู้วิจัยจึงคำนวณค่าความแตกต่างระหว่าง $client_age$ และ lic_age หากผลลัพธ์น้อยกว่า 18 ปี จะถือว่าข้อมูลนั้นผิดพลาดและจะถูกลบออกจากข้อมูล ดังตัวอย่างในตารางที่ 4.2 หลังตรวจสอบแล้วพบว่ามีข้อมูลที่น้อยกว่า 18 อยู่ 2,734 แถว จาก 80,924 แถว เมื่อลบออกแล้วเหลือข้อมูลอยู่ 78,190 แถว

ตารางที่ 4.2 ตัวอย่างข้อมูลที่มีการบันทึกผิดพลาด

$client_age$	lic_age	$client_age - lic_age$
25	10	15
30	22	8

พิจารณาตัวแปร $nclaims_md$ (Y_1) คือ จำนวนการเรียกร้องค่าสินไหมทดแทนประกันภัยรถยนต์ที่มีความเสียหายต่อทรัพย์สิน และ $cost_md$ (Y_2) คือ ความรุนแรงในการเรียกร้องค่าสินไหมทดแทนทั้งหมดสำหรับประกันภัยรถยนต์ที่มีความเสียหายต่อทรัพย์สิน (มูลค่า) เนื่องจากมีข้อมูลที่

nclaims_md มีค่า แต่ cost_md เป็น 0 หมายความว่า จำนวนการเรียกร้องค่าสินไหมขึ้นแต่ไม่มี การจ่ายค่าสินไหมทดแทน ซึ่งอาจมีสาเหตุมาจากการเรียกร้องถูกปฏิเสธเนื่องจากไม่อยู่ในเงื่อนไข ความคุ้มครองของกรมธรรม์ หรือ ค่าเสียหายต่ำกว่า Deductible (ค่าเสียหายส่วนแรกที่ผู้เอา ประกันภัยต้องรับผิดชอบเอง) ทำให้บริษัทประกันภัยไม่ต้องจ่ายเงิน จึงพิจารณาลบข้อมูลนี้ออกก่อน การวิเคราะห์ ซึ่งมีจำนวน 467 แถว จาก 78,190 แถว เมื่อลบแล้วเหลือข้อมูลอยู่ 77,723 แถว

2) ผลการจัดการกับค่าสูญหาย

จากการตรวจสอบข้อมูลในงานวิจัยนี้พบว่า ไม่มีค่าค่าสูญหาย แสดงผลดังภาพที่ 4.10

Number of missing values	
nclaims_md	0
cost_md	0
npol_auto	0
client_sex	0
client_age	0
lic_age	0
client_nother	0
city	0
north	0
rest	0

ภาพที่ 4.10 ผลการตรวจสอบค่าสูญหาย (Missing Value)

3) ผลการตรวจสอบค่านอกเกณฑ์

จากการวิเคราะห์พบว่ามีตัวแปรที่มีค่านอกเกณฑ์มากกว่า 5% ของข้อมูลทั้งหมดดังนี้

nclaims_md (Y_1) มีค่านอกเกณฑ์จำนวน 11,132 แถว (คิดเป็น 14.58%)

cost_md (Y_2) มีค่านอกเกณฑ์จำนวน 11,132 แถว (คิดเป็น 14.58%)

client_nother (X_8) มีค่านอกเกณฑ์จำนวน 11,598 แถว (คิดเป็น 14.92%)

มีรายละเอียดข้อมูลแสดงดังภาพที่ 4.11

Column: nclaims_md
 Lower Bound: 0.0
 Upper Bound: 0.0
 Number of outliers: 11332
 Percentage of outliers: 14.58%

Column: cost_md
 Lower Bound: 0.0
 Upper Bound: 0.0
 Number of outliers: 11332
 Percentage of outliers: 14.58%

Column: npol_auto
 Lower Bound: -2.0
 Upper Bound: 6.0
 Number of outliers: 3385
 Percentage of outliers: 4.36%

Column: client_age
 Lower Bound: 15.5
 Upper Bound: 91.5
 Number of outliers: 0
 Percentage of outliers: 0.00%

Column: lic_age
 Lower Bound: -4.5
 Upper Bound: 63.5
 Number of outliers: 86
 Percentage of outliers: 0.11%

Column: client_nother
 Lower Bound: 0.0
 Upper Bound: 0.0
 Number of outliers: 11598
 Percentage of outliers: 14.92%

ภาพที่ 4.11 ผลการตรวจสอบค่าผิดปกติ

4.3.2 ผลการแปลงข้อมูล

1) จากตัวแปรทั้งหมด มีตัวแปรเพศของผู้ถือกรมธรรม์ เป็นข้อมูลเชิงคุณภาพ มีค่าที่เป็นไปได้คือ Man, Woman เพื่อให้สามารถนำไปวิเคราะห์ในโปรแกรมได้ จึงทำการแปลงข้อมูลให้อยู่ในรูปของตัวเลข โดยกำหนดให้ 0 = Man, 1 = Woman

2) หลังจากการแปลงข้อมูลของตัวแปรตามแล้ว มีผลแสดงตัวอย่างข้อมูล ดังตารางที่ 4.3

ตารางที่ 4.3 ตัวอย่างข้อมูลการประกันภัยรถยนต์ที่ทำการแปลงเรียบร้อยแล้ว

	คนที่ 1	คนที่ 2	คนที่ 3	คนที่ 4	คนที่ 5
nclaims_md ก่อนการแปลงข้อมูล	0	2	1	0	1
nclaims_md_Log หลังการแปลงข้อมูล	-4.6051	0.6981	0.0099	-4.6051	0.0099
cost_md ก่อนการแปลงข้อมูล	0	1764	882	0	4967.47
cost_md_Log หลังการแปลงข้อมูล	-4.6051	7.4753	6.7822	-4.6051	8.5107
client_sex	0	0	0	1	1

4.4 ผลการสร้างแบบจำลอง

4.4.1 ผลการแบ่งข้อมูล

จากการตรวจสอบพบว่า ตัวแปรความถี่ในการเรียกร้องค่าสินไหมทดแทนเป็นข้อมูลไม่สมดุล ผู้วิจัยจึงแบ่งข้อมูลด้วยวิธี stratify ใช้ข้อมูล “Has_claim” ที่มาจากการแปลงข้อมูล nclaims_md ซึ่งตัวแปรนี้จะไม่ถูกนำไปสร้างแบบจำลอง แต่ใช้ในการสร้างความสมดุลในการแบ่งข้อมูลค่าของข้อมูลที่เป็นไปได้คือ

0 หมายถึง ไม่มีการเรียกร้องค่าสินไหมทดแทน

1 หมายถึง มีการเรียกร้องค่าสินไหมทดแทน ตั้งแต่ 1 ครั้งขึ้นไป

มีตัวอย่างการแปลงข้อมูลดังตารางที่ 4.4

ในงานวิจัยนี้จะใช้ library ที่ชื่อว่า model_selection คำสั่ง train_test_split (X, y, test_size=0.2, random_state=42, stratify=data_claims['Has_claim'])

ตารางที่ 4.4 ตัวอย่างการแปลงข้อมูล Has_claim

nclaims_md	Has_claim
0	0
1	1
4	1
0	0

หลังจากทำการแบ่งข้อมูลแล้ว สัดส่วนของกลุ่มข้อมูลทั้งสองชุดยังคงใกล้เคียงกัน โดยร้อยละของข้อมูลจำแนกตามการเรียกร้องค่าสินไหมทดแทนของข้อมูลชุดฝึกแสดงดังตารางที่ 4.5 และข้อมูลชุดทดสอบแสดงดังตารางที่ 4.6

ตารางที่ 4.5 ร้อยละของข้อมูลชุดฝึกจำแนกตามการเรียกร้องค่าสินไหมทดแทน

มีการเรียกร้องค่าสินไหมทดแทนใช่หรือไม่	จำนวน	ร้อยละ
ไม่ใช่	53,112	85.42
ใช่	9,066	14.58
รวม	62,178	100.00

ตารางที่ 4.6 ร้อยละข้อมูลชุดทดสอบจำแนกตามการเรียกร้องค่าสินไหมทดแทน

มีการเรียกร้องค่าสินไหมทดแทนใช่หรือไม่	จำนวน	ร้อยละ
ไม่ใช่	13,279	85.42
ใช่	2,266	14.58
รวม	15,545	100.00

4.4.1 ผลการสร้างแบบจำลองฐาน (Baseline Model)

จากการศึกษาข้อมูล ผู้วิจัยได้ปรับความสมดุลของข้อมูลด้วยฟังก์ชันล๊อคกาลิทึมธรรมชาติ และเลือกใช้ Mean Strategy ในการทำนาย ส่งผลให้ได้ค่าทำนายของตัวแปรผลลัพธ์ทั้งสองดังตารางที่ 4.7

ตารางที่ 4.7 ค่าทำนายพื้นฐานจากค่าเฉลี่ยของตัวแปรตามในชุดฝึกสอน

	nclaims_md	cost_md
Mean Prediction	0.0106	0.0451

4.4.2 ผลการปรับแต่งพารามิเตอร์ของแต่ละแบบจำลอง

ในขั้นตอนการปรับแต่งค่าพารามิเตอร์ด้วยวิธี Manual Search เพื่อหาขอบเขตของพารามิเตอร์ที่เป็นไปได้ของแต่ละแบบจำลอง ได้ขอบเขตพารามิเตอร์ดังตารางที่ 4.8

ตารางที่ 4.8 ขอบเขตของไฮเปอร์พารามิเตอร์กำหนดใน GridSearchCV ของแต่ละแบบจำลอง

แบบจำลอง	ไฮเปอร์พารามิเตอร์	ค่าที่กำหนด
Multi-Output Random Forest Regression	n_estimators max_depth min_samples_split min_samples_leaf	[1100, 1200, 1300] [6, 7, 8, 9] [100, 110, 120, 130, 140] [40, 50, 60]
Multi-Output Extreme Gradient Boosting Regression	n_estimators max_depth learning_rate min_child_weight subsample	[100, 200, 300] [5, 10, 15, 20] [0.05, 0.1, 0.2] [10, 15, 20] [0.6, 0.7, 0.8]
Multi-Output Light Gradient Boosting Regression	n_estimators max_depth learning_rate num_leaves min_data_in_leaf	[100, 200, 300] [-1, 5, 15, 25] [0.05, 0.1, 0.2] [30, 35, 40, 45] [60, 65, 70]
Multi-Output Artificial Neural Network	hidden_layer_sizes solver activation learning_rate	[(64,), (128,), (128,64), (256,128), (512,256)] ['adam'] ['identity', 'relu'] ['constant', 'adaptive']

ขั้นตอนการปรับแต่งค่าพารามิเตอร์ของแบบจำลองดำเนินการโดยใช้วิธี GridSearchCV ซึ่งเป็นกระบวนการที่ทดลองค่าพารามิเตอร์ทุกค่าที่เป็นไปได้ภายในช่วงที่กำหนด พร้อมกับการประเมินผลแบบ 10-fold Cross Validation (CV) เพื่อค้นหาชุดค่าพารามิเตอร์ที่ให้ผลลัพธ์ที่ดีที่สุดสำหรับแต่ละแบบจำลอง ใน ผลการวิเคราะห์จากการปรับแต่งพารามิเตอร์ แสดงไว้ใน ตารางที่ 4.9

ตารางที่ 4.9 พารามิเตอร์ที่ดีที่สุดจากการค้นหาด้วย GridsearchCV ของแต่ละแบบจำลอง

Model Multi-Output	Tuning Parameters	Best Parameters
Random Forest Regression	n_estimators: [1100, 1200, 1300]	1300
	max_depth: [6, 7, 8, 9]	9
	min_samples_split: [100, 110, 120, 130, 140]	130
	min_samples_leaf: [40, 50, 60]	60
Extreme Gradient Boosting Regression	n_estimators: [100, 200, 300]	200
	max_depth: [5, 10, 15, 20]	5
	learning_rate: [0.05, 0.1, 0.2]	0.05
	min_child_weight: [10, 15, 20]	10
Light Gradient Boosting Regression	subsample: [0.6, 0.7, 0.8]	0.6
	n_estimators: [100, 200, 300]	200
	max_depth: [-1, 5, 15, 25]	15
	learning_rate: [0.05, 0.1, 0.2]	0.05
Artificial Neural Network	num_leaves: [30, 35, 40, 45]	30
	min_data_in_leaf: [60, 65, 70]	60
	hidden_layer_sizes: [(64,),(128,),(128,64) (256,128), (512,256)]	(128, 64)
	solver: ['adam']	Adam
	activation: ['identity', 'relu']	relu
	learning_rate: ['constant', 'adaptive']	constant

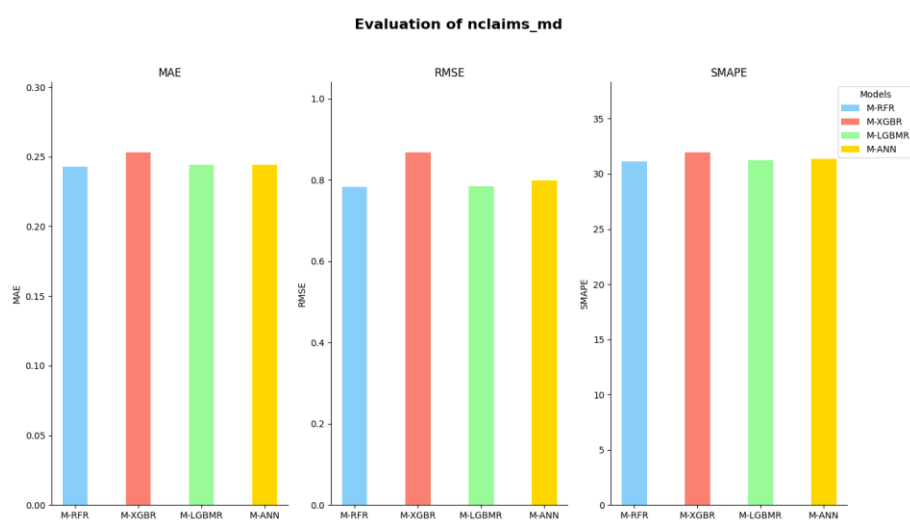
4.5 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลอง

4.5.1 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการถดถอยพหุผลลัพธ์

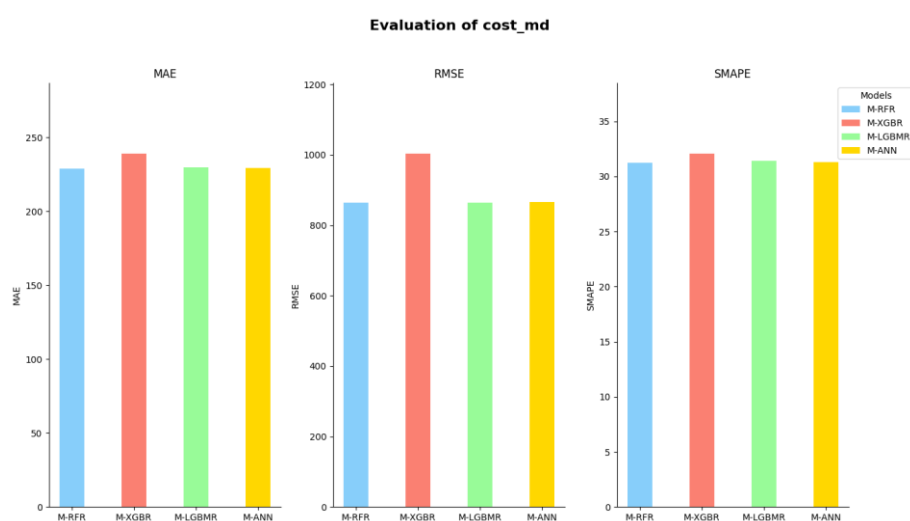
พารามิเตอร์ที่ดีที่สุดของแต่ละแบบจำลองจะถูกทดสอบกับชุดข้อมูลทดสอบเพื่อประเมินและเปรียบเทียบประสิทธิภาพในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ ทั้ง 4 แบบจำลอง โดยมีผลตัวเลขแสดงประสิทธิภาพดังในตารางที่ 4.10 และกราฟเปรียบเทียบประสิทธิภาพในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนดังภาพที่ 4.12 และ 4.13 ตามลำดับ

ตารางที่ 4.10 ผลการทดสอบประสิทธิภาพการทำนายของแต่ละแบบจำลอง

Multi-Output Model	nclaims_md			cost_md			Duration (seconds)
	MAE	RMSE	SMAPE	MAE	RMSE	SMAPE	
Baseline	0.2541	0.8050	199.522	229.8564	869.7132	199.996	3
M-RFR	0.2427	0.7832	28.8311	228.9800	864.6607	29.2370	1658
M-XGBoost	0.2520	0.8621	29.0140	238.9562	1003.0973	29.4286	24
M-LightGBM	0.2440	0.7863	28.9626	229.7467	865.2349	29.2945	35
M-ANN	0.2443	0.7998	28.8917	229.3862	866.5528	29.2581	1937



ภาพที่ 4.12 ประสิทธิภาพการทำนายความถี่ของการเรียกร้องค่าสินไหมทดแทน



ภาพที่ 4.13 ประสิทธิภาพการทำนายความความรุนแรงของการเรียกร้องค่าสินไหมทดแทน

จากตารางที่ 4.9 พบว่า Baseline Model แสดงค่าความคลาดเคลื่อนที่สูงกว่าแบบจำลองการเรียนรู้ของเครื่องในภาพรวม โดยแบบจำลองที่ให้ผลการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ได้ดีที่สุดในชุดข้อมูลทดสอบ คือ Multi-Output Random Forest Regression (M-RFR) สำหรับการทำนายความถี่ของการเรียกร้องค่าสินไหมทดแทน มีค่า MAE เท่ากับ 0.2427 ค่า RMSE เท่ากับ 0.7832 ค่า SMAPE เท่ากับ 31.1281 และการทำนายความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ มีค่า MAE เท่ากับ 228.9800 ค่า RMSE เท่ากับ 864.6607 และค่า SMAPE เท่ากับ 31.2500 และทำการทำนายโดยใช้ข้อมูลชุดทดสอบแสดงดังในตารางที่ 4.11

ตารางที่ 4.11 ค่าจริงและค่าจากการทำนายโดยใช้แบบจำลอง M-RFR จากข้อมูลชุดทดสอบ

ตัวแปร	คนที่ 1	คนที่ 2	คนที่ 3	คนที่ 4	คนที่ 5
client_sex	0	1	0	1	0
client_age	54	56	53	53	53
lic_age	32	37	33	28	35
city	0	1	0	0	0
north	1	0	1	0	0
rest	0	0	0	1	1
client_nother	0	0	0	5	0
npol_auto	13	2	12	24	8
cost_md	0	2170.76	358	3505.01	0
cost_md Predict	250.6056	0	146.5017	508.5177	0
nclaims_md	0	5	1	5	0
nclaims_md Predict	1	0	1	2	0

4.5.2 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการถดถอยพหุผลลัพธ์กับการถดถอยผลลัพธ์เดียว

จากผลการวิจัยในข้อ 4.5.1 พบว่าแบบจำลองที่สามารถทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนได้ดีที่สุดคือ การถดถอยพหุผลลัพธ์แบบสุ่มป่า (Multi-Output Random Forest Regression) จึงสร้างแบบจำลองการถดถอยพหุผลลัพธ์แบบสุ่มป่า ประเภทการถดถอยผลลัพธ์เดียว (Single-Output Regression) โดยทำนายเป็นแยกผลลัพธ์ ได้แก่ 1. ทำนายความถี่ของการเรียกร้องค่าสินไหมทดแทน และ 2. ทำนายความรุนแรงของการเรียกร้องค่าสินไหมทดแทน โดยมีการปรับแต่งพารามิเตอร์และค่าที่ดีที่สุด แสดงดังในตารางที่ 4.11

ตารางที่ 4.12 ค่าพารามิเตอร์จาก GridSearchCV ของ Single-Output Random Forest Regression

Random Forest Model	Tuning Parameters	Best Parameters
cost_md	n_estimators: [100, 200, 300]	100
	max_depth: [20, 22, 24, 26]	26
	min_samples_split: [90, 100, 110, 130, 150]	150
	min_samples_leaf: [50, 55, 60]	50
nclaims_md	n_estimators: [100, 200, 300]	100
	max_depth: [18, 20, 22, 24, 26]	18
	min_samples_split: [80, 100, 120, 140]	80
	min_samples_leaf: [40, 50, 60]	50

นำพารามิเตอร์ที่ดีที่สุดของแบบจำลองการถดถอยผลลัพธ์เดียวมาทดสอบกับข้อมูลชุดทดสอบ และเปรียบเทียบประสิทธิภาพของแบบจำลองการถดถอยผลลัพธ์เดียวกับการทำนายจากแบบจำลองการถดถอยพหุผลลัพธ์ ผลการทดสอบดังแสดงในตารางที่ 4.13

ตารางที่ 4.13 ผลการเปรียบเทียบประสิทธิภาพของแบบจำลองประเภทการถดถอยพหุผลลัพธ์กับการถดถอยผลลัพธ์เดียว

Random Forest Model	MAE	RMSE	SMAPE	Duration (second)
S - nclaims_md	0.2482	0.8074	29.8663	92
S - cost_md	225.5343	849.2932	198.1965	80
M - nclaims_md	0.2427	0.7832	31.1281	1658
M - cost_md	228.9800	864.6607	31.2500	

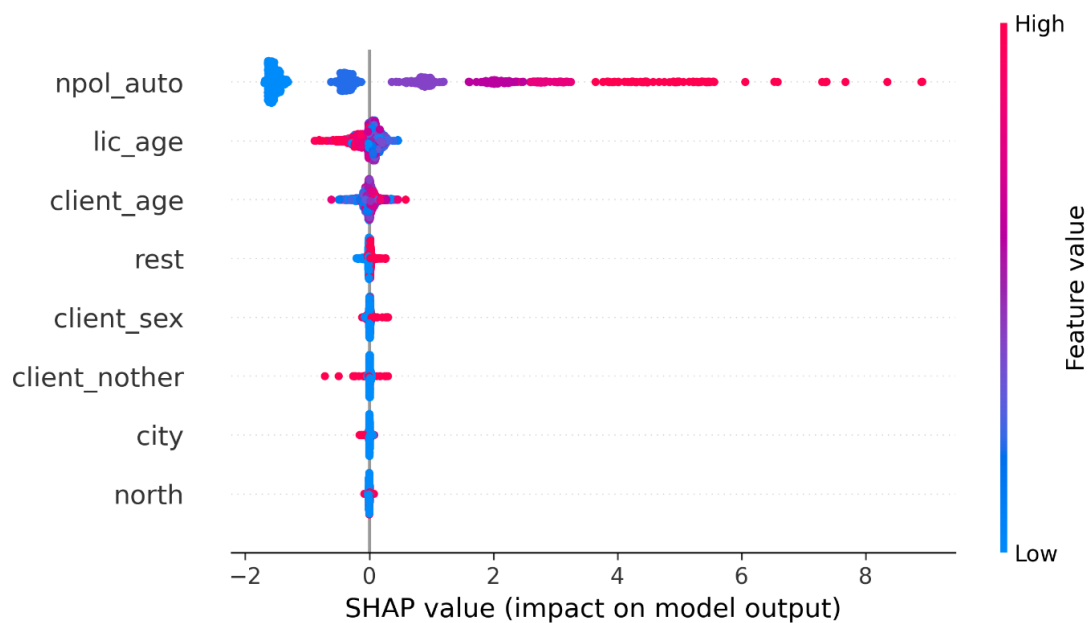
จากตารางที่ 4.13 พบว่าแบบจำลองประเภทการถดถอยพหุผลลัพธ์กับการถดถอยผลลัพธ์เดียวมีประสิทธิภาพใกล้เคียงกัน โดย

- การทำนายความถี่ในการเรียกร้องค่าสินไหมทดแทนแบบจำลองที่สามารถทำนายได้ดีคือ Multi-Output Random Forest Regression มีค่า MAE เท่ากับ 0.2427 ค่า RMSE เท่ากับ 0.7832 และค่า SMAPE เท่ากับ 31.1281

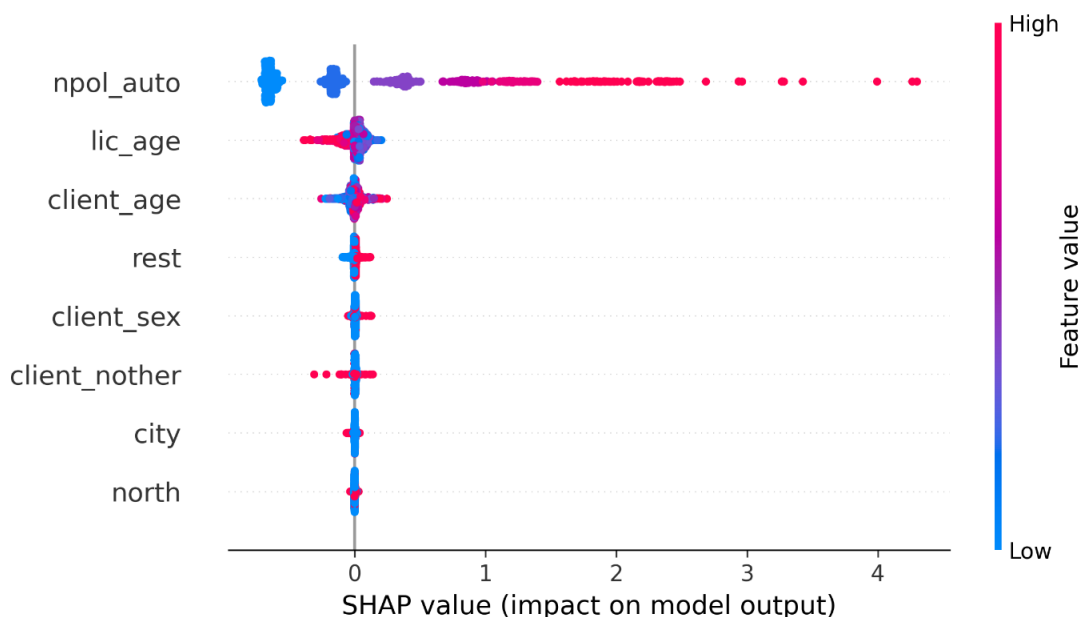
- การทำนายความรุนแรงในการเรียกร้องค่าสินไหมทดแทนแบบจำลองที่สามารถทำนายได้ดีคือ Single-Output Random Forest Regression มีค่า MAE เท่ากับ 225.5343 ค่า RMSE เท่ากับ 849.2932 และค่า SMAPE เท่ากับ 198.1965

4.5.3 ผลการอธิบายแบบจำลองด้วยเทคนิค SHAP

ผลการทำนายจาก Multi-Output Random Forest Regression สามารถอธิบายเพิ่มเติมได้โดยใช้เทคนิค Shapley Additive Explanation (SHAP) เพื่อให้เข้าใจถึงคุณลักษณะที่ส่งผลต่อการทำนายของแบบจำลองทั้งในเชิงบวกและเชิงลบ ผลการวิเคราะห์ SHAP แสดงในภาพที่ 4.14 และ 4.15



ภาพที่ 4.14 ผล SHAP ของแบบจำลอง M-RFR ในการทำนายความถี่ของการเรียกร้องค่าสินไหมทดแทน



ภาพที่ 4.15 ผล SHAP ของแบบจำลอง M-RFR ในการทำนายความรุนแรงของการเรียกร้องค่าสินไหมทดแทน

จากภาพที่ 4.14 และ 4.15 แสดงผล SHAP ที่สะท้อนถึงความสำคัญของคุณลักษณะต่อการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนของแบบจำลอง M-RFR พบว่าคุณลักษณะที่ผลกระทบมากต่อการทำนาย 3 อันดับแรก ได้แก่

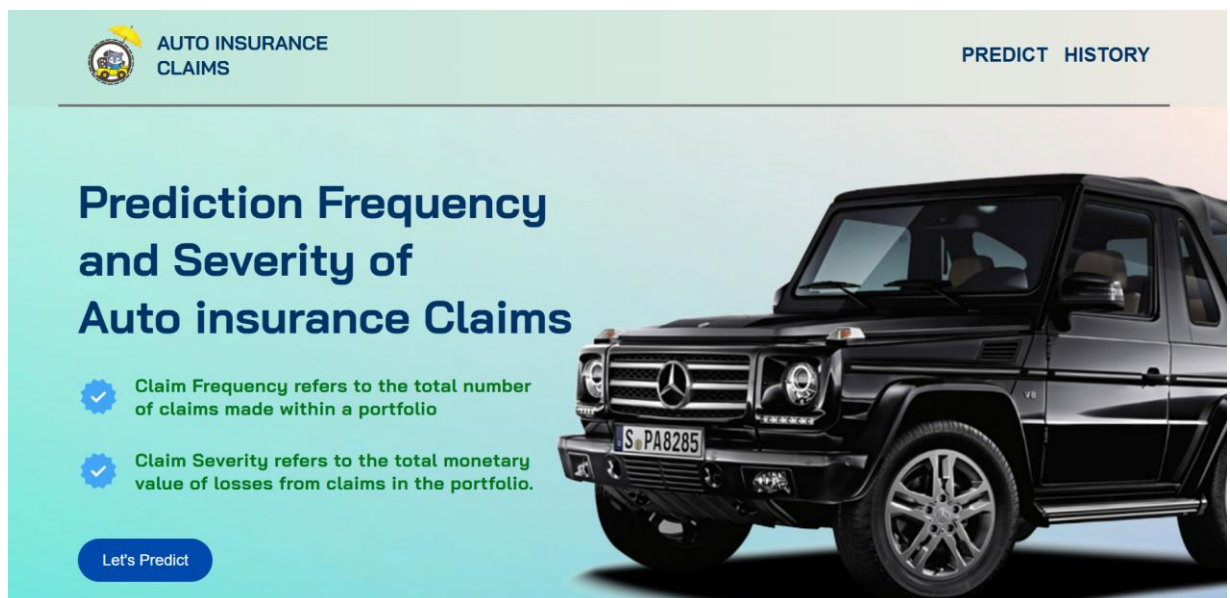
- จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี (npol_auto) มีผลกระทบเชิงบวกกับการทำนาย หมายความว่าผู้ถือกรมธรรม์ประกันภัยรถยนต์จำนวนมากมีแนวโน้มที่จะเรียกร้องค่าสินไหมทั้งความถี่และความรุนแรงมาก
- อายุของใบอนุญาตขับขี่ (lic_age) มีผลกระทบทั้งเชิงลบและเชิงบวก สามารถพิจารณาหลากหลายดังนี้ 1. หากมีประสบการณ์ในการขับขี่มากจะมีแนวโน้มการร้องค่าสินไหมทดแทนน้อยลง 2. หากประสบการณ์ในการขับขี่มากและอายุมาก อาจมีความเสี่ยงเพิ่มขึ้นจากข้อจำกัดด้านสมรรถภาพทางร่างกาย ส่งผลให้ความรุนแรงของการเรียกร้องค่าสินไหมสูง และ 3. หากผู้ถือกรมธรรม์มีประสบการณ์ในการขับขึ้น้อยลง มีแนวโน้มที่จะเรียกร้องค่าสินไหมทดแทนที่สูงขึ้น
- อายุของผู้ถือกรมธรรม์ (client_age) มีผลกระทบทั้งเชิงลบและเชิงบวก ซึ่งสัมพันธ์กับพฤติกรรมการขับขี่ สามารถพิจารณาหลากหลายดังนี้ 1. ผู้ถือกรมธรรม์ที่มีอายุในช่วงวัยกลางคนมักเป็นกลุ่มที่มีความรับผิดชอบสูงและมีพฤติกรรมการขับขี่ที่ระมัดระวัง ทำให้ความถี่และความรุนแรงของการเรียกร้องค่าสินไหมลดลง 2. ผู้ถือกรมธรรม์ที่อายุน้อยบางรายอาจขาดประสบการณ์ในการขับขี่ ส่งผลให้ความเสี่ยงต่อการเกิดอุบัติเหตุเพิ่มขึ้น และ

อาจทำให้การเรียกร้องค่าสินไหมมีความถี่หรือรุนแรงมากขึ้น 3. ผู้ถือกรมธรรม์ที่มีอายุเยอะมากอาจเผชิญข้อจำกัดด้านสมรรถภาพทางร่างกาย เช่น การมองเห็นหรือการตอบสนองที่ลดลง ซึ่งเพิ่มความเสี่ยงต่ออุบัติเหตุและส่งผลให้การเรียกร้องค่าสินไหมรุนแรงขึ้น

คุณลักษณะที่มีผลกระทบต่อการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน ได้แก่ การอาศัยอยู่ในพื้นที่อื่น ๆ ของประเทศ (rest) เพศของผู้ถือกรมธรรม์ (client_sex) จำนวนกรมธรรม์ประเภทอื่นในบริษัทเดียวกัน (client_nothor) พื้นที่อยู่อาศัย (city) การอาศัยอยู่ในภาคเหนือหรือไม่ (north) คุณลักษณะเหล่านี้มีผลกระทบต่อการทำนายในระดับต่ำ ซึ่งอาจพิจารณาตัดออกเพื่อลดความซับซ้อนของแบบจำลองและเพิ่มประสิทธิภาพในอนาคต

4.6 ผลการพัฒนาเว็บแอปพลิเคชัน

ผลการพัฒนาเว็บแอปพลิเคชันสำหรับการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ มีหน้าต่างแสดงดังภาพที่ 4.16 การกรอกข้อมูลดังภาพที่ 4.17, 4.18 และผลการทำนายดังภาพที่ 4.19



ภาพที่ 4.16 หน้าต่างเว็บแอปพลิเคชัน

Enter Individual Policyholder Information

n_pol	Number of motor insurance policies held by the client.
gender	Gender of the client
age	Age of the client.
license age	Number of years since the client obtained their driving license.
client_norther	The number of other types of insurance policies held by the client within the company ▼
city	The client resides in a major city ▼
north	The client resides in the northern region. ▼
rest	The client resides in other areas of the country ▼

*Rest is automatically generated based on other variables; no manual input is required.

ภาพที่ 4.17 หน้าต่างเว็บแอปพลิเคชันสำหรับกรอกข้อมูลผู้เอาประกันภัย (1)

Enter Individual Policyholder Information

16	Number of motor insurance policies held by the client.
Male	Gender of the client
26	Age of the client.
11	Number of years since the client obtained their driving license.
13	The number of other types of insurance policies held by the client within the company ▼
Yes	The client resides in a major city ▼
No	The client resides in the northern region. ▼
No	The client resides in other areas of the country ▼

*North is automatically populated based on City variables. No manual input is required.


*Rest is automatically generated based on other variables; no manual input is required.

ภาพที่ 4.18 หน้าต่างเว็บแอปพลิเคชันหลังจากกรอกข้อมูลผู้เอาประกันภัย (2)

13	The number of other types of insurance policies held by the client within the company ▼
Yes	The client resides in a major city ▼
No	The client resides in the northern region. ▼
No	The client resides in other areas of the country ▼

**North is automatically populated based on City variables. No manual input is required.*

**Rest is automatically generated based on other variables; no manual input is required.*



Predict

Claim Frequency	Claim Severity
1	220.3943

ภาพที่ 4.19 หน้าต่างเว็บแอปพลิเคชันที่แสดงผลการทำนาย



ภาพที่ 4.20 QR code: Web Application



ภาพที่ 4.21 QR code: GitHub

บทที่ 5

สรุปผลการวิจัย

การวิเคราะห์การถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ มีวัตถุประสงค์เพื่อสร้างและเปรียบเทียบแบบจำลองการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ด้วยการเรียนรู้ของเครื่อง และเพื่อพัฒนาเว็บแอปพลิเคชันที่ช่วยในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน ข้อมูลที่ใช้ในการวิจัยครั้งนี้คือข้อมูลการประกันภัยรถยนต์ของบริษัทแห่งหนึ่งในประเทศสเปนตั้งแต่ปี ค.ศ. 2006 ถึง ค.ศ. 2015 เป็นระยะเวลา 10 ปี รวมทั้งสิ้น 80,924 ราย มีสรุปผลการวิจัยที่ได้จากการวิเคราะห์แบบจำลองและประสิทธิภาพของเว็บแอปพลิเคชัน พร้อมทั้งอภิปรายผลลัพธ์และเสนอแนะแนวทางในการนำไปประยุกต์ใช้ รวมถึงแนวทางสำหรับการศึกษาวิจัยในอนาคตได้ มีการสรุปผลการวิจัยดังนี้

5.1 สรุปผลการวิจัย

5.2 อภิปรายผลการวิจัย

5.3 ประโยชน์ของสถิติ/สารสนเทศสถิติที่ใช้ในการวิจัย

5.3.1 การแปลงข้อมูลก่อนการวิเคราะห์

5.3.2 การเรียนรู้ของเครื่อง

5.3.3 การพัฒนาเว็บแอปพลิเคชัน

5.4 ข้อเสนอแนะ

5.4.1 ข้อเสนอแนะในการนำผลการวิจัยไปใช้

5.4.2 ข้อเสนอแนะในการทำวิจัยครั้งต่อไป

5.1 สรุปผลการวิจัย

งานวิจัยนี้ใช้แบบจำลองการถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงในการเรียกร้อยค่าสินไหมทดแทนทั้งหมด 4 แบบจำลอง ได้แก่ Random Forest, Extreme Gradient Boosting, Light Gradient Boosting Machine และ Artificial Neural Network พบว่าประสิทธิภาพของแบบจำลองการเรียนรู้ของเครื่องทั้งหมดมีประสิทธิภาพดีกว่าแบบจำลองฐานในภาพรวม และแบบจำลอง Random Forest ให้ผลการทำนายที่แม่นยำที่สุดเมื่อเปรียบเทียบกับแบบจำลองอื่น ๆ โดยมีผลคือ

1. การทำนายความถี่ของการเรียกร้อยค่าสินไหมทดแทนในการประกันภัยรถยนต์มีค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (MAE) เท่ากับ 0.2427 ค่ารากที่สองของความคลาดเคลื่อนกำลังสองเฉลี่ย (RMSE) เท่ากับ 0.7832 ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (SMAPE) เท่ากับ 28.8311
2. การทำนายความรุนแรงของการเรียกร้อยค่าสินไหมทดแทนในการประกันภัยรถยนต์ มีค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (MAE) เท่ากับ 228.9800 ค่ารากที่สองของความคลาดเคลื่อนกำลังสองเฉลี่ย (RMSE) เท่ากับ 864.6607 ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (SMAPE) เท่ากับ 29.2370

ผลการอธิบายแบบจำลองด้วย Shapley Additive Explanation (SHAP) พบว่าตัวแปรที่มีผลต่อการทำนายทั้งความถี่และความรุนแรงในการเรียกร้อยค่าสินไหมทดแทนมากที่สุดสามอันดับแรก ได้แก่ จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี อายุของใบอนุญาตขับขี่ และอายุของผู้ถือกรมธรรม์ ตามลำดับ

จากผลการวิจัยพบว่า แบบจำลองการถดถอยพหุผลลัพธ์ที่ใช้ Random Forest มีประสิทธิภาพในการทำนายสูงสุด จึงได้พัฒนาแบบจำลองการถดถอยผลลัพธ์เดียวโดยใช้ Random Forest เช่นเดียวกัน เพื่อเปรียบเทียบประสิทธิภาพกับแบบจำลองการถดถอยพหุผลลัพธ์ ผลการเปรียบเทียบพบว่าแบบจำลองการถดถอยพหุผลลัพธ์มีความแม่นยำในการทำนายความถี่ของการเรียกร้อยค่าสินไหมทดแทนดีกว่า ขณะที่แบบจำลองการถดถอยผลลัพธ์เดียวมีความแม่นยำในการทำนายความรุนแรงของการเรียกร้อยค่าสินไหมทดแทนสูงกว่า ซึ่งทั้งสองแบบจำลองมีประสิทธิภาพในการทำนายโดยรวมใกล้เคียงกัน โดยผลการวิเคราะห์แบบจำลองผลลัพธ์เดียวพบว่า

1. การทำนายความถี่ในการเรียกร้อยค่าสินไหมทดแทน มีค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (MAE) เท่ากับ 0.2482 ค่ารากที่สองของความคลาดเคลื่อนกำลังสองเฉลี่ย (RMSE) เท่ากับ

0.8074 ค่าร้อยละความคลาดเคลื่อน สัมบูรณ์เฉลี่ยแบบสมมาตร (SMAPE) เท่ากับ 29.8663

2. การทำนายความรุนแรงในการเรียกร้องค่าสินไหมทดแทน มีค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ย (MAE) เท่ากับ 225.5343 ค่ารากที่สองของความคลาดเคลื่อน กำลังสองเฉลี่ย (RMSE) เท่ากับ 849.2932 ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (SMAPE) เท่ากับ 198.1965

และผู้วิจัยได้พัฒนาเว็บแอปพลิเคชันเพื่อให้บริษัทประกันภัยรถยนต์สามารถเข้าถึงและทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนได้สะดวกและรวดเร็ว การทำงานของระบบเริ่มจากผู้ใช้กรอกข้อมูลของผู้เอาประกันภัยที่หน้าเว็บแอปพลิเคชัน จากนั้นกดทำนายเพื่อให้ระบบประมวลผลด้วยแบบจำลอง Multi-Output Random Forest Regression และผลลัพธ์จะแสดง 2 ค่า คือ ค่าความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์

5.2 อภิปรายผลการวิจัย

การสร้างแบบการถดถอยพหุผลลัพธ์ในการทำนายความถี่และความรุนแรงในการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์ โดยใช้แบบจำลองการเรียนรู้ของเครื่อง จำนวน 4 แบบจำลอง ได้แก่ Random Forest, Extreme Gradient Boosting (XGboost), Light Gradient Boosting Machine (LightGBM), Artificial Neural Network (ANN) ผลการวิจัยพบว่าแบบจำลองทั้งหมดมีค่าความคลาดเคลื่อนต่ำกว่าแบบจำลองฐานในภาพรวม แม้ว่าความแตกต่างจะไม่มากแต่แสดงให้เห็นว่าแบบจำลองสามารถเรียนรู้และทำนายข้อมูลที่มีความซับซ้อนได้ดีกว่า และผลการวิจัยพบว่า Random Forest มีประสิทธิภาพในการทำนายสูงสุดเมื่อเปรียบเทียบกับแบบจำลองอื่น เนื่องจากความสามารถในการลดผลกระทบจากค่านอกเกณฑ์ด้วยวิธีการสุ่มตัวอย่างข้อมูลและการสุ่มเลือกคุณลักษณะ ในกระบวนการสร้างต้นไม้แต่ละต้น ส่งผลให้ข้อมูลที่มีค่านอกเกณฑ์ไม่ถูกนำมาใช้ในการสร้างต้นไม้ทุกต้น จึงช่วยลดความคลาดเคลื่อนและเพิ่มความแม่นยำในการทำนาย นอกจากนี้ Random Forest ยังมีความยืดหยุ่นสูงและสามารถจัดการกับข้อมูลที่มีความไม่สมดุลได้ดี (Liaw and Wiener, 2002) และเพื่อทำความเข้าใจตัวแปรที่มีผลต่อการทำนาย ผู้วิจัยใช้ Shapley Additive Explanation (SHAP) ในการวิเคราะห์ความสำคัญของตัวแปร พบว่าตัวแปรที่ส่งผลต่อการทำนายทั้งความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทน 3 อันดับแรก ได้แก่ จำนวนกรมธรรม์ของผู้ถือกรมธรรม์ อายุของผู้ถือกรมธรรม์ และอายุของใบอนุญาตขับขี่ โดยอายุของผู้ถือกรมธรรม์ที่

เพิ่มขึ้นส่งผลให้ค่าทำนายสูงขึ้นซึ่งสอดคล้องกับงานวิจัยของ Jin (2021) ที่ระบุว่าอายุของผู้เอาประกันภัยเป็นหนึ่งในตัวแปรที่มีความสำคัญสูงสุด 5 อันดับแรก ที่ส่งผลต่อค่าทำนายในการเรียกร้องค่าสินไหมทดแทน

ระยะเวลาในการสร้างแบบจำลองพบว่า XGBoost ใช้เวลาสร้างแบบจำลองน้อยที่สุดที่ 24 วินาที รองลงมาคือ LightGBM ที่ 35 วินาที เนื่องจากทั้งสองใช้เทคนิค Gradient Boosting ที่มีการทำงานแบบขนาน (Parallel Processing) ซึ่งช่วยให้การคำนวณและการฝึกแบบจำลองเกิดขึ้นได้รวดเร็วขึ้น โดยเฉพาะในข้อมูลขนาดใหญ่ ในขณะที่ Random Forest ต้องสร้างต้นไม้และเลือกจุดแบ่งที่ดีที่สุดในแต่ละต้นไม้ ทำให้เวลาในการประมวลผลมากถึง ใช้เวลาที่ 27 นาที 38 วินาที และแบบจำลองที่ใช้เวลามากที่สุดคือ ANN ใช้เวลา 32 นาที 17 วินาที เนื่องจากต้องคำนวณฟังก์ชันการกระตุ้นในแต่ละ layer จึงทำให้ใช้เวลามากในการฝึกแบบจำลอง นอกจากนี้ปัจจัยที่ส่งผลต่อระยะเวลาการสร้างแบบจำลอง ได้แก่ ขนาดข้อมูล จำนวนพารามิเตอร์ และโครงสร้างแบบจำลอง

การเปรียบเทียบประสิทธิภาพการทำนายของแบบจำลองการถดถอยพหุผลลัพธ์ (Multi-Output Regression) กับแบบจำลองการถดถอยผลลัพธ์เดียว (Single-Output Regression) พบว่าประสิทธิภาพของทั้งสองแบบจำลองมีความใกล้เคียงกันในภาพรวม แต่แบบจำลอง Multi-Output Regression มีข้อได้เปรียบในด้านความสะดวก เนื่องจากสามารถทำนายหลายผลลัพธ์ได้พร้อมกันในครั้งเดียวนั้นช่วยลดขั้นตอนในการฝึกแบบจำลอง ดังนั้นเพื่อประยุกต์ในการใช้งานจริง ผู้วิจัยได้พัฒนาเว็บแอปพลิเคชันเชื่อมต่อกับแบบจำลอง Multi-Output Random Forest Regression ผ่าน Flask ซึ่งทำหน้าที่เป็น API รับข้อมูลจากผู้ใช้และส่งไปยังแบบจำลองเพื่อคำนวณค่าทำนาย ก่อนแสดงผลลัพธ์บนหน้าเว็บแอปพลิเคชัน ระบบนี้ช่วยให้บริษัทประกันภัยรถยนต์สามารถใช้งานได้สะดวกและรวดเร็ว โดยผู้ใช้กรอกข้อมูลของผู้เอาประกันภัยลงในระบบ และกดทำนาย จากนั้นระบบจะประมวลผลและแสดงผลลัพธ์ประกอบด้วยค่าทำนายสองค่าได้แก่ ความถี่และความรุนแรงของการเรียกร้องค่าสินไหมในการประกันภัยรถยนต์

5.3 ประโยชน์ของสถิติ/สารสนเทศสถิติที่ใช้ในการวิจัย

5.3.1 การแปลงข้อมูลก่อนการวิเคราะห์

การแปลงข้อมูลด้วยฟังก์ชันลอการิทึมธรรมชาติ (Logarithmic Transformation) ช่วยหลีกเลี่ยงค่าทำนายที่ติดลบ ลดการกระจายตัวของข้อมูลที่มีการเบี่ยงเบนมาก (skewness) และผลกระทบของค่าผิดปกติ (outliers) ทำให้มีการกระจายตัวของข้อมูลมีความสมดุลยิ่งขึ้น ส่งผลให้ผลลัพธ์สอดคล้องกับข้อมูลจริงและการทำนายมีความแม่นยำมากขึ้น

5.3.2 การเรียนรู้ของเครื่อง

การสร้างแบบจำลองทำนายด้วยการเรียนรู้ของเครื่อง (Machine Learning) สามารถทำนายหลายผลลัพธ์พร้อมกันได้ และมีการปรับปรุงประสิทธิภาพของแบบจำลองหลายวิธี เช่น การปรับแต่งพารามิเตอร์ (Hyperparameter Tuning) เพื่อหาค่าพารามิเตอร์ที่เหมาะสมที่สุดสำหรับแต่ละแบบจำลอง และการใช้เทคนิค Cross-Validation (CV) เพื่อประเมินความสามารถของแบบจำลองในชุดข้อมูลที่แตกต่างกัน วิธีดังกล่าวช่วยเพิ่มความแม่นยำในการทำนายและทำให้แบบจำลองทำนายได้อย่างมีประสิทธิภาพในสถานการณ์ที่หลากหลาย

5.3.3 การพัฒนาเว็บแอปพลิเคชัน

การพัฒนาเว็บแอปพลิเคชันด้วย Python, HTML, CSS และ JavaScript โดยเชื่อมต่อแบบจำลองผ่าน Flask ช่วยให้สามารถรับส่งข้อมูลระหว่างผู้ใช้ (Users) และเซิร์ฟเวอร์ได้อย่างมีประสิทธิภาพ พร้อมทั้งส่งค่าทำนายจากแบบจำลองกลับมาแสดงผลบนหน้าเว็บได้โดยอัตโนมัติ ซึ่งช่วยให้การใช้งานแบบจำลองมีความสะดวกและเข้าถึงได้ง่ายยิ่งขึ้น

5.4 ข้อเสนอแนะ

5.4.1 ข้อเสนอแนะในการนำผลการวิจัยไปใช้

1) บริษัทประกันภัยสามารถนำผลการวิจัยไปใช้ในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนพร้อมกัน ที่มาจากผู้เอาประกันรายเดียวกัน ซึ่งมีข้อมูลต่าง ๆ สอดคล้องกับงานวิจัยนี้ โดยเฉพาะอย่างยิ่งในประเทศสเปนหรือภูมิภาคยุโรปที่มีความเสี่ยงที่คล้ายกัน สำหรับบริษัทที่มีข้อมูลของตัวแปรอิสระแตกต่างจากงานวิจัยนี้ อาจพิจารณาสร้างแบบจำลอง Random Forest โดยใช้ตัวแปรอิสระที่มีมาใช้ในการทำนายการเรียกร้องค่าสินไหมทดแทน ซึ่งสามารถประยุกต์ใช้วิธีการเตรียมข้อมูล และขั้นตอนการสร้างแบบจำลองจากงานวิจัยนี้ ซึ่งจะช่วยให้บริษัทประกันภัยสามารถคาดการณ์แนวโน้มของการเรียกร้องค่าสินไหมทดแทนในอนาคตได้อย่างแม่นยำและเตรียมการรับมือได้อย่างเหมาะสม โดยแบบจำลอง Random Forest จะใช้เวลาในการประมวลผลค่อนข้างนานหากมีข้อมูลจำนวนมาก บริษัทประกันภัยอาจใช้พิจารณาใช้แบบจำลองที่มีระยะเวลาในการฝึกน้อยกว่าแต่ประสิทธิภาพใกล้เคียงกับ Random Forest ได้แก่ แบบจำลอง LightGBM

2) การวิเคราะห์ผลกระทบของปัจจัยต่าง ๆ ด้วยเทคนิค Shapley Additive Explanation (SHAP) พบว่า จำนวนกรมธรรม์ประกันภัยรถยนต์ที่ผู้ถือกรมธรรม์มี อายุของใบอนุญาตขับขี่ และอายุของผู้ถือกรมธรรม์ เป็นปัจจัยสำคัญที่ส่งผลต่อการเรียกร้องค่าสินไหมทดแทน

เป็นข้อมูลช่วยสนับสนุนการตัดสินใจในการกำหนดเบี้ยประกันภัยของผู้เอาประกันภัยแต่ละรายอย่างเหมาะสม รวมถึงนำไปใช้ในการวางแผนบริหารจัดการความเสี่ยงอื่น ๆ ในลำดับต่อไป

3) เนื่องจากข้อมูลที่ใช้ในการวิจัยเป็นข้อมูลในอดีต ซึ่งอาจไม่สอดคล้องกับบริบทสภาพแวดล้อม หรือพฤติกรรมของผู้เอาประกันภัยเมื่อเวลาเปลี่ยนแปลงไป ดังนั้นควรมีการอัปเดตแบบจำลองเมื่อเวลาผ่านไประยะหนึ่ง นอกจากนี้อาจพิจารณานำปัจจัยอื่น ๆ เพิ่มเติมมาใช้ในการสร้างแบบจำลอง เช่น การเปลี่ยนแปลงของปัจจัยทางเศรษฐกิจ สังคม และพฤติกรรมของผู้เอาประกันภัยในช่วงเวลาต่าง ๆ เพื่อให้สอดคล้องกับสถานการณ์ปัจจุบัน ซึ่งอาจส่งผลต่อความถี่และมูลค่าของการเรียกร้องค่าสินไหมทดแทน

5.4.2 ข้อเสนอแนะในการทำวิจัยครั้งต่อไป

1) ข้อมูลที่ใช้ในการสร้างแบบจำลองเป็นข้อมูลในอดีต ซึ่งอาจไม่สะท้อนถึงการเปลี่ยนแปลงของตลาดหรือพฤติกรรมของผู้เอาประกันภัยเมื่อเวลาเปลี่ยนแปลงไป ดังนั้น ในการวิจัยครั้งต่อไป ควรพิจารณาใช้ข้อมูลที่เป็นปัจจุบันมากขึ้น และเพิ่มตัวแปรที่เกี่ยวข้อง เช่น ปัจจัยทางเศรษฐกิจและสังคม เพื่อให้แบบจำลองมีความแม่นยำและสามารถสะท้อนสภาพแวดล้อมที่เปลี่ยนแปลงได้ดียิ่งขึ้น

2) ในการวิจัยครั้งนี้ได้ดำเนินการทดลองเฉพาะแบบจำลองต้นไม้ประเภท Ensemble ได้แก่ Random Forest, XGBoost และ LightGBM และโครงข่ายประสาทเทียม (ANN) เท่านั้น ในการวิจัยครั้งต่อไป อาจพิจารณานำแบบจำลองประเภทอื่น ๆ เช่น แบบจำลองการถดถอยเชิงเส้น (Linear Regression), การถดถอยด้วยเวกเตอร์สนับสนุน (SVR) และแบบจำลอง K-Nearest Neighbors (KNN) เพื่อขยายขอบเขตของการศึกษาและเปรียบเทียบประสิทธิภาพของแบบจำลองที่หลากหลาย

3) ควรพิจารณาจัดการข้อมูลก่อนนำมาสร้างแบบจำลองเพื่อปรับปรุงประสิทธิภาพของแบบจำลองให้ดียิ่งขึ้น เนื่องจากในการวิจัยครั้งนี้ได้ใช้ข้อมูลจริง โดยไม่ได้ดำเนินการตัดหรือจัดการค่าที่อยู่นอกเกณฑ์ (outliers) ซึ่งอาจส่งผลต่อความแม่นยำของผลลัพธ์ อย่างไรก็ตาม จากการทดลองเบื้องต้นพบว่า เมื่อมีการตัดค่าที่อยู่นอกเกณฑ์ออกแล้ว ค่าร้อยละความคลาดเคลื่อนสัมบูรณ์เฉลี่ยแบบสมมาตร (SMAPE) ลดลง ซึ่งบ่งชี้ว่าการจัดการค่าที่อยู่นอกเกณฑ์อาจช่วยเพิ่มความแม่นยำของแบบจำลองได้

4) การใช้เทคนิคการเลือกคุณลักษณะ (feature selection) อาจช่วยเพิ่มประสิทธิภาพของแบบจำลองได้ โดยการคัดเลือกเฉพาะคุณสมบัติที่มีความสำคัญต่อการทำนายจะช่วยลดความซับซ้อนของแบบจำลองและอาจส่งผลให้มีความแม่นยำสูงขึ้น นอกจากนี้อาจพิจารณาตัว

แปรอื่นเข้ามาร่วมด้วยเช่น previous claims severity, past claim frequency เพื่อเพิ่มความแม่นยำของแบบจำลอง

5) ควรพิจารณาการสร้างแบบจำลองโดยแบ่งกลุ่มข้อมูลออกเป็น 3 กลุ่ม ได้แก่ Non-claim, Low-claim, High-claim เพื่อช่วยให้สามารถจำแนกพฤติกรรมการเรียกร้องค่าสินไหมทดแทนของลูกค้าได้อย่างละเอียดมากขึ้น นอกจากนี้ การใช้เทคนิคการกำหนดค่าน้ำหนักตัวอย่าง (sample weight) อาจช่วยเพิ่มประสิทธิภาพของแบบจำลอง โดยให้ความสำคัญกับกลุ่มข้อมูลที่มีความสำคัญมาก ส่งผลให้แบบจำลองมีความแม่นยำและสอดคล้องกับโครงสร้างข้อมูลมากขึ้น

เอกสารอ้างอิง

- กิตติศักดิ์ จังพานิช, ศุภเจษฎา สีวันนา, วรรณพร เซวน์ชวานิล, กุลภัสสรณ์ ชีวลักษณะณาสีทธิ์, และธีระวัฒน์ สีมากันทร. (2564). การประยุกต์ใช้การเรียนรู้ของเครื่องในการตรวจจับการเรียกร้องค่าสินไหมทดแทนประกันภัยรถยนต์ที่ไม่มีการออกสำรวจภัย. *วารสารวิทยาศาสตร์มหาวิทยาลัยหัวเฉียวเฉลิมพระเกียรติ*, 27(2), 50-63.
- นพัตถ วิระมิตรชัย (2562). การพยากรณ์ยอดขายประกันภัย กรณีศึกษา ธุรกิจประกันภัย. *วารสารวิจัยมหาวิทยาลัยธุรกิจบัณฑิต*, 9(3), 21-35.
- บริษัท ไรจ์ จำกัด. (2566). **ประกันรถยนต์คืออะไร**. ค้นเมื่อ 1 กรกฎาคม 2567, จาก <https://www.roojai.com/insurance-glossary/what-is-car-insurance/>
- ปวีศา สุขเรื่อย และ สำรวม จงเจริญ. (2561). ตัวแบบการถดถอยที่มีผลกระทบจากค่าศูนย์ประยุกต์ใช้กับจำนวนครั้งของการเรียกร้องค่าสินไหมทดแทนในประกันภัยรถยนต์ภาคสมัครใจ. *วารสารมหาวิทยาลัยศรีนครินทรวิโรฒ (สาขาวิทยาศาสตร์และเทคโนโลยี)* ปีที่ 10 ฉบับที่ 20 กรกฎาคม - ธันวาคม 2561.
- ปฏิภาณ ประเสริฐสม และ พีรตล สามะศิริ. (2566). **ตีความโมเดล Machine Learning: ตัวอย่างและการตีความ Shapley value**. ค้นเมื่อ 10 ตุลาคม 2567, จาก <https://bdi.or.th/big-data-101/shapley-value-example/>
- สมาคมประกันวินาศภัยไทย. (2023). **คาดผลประกอบการธุรกิจประกันวินาศภัย ปี 66**. The General Insurance Association of Thailand. สืบค้นเมื่อ 8 กรกฎาคม 2567. https://www.tgia.org/newsandevents-detail-EN_1344_1
- อัศรพล พรหมพิริยะพงษ์. (2566). **การวิเคราะห์คุณลักษณะที่สามารถคาดการณ์ปริมาณการเพาะปลูกแบบเกษตรแม่นยำ: กรณีศึกษาผลผลิตข้าว.วิทยานิพนธ์ปริญญาโท สาขาวิชาการแปรรูปธุรกิจทางดิจิทัล คณะ วิทยาศาสตร์มหาบัณฑิต มหาวิทยาลัยธรรมศาสตร์**
- Athiwat. (2019). **Machine Learning คืออะไร?**. Medium. Retrieved July 1, 2024. from <https://bit.ly/3Un9MhJ>
- Ahmed, M. W. (2023). **Understanding mean absolute error (MAE) in regression: A practical guide**. Retrieved July 5, 2024, from <https://medium.com/@m.waqar.ahmed/under-standing-mean-absolute-error-mae-in-regression-a-practical-guide-26e80ebb97df>

- Amarin TV. (2567). **สินเชื่อเงินด่วน สมัครง่าย ไม่ต้องค้ำประกัน เช็กได้ที่นี่**. Amarin TV. สืบค้นเมื่อ 8 กรกฎาคม พ.ศ. 2567. <https://www.amarintv.com/spotlight/money-product/detail/37444>
- Abebe, M., Shin, Y., Noh, Y., Lee, S. & Lee, I. (2020). Machine Learning Approaches for Ship Speed Prediction toward Energy Efficient Shipping. **Applied Sciences**, 10(7).
- BDI. (2020). **Introduction to reinforcement learning**. Retrieved September 1, 2024, from <https://bdi.or.th/big-data-101/introduction-to-reinforcement-learning/>
- Binariks. (2024). **The Role of Big Data in Personalizing Insurance**. Retrieved November 8, 2024. https://binariks.com/blog/big-data-in-insurance-personalization/?fbclid=IwY2xjawGanHRleHRuA2FlbQlxMAABHdU2ecPwChD6Wka4xMOZYVW8_GpdUENV18cCf8rqCc7dvqwVtmGVex6wPw_aem_sHYGaj8WVrWHM3a6F5lW-w
- Bentejac, C., Csorgo, A., Munoz, M. G. (preprint). **A Comparative Analysis of XGBoost**. <https://doi.org/10.48550/arXiv.1911.01914>
- Bolancé, C., Vernic, R. (2017). Multivariate count data generalized linear models: Three approaches based on the Sarmanov distribution. **IREA Working Papers 201718**. University of Barcelona
- Borchani, H., Varando, G., Bielza, C., & Larrañaga, P. (2015). A survey on multi-output regression. **WIREs Data Mining and Knowledge Discovery**, 5(5), 216-233
- Chen, Y., Hu, M., Xie, Y., Qiu, R. (2020). Claim frequency predicting based on LightGBM. **Journal of Nonlinear and Analysis**, 22(8), 1757-1770.
- Daroontham, W. (2563). **เจาะลึก Random Forest (Part 2 of รู้จัก Decision Tree, Random Forest และ XGBoost)**. ค้นเมื่อ 5 กรกฎาคม 2567, จาก <https://bit.ly/40hcc5f>
- Daroontham, W. (2563). **รู้จัก Decision Tree, Random Forrest และ XGBoost (Part 1)**. ค้นเมื่อ 9 สิงหาคม 2567, จาก <https://bit.ly/4eTTsgM>
- Eastgate Software. (2024). **What is unsupervised learning?**. Retrieved July 25, 2024, from <https://eastgate-software.com/what-is-unsupervised-learning/>
- FINNOMENA. (2565). **รู้จัก Light Gradient Boosting Machine (LightGBM) โมเดลสุดล้ำสำหรับงานด้านการเงิน**. สืบค้นเมื่อ 27 สิงหาคม, 2567. <https://www.finnomena.com/finnomena-ic/light-gradient-boosting-machine-model/>

- GeeksforGeeks. (2023). **Multioutput regression in machine learning**. Retrieved August 20, 2024, from <https://www.geeksforgeeks.org/multioutput-regression-in-machine-learning/>
- Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow (2nd ed.). O'Reilly Media.
- Hmong.in.th. (n.d.). **Symmetric mean absolute percentage error**. Retrieved January 28, 2025, from https://hmong.in.th/wiki/Symmetric_mean_absolute_percentage_error
- Hong, L., Kuffner, T., Martin, R. (2023). **On Prediction of Future Insurance Claims When the Model Is Uncertain**, 12(1), 90-99.
- Inn Why. (2567). Insure World : จัปตา!วินาศภัยไทย เบี้ยทะลุ3.7แสนล้าน! ประกันทรัพย์สิน/รถ/ฟೀเอ-สุขภาพตันเตบโต. คั่นเมื่อ 8 พฤศจิกายน 2567. <https://shorturl.at/YLUmX>
- Jin, F. F. (2021). Using decision tree ensemble methods for the estimation of individual claims reserving (Master's thesis, Erasmus School of Economics). Erasmus University.
- Johnson, A. (2024). **Car & Moto Insurance In Spain >> Get The Right Cover At The Right Price**. Moving TO SPAIN. Retrieved November 7, 2024. https://movingtospain.com/car-insurance-in-spain/?fbclid=IwY2xjawGanF5leHRuA2FlbQlXMAABHa-v4lZFq4OjvzTr5r0FxeR10J59VwSG1qDdvMfTJJJug284wjqeXOQurA_aem_Q7bB-2P6vwhnKV_DiVXvNQ
- Krzyk, K. (2023). **Coding deep learning for beginners – Types of machine learning**. Retrieved July 3, 2024, from <https://resources.experfy.com/ai-ml/coding-deep-learning-for-beginners-types-of-machine-learning/>
- Kutner, H. M., Nachtsheim, J. C., Neter, J., Li, W. (2004). **Applied Linear Statistical Models**. 5 th ed. [n.p.]: McGraw-Hill.
- Kumar, V. S., Satpathi, D. K., Kumar, P. T. V. P., Haragopal, V. V. (2020). Modeling and Predicting of Motor Insurance Claim Amount using Artificial Neural Network. **International Journal of Recent Technology and Engineering** ISSN:2277-3878
- Liaw, A., Wiener, M. (2002). Classification and Regression by randomForest. **R News** ISSN:1609-3631 Vol. 2(3), 18-22.

- Linusson, H. (2013). MULTI-OUTPUT RANDOM FOREST (Magisteruppsats i Informatik, University of Borås). DiVA portal.
- Naqa, E. I., Murphy, M. J., & Li, H. (2022). **What Are Machine and Deep Learning**. In Machine and Deep Learning in Oncology, Medical Physics and Radiology (pp. 3-15). https://doi.org/10.1007/978-3-030-83047-2_1
- Patel, H. (n.d.). How To Implement Baselines In ML Modeling And Why We Need Them?. Census blog. <https://census.ai/blogs/how-to-implement-baselines-in-ml-modeling#blogpost-toc-0>
- Pathmind. (n.d.). **Decision tree**. Pathmind Wiki. Retrieved July 27, 2024. <https://wiki.pathmind.com/decision-tree>
- pigabyte. (2022, December 1). เผย 3 เทรนด์สำคัญ “ธุรกิจประกันรถยนต์” ในอนาคต พร้อมทิศทางการเข้ามาของ EV ที่จะพลิกโฉมประกันรถยนต์รูปแบบใหม่. Marketing Oops! <https://www.marketingoops.com/reports/industry-insight/priceza-money-car-insurance/>
- Poufnas, T., Gogas, P., Papadimitriou, T., Zaganidis, E. (2023). Machine Learning in Forecasting Motor Insurance Claims. **Risks** 2023, 11(9), 164. <https://doi.org/10.3390/risks11090164>
- Schagen, V. S. (2023). **Multi Target XGBoost Cash Flow Prediction An Efficient Machine Learning Algorithm For Future Liability Projections** (Master's thesis, Master of Science). The Delft University of Technology.
- Scikit learn. (2025). DummyRegressor. [https://scikit-learn.org/stable/modules/generated/sklearn.dummy.DummyRegressor.html#](https://scikit-learn.org/stable/modules/generated/sklearn.dummy.DummyRegressor.html#arn.dummy.DummyRegressor.html#)
- Segal, M. R. (1992). Tree-structured methods for longitudinal data. Journal of the American Statistical Association, 87(418), 407–418.
- Sun, T., Yang, J., Li, J., Chen, J., Liu, M., Fan, L., & Wang, X. (2024). Enhancing auto insurance risk evaluation with transformer and SHAP. IEEE Access, PP(99), 1–1
- WHO. (2024). Road traffic injuries. Retrieved November 7, 2024. <https://shorturl.asia/H1C0n>
- XGBoost. (n.d.). Multi-output. Retrieved August 27, 2024. <https://xgboost.readthedocs.io/en/stable/tutorials/multioutput.html>

ภาคผนวก

ภาคผนวก ก
ตารางการดำเนินโครงการ

ตารางที่ ก.1 ตารางการดำเนินโครงการ

ขั้นตอนการดำเนินงาน	ระยะเวลาดำเนินงานวิจัย ระหว่าง มิถุนายน 2567 - มีนาคม 2568																			
	มิถุนายน 2567				กรกฎาคม 2567				สิงหาคม 2567				กันยายน 2567				ตุลาคม 2567			
	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4
1. ศึกษาค้นคว้างานวิจัยที่สนใจ																				
2. เข้าพบอาจารย์ที่ปรึกษาวิจัย เพื่อ กำหนดหัวข้อวิจัย																				
3. ศึกษาบทความและงานวิจัยที่เกี่ยวข้อง																				
4. จัดทำบทที่ 1 บทนำ ประกอบด้วย ความเป็นมาและความสำคัญของปัญหา วัตถุประสงค์ ขอบเขตของงานวิจัย ประโยชน์ที่คาดว่าจะได้รับ																				
5. จัดทำบทที่ 2 ทบทวนวรรณกรรมที่ เกี่ยวข้อง																				
5.1 การเรียนรู้ของเครื่อง																				
5.2 การเรียนรู้เชิงลึก																				
5.3 การฝึกสอนแบบจำลอง																				
5.4 แบบจำลองการเปรียบเทียบ ประสิทธิภาพแบบจำลอง																				
5.5 งานวิจัยที่เกี่ยวข้อง																				
6. จัดทำบทที่ 3 วิธีการดำเนินงานวิจัย																				
6.1 การทำความเข้าใจธุรกิจ																				
6.2 การศึกษาและทำความเข้าใจข้อมูล																				

ตารางที่ ก.2 ตารางการดำเนินงาน

ขั้นตอนการดำเนินงาน	ระยะเวลาดำเนินงานวิจัย ระหว่าง มิถุนายน 2567 - มีนาคม 2568																			
	มิถุนายน 2567				กรกฎาคม 2567				สิงหาคม 2567				กันยายน 2567				ตุลาคม 2567			
	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4
6.3 การเตรียมข้อมูล																				
- การแปลงข้อมูล																				
- การทำความสะอาดข้อมูล																				
- ทำ EDA																				
6.4 การวิเคราะห์ข้อมูล																				
- Random Forest																				
- ExtreemGradientBoosting																				
- LightGradientBoosting																				
- Artificial Neural Network																				
6.5 ประเมินประสิทธิภาพแบบจำลอง																				
6.6 การนำแบบจำลองไปใช้งานจริง																				
7. เข้าพบอาจารย์ที่ปรึกษาเพื่อปรับปรุงแก้ บทที่ 1-3																				
8. เตรียมสอบเค้าโครงวิจัย																				
9. นำเสนอเค้าโครงงานวิจัยบทที่ 1-3																				

ตารางที่ ก.3 ตารางการดำเนินโครงการ

ขั้นตอนการดำเนินงาน	ระยะเวลาดำเนินงานวิจัย ระหว่าง มิถุนายน 2567 - มีนาคม 2568																			
	พฤศจิกายน 2567				ธันวาคม 2567				มกราคม 2568				กุมภาพันธ์ 2568				มีนาคม 2568			
	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4
10. จัดทำบทที่ 4 ผลการวิจัย																				
10.1 วิเคราะห์ผลข้อมูลแบบจำลอง																				
10.2 ตารางสรุปผลการเปรียบเทียบประสิทธิภาพแบบจำลอง																				
11. จัดทำบทที่ 5 สรุปผลการวิจัย																				
11.1 สรุปผลการวิจัย																				
11.2 อภิปรายผลการวิจัย																				
11.3 ข้อเสนอแนะ																				
12. เข้าพบอาจารย์ที่ปรึกษาเพื่อปรับปรุงแก้บทที่ 4-5																				
13. จัดทำรูปเล่มโครงการวิจัย																				
14. จัดทำสื่อสำหรับนำเสนอโครงการวิจัย																				
15. เตรียมตัวนำเสนอโครงการวิจัย																				
16. นำเสนอโครงการวิจัย																				
17. แก้ไขและปรับปรุงรูปเล่มโครงการตามข้อเสนอแนะของกรรมการ																				

ภาคผนวก ข
งบประมาณค่าใช้จ่ายในการดำเนินโครงการ

ตารางที่ ข.1 งบประมาณค่าใช้จ่ายในการดำเนินโครงการ

รายการ	จำนวนเงิน (บาท)
ค่าจัดพิมพ์เอกสารโครงการ	1000
รวม	1000

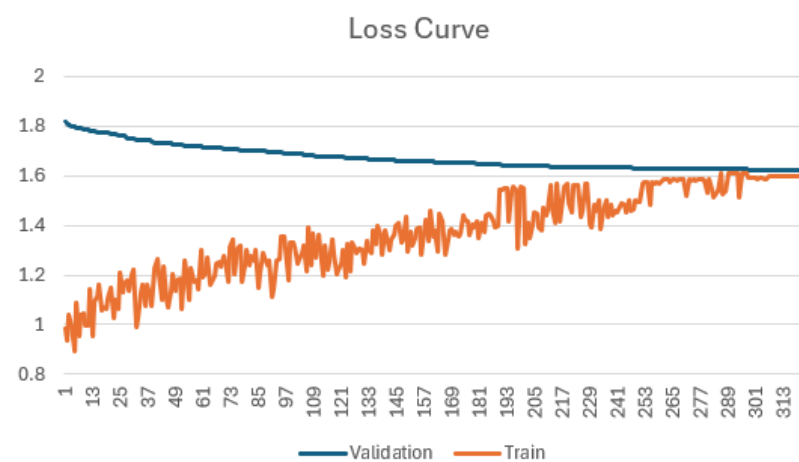
ภาคผนวก ค

กราฟแสดงผลการเรียนรู้ของแบบจำลอง

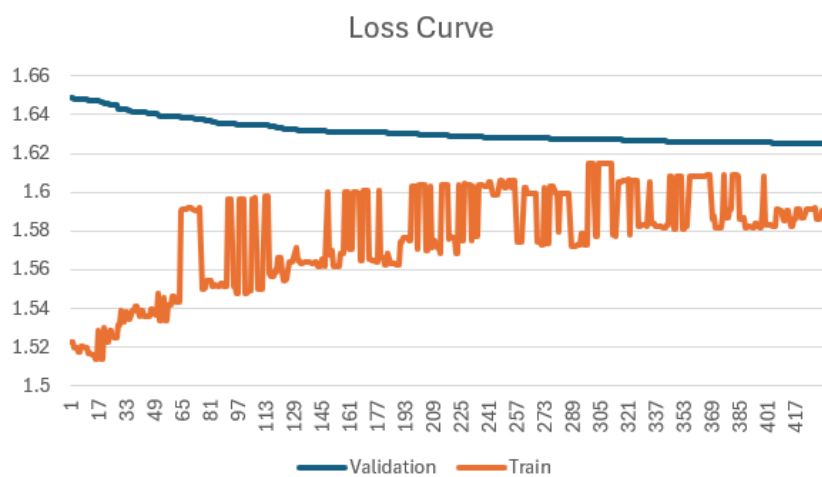
ใช้ Scoring คือ MAE



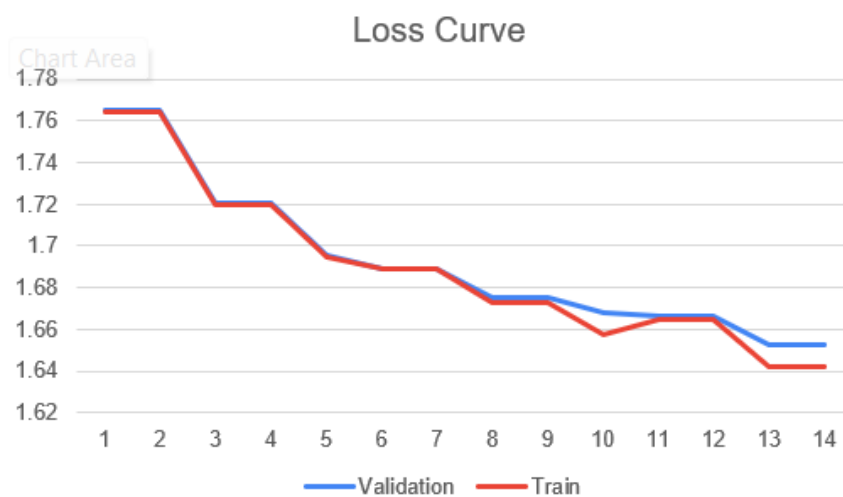
ภาพที่ ค.1 Loss Curve ของแบบจำลอง Random Forest



ภาพที่ ค.2 Loss Curve ของแบบจำลอง XGBoost



ภาพที่ ค.3 Loss Curve ของแบบจำลอง LightGBM



ภาพที่ ค.4 Loss Curve ของแบบจำลอง ANN

ภาคผนวก ง

Implementation Code

ตัวอย่างโค้ดที่ใช้ในการพัฒนาแบบจำลองในการทำนายความถี่และความรุนแรงของการเรียกร้องค่าสินไหมทดแทนในการประกันภัยรถยนต์

จ.1 การแบ่งข้อมูลสำหรับพัฒนาแบบจำลอง

```
# สมมุติว่า 'Has_claim' เป็นคอลัมน์ที่บอกว่ามีการเรียกร้องค่าสินไหมทดแทนประกันหรือไม่
X = data_claims_cleaned[['npol_auto', 'client_sex', 'client_age', 'lic_age',
'client_nother', 'cities2', 'north', 'rest']]
y = data_claims_cleaned[['nclaims_md', 'cost_md']]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42, stratify=data_claims_cleaned['Has_claim'])
# สร้าง DataFrame สำหรับชุดฝึก
train_data = pd.concat([X_train, y_train], axis=1)

# สร้าง DataFrame สำหรับชุดทดสอบ
test_data = pd.concat([X_test, y_test], axis=1)
```

จ.2 ตัวอย่างการฝึกแบบจำลอง Multi-output Random Forest

```
# กำหนดพารามิเตอร์สำหรับ Grid Search
param_grid = {
    "estimator__n_estimators": [1400,1500],
    "estimator__max_depth": [10],
    "estimator__min_samples_split": [130],
    "estimator__min_samples_leaf": [60]
}
kf = KFold(n_splits=10, shuffle=True, random_state=42)

# ฟังก์ชันสำหรับการ Fit GridSearchCV และบันทึกผลลัพธ์
def run_multioutput_grid_search_optimized(X_train, y_train, param_grid):
    rf = MultiOutputRegressor(RandomForestRegressor(random_state=42))
    grid_search = GridSearchCV(
```



```

estimator=rf,
param_grid=param_grid,
scoring="neg_mean_absolute_error",
cv=kf,
return_train_score=True,
verbose=1,
n_jobs=-1,
)

# รัน GridSearch โดยใช้ฟังก์ชัน fit_with_sample_weights
start_time = time.time()
grid_search.fit(X_train, y_train) # ส่ง sample_weight ตรงๆ ไปที่ fit()
total_time = time.time() - start_time

# สรุปผล
results_df = pd.DataFrame(grid_search.cv_results_)
results_df["total_time"] = total_time
results_df["avg_fold_time"] = total_time / kf.get_n_splits()

return results_df, grid_search.best_params_

# รัน Multi-Output GridSearchCV
print("Running Multi-Output GridSearchCV...")
results, best_params = run_multioutput_grid_search_optimized(
    X_train, y_train, param_grid
)
print(f"Best parameters: {best_params}")

# บันทึกเฉพาะพารามิเตอร์ที่ดีที่สุดและค่าเฉลี่ยใน result_summary
result_summary = {
    "Best_Parameters": [str(best_params)], # Best hyperparameters

```

```

    "Average_Train_Score":
[results.filter(like="mean_train_score").mean(axis=0).iloc[0]],
    "Average_Validation_Score":
[results.filter(like="mean_test_score").mean(axis=0).iloc[0]],
    "Average_Fold_Time": [results["avg_fold_time"].mean()]
}
print("Detailed results saved to 'multioutput_results.csv' and 'detailed_results.csv'.")

# ส่วนที่เพิ่ม: คำนวณ RMSE, MAE, และ SMAPE
def calculate_smape(y_true, y_pred):
    return 100 * np.mean(2 * np.abs(y_pred - y_true) / (np.abs(y_true) +
np.abs(y_pred)))

# ใช้พารามิเตอร์ที่ดีที่สุดจาก GridSearchCV
cleaned_best_params = {key.replace("estimator__", ""): value for key, value in
best_params.items()}

# ใช้พารามิเตอร์ที่ดีที่สุดจาก GridSearchCV
rf_best = MultiOutputRegressor(RandomForestRegressor(random_state=42,
**cleaned_best_params))

# ฟิตโมเดลด้วยข้อมูล train
rf_best.fit(X_train, y_train)

# ทำนายค่าจากข้อมูล test
y_pred = rf_best.predict(X_test)

# คำนวณค่าตัวชี้วัดแยกตาม y1 และ y2
metrics = {}
for i, target_name in enumerate(y.columns):
    y_true = y_test.iloc[:, i]

```

```

y_pred_target = y_pred[:, i]

metrics[f"MAE_{target_name}"] = mean_absolute_error(y_true, y_pred_target)
metrics[f"RMSE_{target_name}"] = mean_squared_error(y_true, y_pred_target) **
0.5
metrics[f"SMAPE_{target_name}"] = calculate_smape(y_true.values,
y_pred_target)

# สร้าง DataFrame สำหรับผลลัพธ์ทั้งหมดในแถวเดียว
result_row = {"Best_Parameters": str(best_params)}
result_row.update(metrics)
df_metrics = pd.DataFrame([result_row])
print("Test metrics saved to 'model_test_metrics.csv'.")

y_pred_df = pd.DataFrame(y_pred, columns=["Predicted_y1", "Predicted_y2"],
index=y_test.index)

# สร้างคอลัมน์ใหม่สำหรับค่าที่ปรับตามเงื่อนไข
adjust_predict_y1 = y_pred_df["Predicted_y1"].copy()
adjust_predict_y2 = y_pred_df["Predicted_y2"].copy()

# เงื่อนไข: ถ้า adjust_predict_y1 < 0 หรือ adjust_predict_y2 < 0 ให้ทั้งคู่เป็น 0
adjust_predict_y1 = np.where((adjust_predict_y1 < 0) | (adjust_predict_y2 < 0), 0,
adjust_predict_y1)
adjust_predict_y2 = np.where((adjust_predict_y1 < 0) | (adjust_predict_y2 < 0), 0,
adjust_predict_y2)

# ปัดทศนิยมของ adjust_predict_y1: ถ้ามากกว่า 0.5 ให้ปัดขึ้น ถ้าน้อยกว่า 0.5 ให้ปัดลง
adjust_predict_y1 = np.where(adjust_predict_y1 > 0.5, np.ceil(adjust_predict_y1),
np.floor(adjust_predict_y1))

```

```

# เงื่อนไข: ถ้า adjust_predict_y1 == 0 ให้ adjust_predict_y2 = 0
adjust_predict_y2 = np.where(adjust_predict_y1 == 0, 0, adjust_predict_y2)

# เงื่อนไข: ถ้า adjust_predict_y2 == 0 ให้ adjust_predict_y1 = 0
adjust_predict_y1 = np.where(adjust_predict_y2 == 0, 0, adjust_predict_y1)

# เพิ่มคอลัมน์ใหม่เข้าไปใน DataFrame
y_pred_df["adjust_predict_y1"] = adjust_predict_y1
y_pred_df["adjust_predict_y2"] = adjust_predict_y2
# รวมข้อมูล X_test, y_test และ y_pred
y_test_with_pred = pd.concat([X_test_df, y_test.reset_index(drop=True),
y_pred_df.reset_index(drop=True)], axis=1)
# เปลี่ยนชื่อคอลัมน์
y_test_with_pred.columns = list(X_test_df.columns) + ["nclaims_md", "cost_md",
"Predicted_y1", "Predicted_y2", "adjust_predict_y1", "adjust_predict_y2"]
# บันทึกไฟล์ CSV
y_test_with_pred.to_csv("A_nes_data_best_parameter_ของแบบจำลอง_XGBoost.csv",
index=False)

```

ง.3 Shapley Additive Explanation (SHAP)

```

# สร้าง SHAP Explainer
explainer = shap.Explainer(rf_best.predict, X_train, feature_names=X_train.columns)
shap_values = explainer(X_test) # คำนวณ SHAP values

# แสดงค่า SHAP value ของแต่ละตัวแปรแยกตามเป้าหมาย
for i, target_name in enumerate(y.columns): # y.columns คือชื่อ Target เช่น Target1,
Target2
    print(f"\nSHAP values summary for {target_name}:")

# สร้าง DataFrame สำหรับ SHAP values

```

```
shap_df = pd.DataFrame(shap_values[..., i].values, columns=X_train.columns,
index=X_test.index)

# คำนวณค่าเฉลี่ย (Mean Absolute SHAP Value) ของแต่ละตัวแปร
mean_shap_values = shap_df.abs().mean().sort_values(ascending=False)
print(mean_shap_values) # แสดงค่า Mean Absolute SHAP Value ของแต่ละตัวแปร

# แสดง SHAP summary plot และบันทึกผล
plt.figure()

shap.summary_plot(shap_values[..., i], X_test, feature_names=X_train.columns,
show=False) # show=False เพื่อบันทึกก่อนแสดง
plt.title(f"SHAP Summary Plot for {target_name}")
```

ภาคผนวก จ

วิธีการใช้งาน Web Application

จ.1 วิธีการใช้งาน Web Application

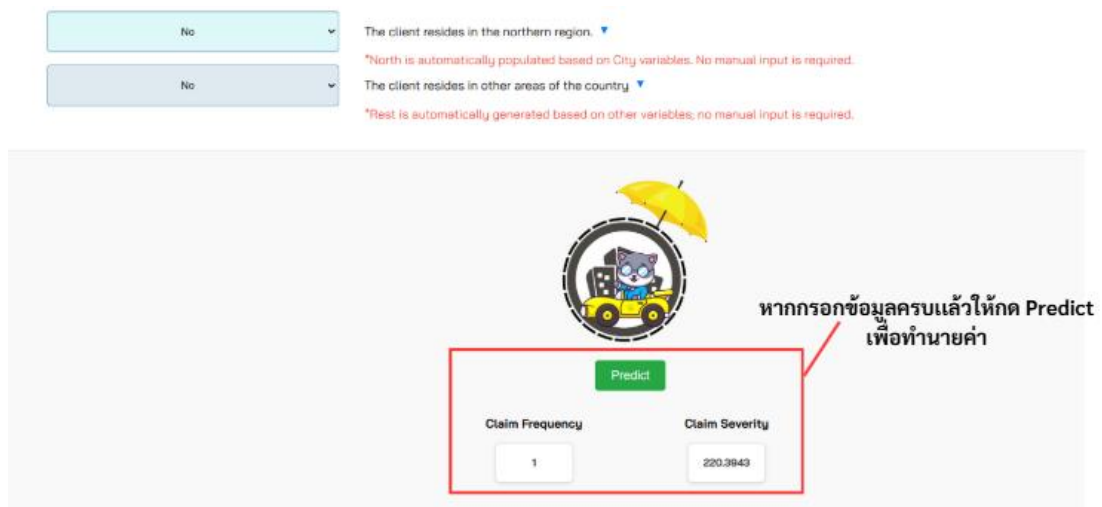


ภาพที่ จ.1 หน้าต่าง Web Application

เมื่อกดคลิกจะมีค่าของตัวแปรให้เลือก

คำอธิบายตัวแปร

ภาพที่ จ.2 หน้าต่าง Web Application ส่วนของการทำนาย (1)



No

The client resides in the northern region.

*North is automatically populated based on City variables. No manual input is required.

No

The client resides in other areas of the country

*Rest is automatically generated based on other variables; no manual input is required.

หากกรอกข้อมูลครบแล้วให้กด Predict เพื่อทำนายค่า

Predict

Claim Frequency

1

Claim Severity

220.3943

ภาพที่ จ.3 หน้าต่าง Web Application ส่วนของการทำนาย (2)

#	date	npol	gender	age	lic_age	client_norther ▲	city	North	rest	claim_frequency	claim_severity
1	16-Mar-2025	15	1	28	14	14	1	0	0	1	156.412
2	16-Mar-2025	16	0	26	11	13	1	0	0	1	220.3943

ภาพที่ จ.4 หน้าต่าง Web Application ส่วนของประวัติการทำนาย

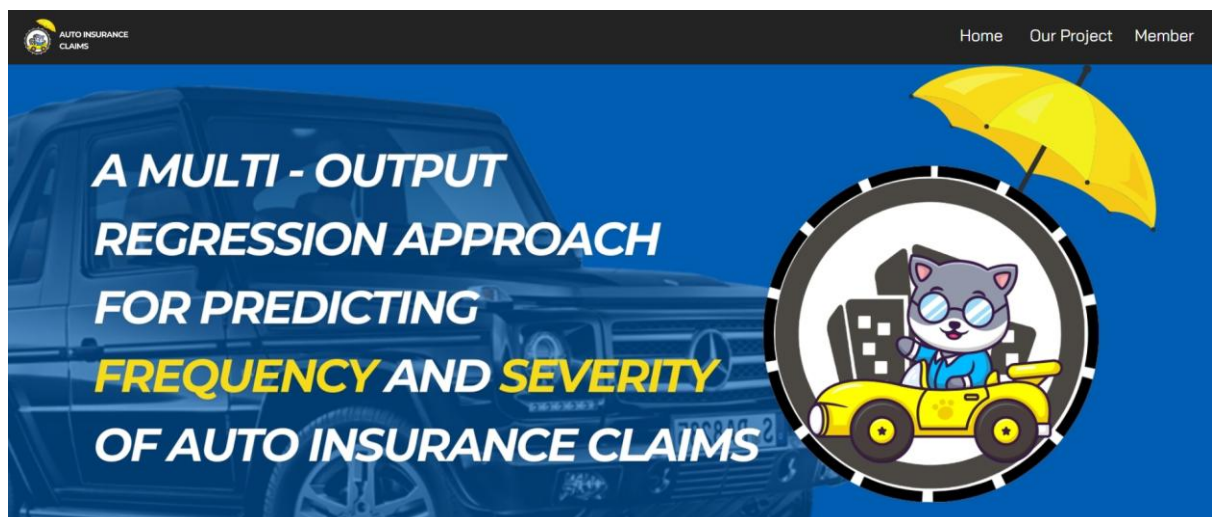
ภาคผนวก ฉ

หน้าต่าง Web Page ที่ผู้วิจัยพัฒนาสำหรับงานวิจัย

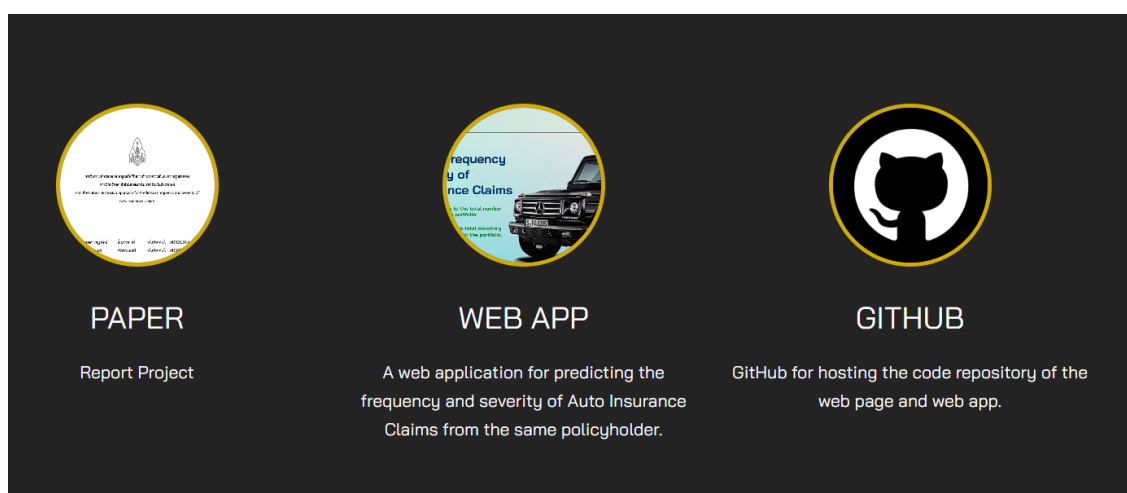
ฉ.1 Web Page ที่ผู้วิจัยพัฒนาสำหรับงานวิจัย



ฉ.2 รายละเอียดของ Web Page



ภาพที่ ฉ.1 หน้าต่าง Web Page (1)



ภาพที่ ฉ.2 หน้าต่าง Web Page (2)

จากภาพที่ ฉ.2 แสดงรายละเอียดงานที่เกี่ยวข้องกับงานวิจัย โดยเมื่อกดที่ภาพจะแสดงรายละเอียดต่างๆ ดังนี้ 1. ภาพ PAPER แสดงรูปเล่มโครงงานวิจัย 2. ภาพ WEB APP แสดง Web Application ที่ใช้ในการทำนายการเรียกร้องค่าสินไหมทดแทน 3. ภาพ GITHUB แสดง Github ที่รวบรวมโค้ดในการพัฒนาเว็บแอปพลิเคชันและเว็บเพจ นอกจากนี้ใน Web Page ยังมีบทความ และ วิดีโอที่แสดงวิธีการใช้งานของ Web Application

ฉ.3 Github ที่รวบรวมโค้ดสำหรับ Web Page & Web Application



ประวัติย่อของผู้วิจัย

ชื่อ นางสาวชฎารัตน์ อิ่มสารพวงค์
 วันเกิด วันที่ 12 กันยายน พ.ศ. 2545
 สถานที่เกิด อำเภอเมือง จังหวัดตราด
 สถานที่อยู่ปัจจุบัน 69 หมู่ 3 ตำบลไม้รูด อำเภอคลองใหญ่ จังหวัดตราด
 ประวัติการศึกษา พ.ศ. 2557 ประถมศึกษาปีที่ 6 โรงเรียนนครศึกษา
 อำเภอคลองใหญ่ จังหวัดตราด
 พ.ศ. 2560 มัธยมศึกษาปีที่ 3 โรงเรียนสตรีประเสริฐศิลป์
 อำเภอเมือง จังหวัดตราด
 พ.ศ. 2563 มัธยมศึกษาปีที่ 6 โรงเรียนสตรีประเสริฐศิลป์
 อำเภอเมือง จังหวัดตราด
 พ.ศ. 2567 ปริญญาวิทยาศาสตรบัณฑิต (วท.บ.)
 อำเภอเมือง จังหวัดขอนแก่น

ชื่อ นายปารเมศ ศิริพรรณนนท์
 วันเกิด วันที่ 30 พฤศจิกายน พ.ศ. 2545
 สถานที่เกิด แขวงศิริราช กรุงเทพมหานคร
 สถานที่อยู่ปัจจุบัน 196 หมู่ 6 ตำบลโนนสัง อำเภอโนนสัง จังหวัดหนองบัวลำภู
 ประวัติการศึกษา พ.ศ. 2557 ประถมศึกษาปีที่ 6 โรงเรียนอนุบาลสุดา
 อำเภอเมือง จังหวัดหนองบัวลำภู
 พ.ศ. 2560 มัธยมศึกษาปีที่ 3 โรงเรียนหนองบัวพิทยาคาร
 อำเภอเมือง จังหวัดหนองบัวลำภู
 พ.ศ. 2563 มัธยมศึกษาปีที่ 6 โรงเรียนหนองบัวพิทยาคาร
 อำเภอเมือง จังหวัดหนองบัวลำภู
 พ.ศ. 2567 ปริญญาวิทยาศาสตรบัณฑิต (วท.บ.)
 อำเภอเมือง จังหวัดขอนแก่น