

ระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง

พรทิwa วิศิษฏ์สรอรรถ

สารนิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรวิศวกรรมศาสตรมหาบัณฑิต

สาขาวิชาวิศวกรรมข้อมูลขนาดใหญ่

วิทยาลัยนวัตกรรมด้านเทคโนโลยีและวิศวกรรมศาสตร์

มหาวิทยาลัยธุรกิจบัณฑิต

ปีการศึกษา 2564

A WEB-BASED SYSTEM FOR TIME SERIES FORECAST

PORNTIVA VISITSORA-AT

**An Independent Study Submitted in Partial Fulfillment of the
Requirements for the Degree of Master of Big Data Engineering,
College of Innovative Technology and Engineering,
Dhurakij Pundit University
Academic Year 2021**

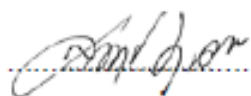


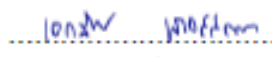
ใบรับรองงานสารนิพนธ์

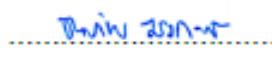
วิทยาลัยนวัตกรรมการด้านเทคโนโลยีและวิศวกรรมศาสตร์ มหาวิทยาลัยธุรกิจบัณฑิต
ปริญญา วิศวกรรมศาสตรมหาบัณฑิต

หัวข้อสารนิพนธ์ ระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง
เสนอโดย พรทิพา วิเศษสุวรรรณ
สาขาวิชา วิศวกรรมข้อมูลขนาดใหญ่
อาจารย์ที่ปรึกษาสารนิพนธ์ ดร.เอกสิทธิ์ พัทธวงค์ศักดิ์

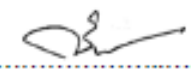
ได้พิจารณาเห็นชอบ โดยคณะกรรมการสอบสารนิพนธ์แล้ว

.....ประธานกรรมการ
(ดร.สรพรพฤทธิ มฤคทัต)

.....กรรมการและอาจารย์ที่ปรึกษา
(ดร.เอกสิทธิ์ พัทธวงค์ศักดิ์)

.....กรรมการ
(ดร.ธนภัทร จังคะจิตร)

วิทยาลัยนวัตกรรมการด้านเทคโนโลยีและวิศวกรรมศาสตร์รับรองแล้ว

.....
(ดร.ชัยพร เขมะภาคะพันธ์)

คณบดีวิทยาลัยนวัตกรรมการด้านเทคโนโลยีและวิศวกรรมศาสตร์

วันที่ 30 เดือน พฤศจิกายน พ.ศ. 2564

| | |
|------------------|---|
| หัวข้อสารนิพนธ์ | ระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง |
| ชื่อผู้เขียน | พรทิชา วิศิษฐ์สรอรรถ |
| อาจารย์ที่ปรึกษา | ดร.เอกสิทธิ์ พัทธวงษ์ศักดิ์ |
| สาขาวิชา | วิศวกรรมข้อมูลขนาดใหญ่ |
| ปีการศึกษา | 2564 |

บทคัดย่อ

การวิจัยนี้เป็นการศึกษามีวัตถุประสงค์เพื่อพัฒนาระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง โดยเปรียบเทียบวิธีพยากรณ์ออกเป็น 5 เทคนิคประกอบไปด้วยเทคนิค AutoRegressive Integrated Moving Average (ARIMA), Seasonal AutoRegressive Integrated Moving Average (SARIMA), วิเคราะห์การถดถอยพหุคูณ (Multiple Linear Regression), ต้นไม้ตัดสินใจ (Decision Tree) และ ป่าสุ่ม (Random Forest) ซึ่งพัฒนาขึ้นมาเป็นรูปแบบของเว็บแอปพลิเคชันที่ช่วยพยากรณ์ข้อมูลในอนาคตได้หลากหลายรูปแบบสำหรับงานวิจัยนี้ได้นำข้อมูลจากกรณีศึกษาต่าง ๆ มาประยุกต์ใช้สร้างแบบจำลองดังนี้ 1. ข้อมูลยอดขายรายวันของร้านกาแฟ 2. ข้อมูลจำนวนผู้ใช้บริการโครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม 3. ข้อมูลจำนวนผู้ป่วยและผู้เสียชีวิตรายใหม่จากสถานการณ์ COVID-19 ในประเทศไทย เมื่อแบ่งข้อมูลเพื่อสร้างแบบจำลองพยากรณ์ยอดขาย, จำนวนผู้ใช้บริการและจำนวนยอดผู้เสียชีวิตรายใหม่โดยใช้การคัดเลือกเทคนิคที่เหมาะสมที่สุดสำหรับการพยากรณ์โดยใช้เกณฑ์ค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน (Mean Absolute Percentage Error, MAPE) ต่ำที่สุดผลวิจัยพบว่าหากข้อมูลที่นำมาสร้างแบบจำลองมีอิทธิพลต่อฤดูกาลเทคนิค SARIMA ให้ผลการพยากรณ์ที่ดีกว่าเทคนิค ARIMA เช่น กรณีศึกษาการพยากรณ์จำนวนผู้ใช้บริการรถไฟฟ้ามหานครสายฉลองรัชธรรมรายเดือน ให้ค่า MAPE สำหรับเทคนิค SARIMA, ARIMA เท่ากับ 25.90% , 35.15% ตามลำดับ ส่วนการเปรียบเทียบเทคนิคระหว่าง Decision Tree, Random Forest และ Multiple Linear Regression สำหรับการพยากรณ์ยอดขายรวมรายวันแยกตามรายสินค้า และจำนวนผู้เสียชีวิตรายใหม่จากสถานการณ์ COVID-19 ในประเทศไทยพบว่าเทคนิคให้ค่า MAPE เฉลี่ยต่ำที่สุดคือเทคนิค Random Forest 18.74% - 54.15% เมื่อเทียบกับเทคนิค Decision Tree 22.45% - 57.34% และ Multiple Linear Regression 28.08% - 156.86% ตามลำดับ

| | |
|------------------------------|---|
| An Independent Study Title | A WEB-BASED SYSTEM FOR TIME SERIES FORECAST |
| Author | Porntiva Visitsora-at |
| An Independent Study Advisor | Dr. Eakasit Pacharawongsakda |
| Department | Big Data Engineering |
| Academic Year | 2021 |

ABSTRACT

This research was conducted to develop a web-based system for time series analysis using statistic and machine learning techniques. The five compared models of time series analysis included AutoRegressive Integrated Moving Average (ARIMA) and Seasonal AutoRegressive Integrated Moving Average (SARIMA), Multiple Linear Regression, Decision Tree, and Random Forest. The system was developed as a web application which could forecast future information in various datasets. In this study, three datasets were utilized for model development: 1. data on daily sales of a coffee shop, 2. data on passengers of the M.R.T Chalong Ratchadham Line, and 3. data on new cases and deaths of COVID-19 in Thailand. Having separated data for modeling sales forecast, the number of passengers, and the number of new COVID-19 deaths in Thailand, through selection of the most appropriate forecasting techniques based on the lowest Mean Absolute Percentage Error (MAPE), the results showed that when the data used for model development influenced SARIMA, it yielded better forecast in comparison to ARIMA. For example, a study forecasting the number of passengers of the M.R.T Chalong Ratchadham Line each month had MAPE of SARIMA, ARIMA representing 25.90%, 35.15% respectively. Moreover, the comparison between Decision Tree, Random Forest, and Multiple Linear Regression for forecasting the daily sales by product and new cases and deaths of COVID-19 in Thailand indicated that Random Forest had the lowest average of MAPE which was 18.74% - 54.15% compared to Decision Tree which was 22.45% - 57.34% and Multiple Linear Regression which was 28.08% - 156.86%, respectively.

กิตติกรรมประกาศ

สารนิพนธ์ฉบับนี้สำเร็จสมบูรณ์ลุล่วงไปได้ด้วยดีเพราะได้รับความกรุณาชี้แนะและช่วยเหลืออย่างดียิ่งจาก ดร.เอกสิทธิ์ พัทธวงษ์ศักดิ์ดาซึ่งเป็นอาจารย์ที่ปรึกษาสารนิพนธ์ ที่ได้กรุณาให้คำแนะนำ ตรวจสอบ และแก้ไขข้อบกพร่องต่าง ๆ มาโดยตลอดเพื่อให้สารนิพนธ์ฉบับนี้สมบูรณ์ ผู้เขียนจึงขอกราบขอบพระคุณไว้ ณ โอกาสนี้

ผู้เขียนขอกราบขอบพระคุณผู้ช่วยศาสตราจารย์ ดร.สรรพทุทธิ มฤคทัต ที่กรุณาให้เกียรติเป็นประธาน โดยมี ดร.ชนภัทร ชังคะจิตรเป็นกรรมการในการสอบสารนิพนธ์ ซึ่งได้กรุณาตรวจแก้ไขสารนิพนธ์ฉบับนี้ให้ถูกต้องสมบูรณ์ยิ่งขึ้นและ นางสาวกุลธิดา รอดบุญ รวมถึงเจ้าหน้าที่บัณฑิตมหาวิทยาลัยธุรกิจบัณฑิตทุกท่านที่ให้ความสะดวกด้านอำนวยความสะดวกและประสานงาน ในการทำสารนิพนธ์ให้กับผู้เขียน ทำให้การจัดทำสารนิพนธ์ของผู้เขียนในครั้งนี้สำเร็จลุล่วงไปด้วยดี

ผู้เขียนขอขอบพระคุณนางสาวพรทิพย์ วิศิษฎ์สรอรรถ ที่กรุณาให้เก็บรวบรวมข้อมูลของร้านกาแฟเพื่อนำมาใช้ทดลองจริงจนทำให้สารนิพนธ์ฉบับนี้สำเร็จลุล่วงด้วยดี

ขอขอบพระคุณบิดามารดาที่สนับสนุนและให้กำลังใจในงานสารนิพนธ์สำเร็จด้วยดี ผู้เขียนขอโน้มบงกชพระคุณบิดามารดาและบูรพาจารย์ทุกท่านที่ได้อบรมสั่งสอนวิชาความรู้ และให้ความเมตตาผู้เขียนมาโดยตลอดและเป็นเป็นกำลังใจสำคัญที่ทำให้สารนิพนธ์สำเร็จลุล่วงไปด้วยดี ขอขอบคุณเพื่อน ๆ BD4 ทุกท่านที่คอยเป็นกำลังใจและช่วยเหลือในงานสารนิพนธ์สำเร็จด้วยดี

พรทิวา วิศิษฎ์สรอรรถ

สารบัญ

| | หน้า |
|---|-------------|
| บทคัดย่อภาษาไทย..... | ฅ |
| บทคัดย่อภาษาอังกฤษ | ง |
| กิตติกรรมประกาศ..... | จ |
| สารบัญตาราง..... | ช |
| สารบัญภาพ | ฌ |
| บทที่ | |
| 1. บทนำ..... | 1 |
| 1.1 ที่มาและความสำคัญของปัญหา..... | 1 |
| 1.2 วัตถุประสงค์ของงานวิจัย..... | 2 |
| 1.3 ขอบเขตงานวิจัย..... | 2 |
| 1.4 ประโยชน์ที่คาดว่าจะได้รับ..... | 3 |
| 2. แนวคิด ทฤษฎี และผลงานวิจัยที่เกี่ยวข้อง..... | 4 |
| 2.1 ความหมาย ประโยชน์และความสำคัญของการพยากรณ์..... | 4 |
| 2.2 องค์ประกอบ ประเภทและคุณลักษณะของการพยากรณ์..... | 4 |
| 2.3 เทคนิคการวิเคราะห์การถดถอย..... | 6 |
| 2.4 การพยากรณ์ในรูปแบบอนุกรมเวลา..... | 8 |
| 2.5 เปรียบเทียบค่าความคลาดเคลื่อนในการประเมินผลการพยากรณ์..... | 12 |
| 2.6 งานวิจัยที่เกี่ยวข้อง..... | 12 |
| 3. ระเบียบวิธีวิจัย..... | 14 |
| 3.1 แนวทางการวิจัย..... | 14 |
| 3.2 เครื่องมือที่ใช้ในการวิจัย..... | 29 |
| 4. ผลการศึกษา | 31 |
| 4.1 ผลการวัดประสิทธิภาพความถูกต้องของแบบจำลอง..... | 31 |
| 4.2 การวัดผลประสิทธิภาพความถูกต้องจากการนำแบบจำลองพยากรณ์ในอนาคต เทียบกับข้อมูลจริง..... | 33 |

สารบัญ (ต่อ)

| บทที่ | หน้า |
|-----------------------------|------|
| 5. บทสรุปและข้อเสนอแนะ..... | 42 |
| 5.1 สรุปผลการวิจัย..... | 42 |
| 5.2 อภิปรายผลการวิจัย..... | 43 |
| 5.3 ข้อเสนอแนะ..... | 43 |
| บรรณานุกรม..... | 45 |
| ภาคผนวก..... | 49 |
| ก | 50 |
| ข | 56 |
| ประวัติผู้เขียน | 74 |

สารบัญตาราง

| ตารางที่ | หน้า |
|---|------|
| 3.1 ตัวอย่างการแปลงแอตทริบิวต์วันหยุด (Holiday), การจัด โปรโมชัน (Promotion) | 20 |
| 3.2 ตัวอย่างการแปลงแอตทริบิวต์วันที่ (Date)..... | 20 |
| 3.3 ตัวอย่างการจัดกลุ่มซื้อสินค้าก่อนนำเข้าวิเคราะห์..... | 21 |
| 3.4 ตัวอย่างค่าพยากรณ์ที่ได้จากเทคนิค Random Forest (ข้อมูลทดสอบ แบบจำลอง)..... | 27 |
| 4.1 เปรียบเทียบค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของการพยากรณ์ด้วยกรณีศึกษาและเทคนิคในรูปแบบต่าง ๆ..... | 32 |
| 4.2 ตัวอย่างข้อมูลพยากรณ์ยอดขายรวมรายวันด้วยเทคนิค ARIMA เมื่อกำหนดคาบเวลาเท่ากับ 7 สำหรับ Export ไฟล์ในรูปแบบ.xlsx..... | 34 |
| 4.3 ตัวอย่างข้อมูล Test สำหรับการพยากรณ์ข้อมูลใหม่ ด้วยเทคนิค Random Forest..... | 35 |
| 4.4 ตัวอย่างข้อมูลพยากรณ์จำนวนผู้ใช้บริการรายเดือนด้วยเทคนิค SARIMA เมื่อกำหนดคาบเวลาเท่ากับ 7 สำหรับ Export ไฟล์ในรูปแบบ.xlsx..... | 38 |
| 4.5 ตัวอย่างข้อมูล Test สำหรับการพยากรณ์ข้อมูลใหม่ ด้วยเทคนิค Random Forest..... | 39 |

สารบัญภาพ

| ภาพที่ | หน้า |
|---|------|
| 2.1 ลักษณะของต้นไม้ตัดสินใจในรูปแบบของ ต้นไม้ไบนารี (Binary tree)..... | 8 |
| 3.1 กระบวนการวิเคราะห์ข้อมูลด้วย CRISP-DM..... | 14 |
| 3.2 Business Canvas ของร้านกาแฟที่ทำการวิจัย..... | 15 |
| 3.3 ตัวอย่างข้อมูลยอดขายสินค้ารายวันที่จะนำมาวิเคราะห์..... | 17 |
| 3.4 ตัวอย่างข้อมูลจำนวนผู้โดยสารโครงการรถไฟฟ้ามหานครสายเฉลิมรัชมงคล ที่นำมาวิเคราะห์..... | 17 |
| 3.5 ตัวอย่างข้อมูลจำนวนป่วย COVID-19 ระลอก 3 ที่จะนำมาวิเคราะห์..... | 18 |
| 3.6 การคัดเลือกข้อมูลสำหรับการสร้างแบบจำลองสำหรับกรณีศึกษาและเทคนิค ต่าง ๆ..... | 19 |
| 3.7 คอลัมภ์ในข้อมูลไฟล์ train, test สำหรับการสร้างแบบจำลองสำหรับ กรณีศึกษาและเทคนิคต่าง ๆ..... | 22 |
| 3.8 รูปแบบของข้อมูลจริงก่อนนำเข้าสู่สร้างแบบจำลองด้วยเทคนิค ARIMA, SARIMA กรณีศึกษาข้อมูลประเภทยอดขายของร้านกาแฟ..... | 23 |
| 3.9 รูปแบบของข้อมูลจริงก่อนนำเข้าสู่สร้างแบบจำลองด้วยเทคนิค ARIMA, SARIMA กรณีศึกษาข้อมูลจำนวนผู้ใช้บริการโครงการรถไฟฟ้ามหานคร สาย ฉลองรัชมงคล..... | 23 |
| 3.10 รูปแบบของข้อมูลจริงก่อนนำเข้าสู่สร้างแบบจำลองด้วยเทคนิค Multiple Linear Regression, Decision Tree, Random Forest กรณีศึกษาข้อมูลประเภท ยอดขายของร้านกาแฟแบ่งตามรหัสสินค้า..... | 23 |
| 3.11 รูปแบบของข้อมูลจริงก่อนนำเข้าสู่สร้างแบบจำลองด้วยเทคนิค Multiple Linear Regression, Decision Tree, Random Forest กรณีศึกษาข้อมูลรายงาน สถานการณ์ COVID-19 ระลอก 3..... | 24 |
| 3.12 ขั้นตอนการทำงานของระบบเพื่อสร้างแบบจำลอง..... | 24 |

สารบัญภาพ (ต่อ)

| ภาพที่ | หน้า |
|--|------|
| 3.13 ตัวอย่างการหารูปแบบ ARIMA (p,d,q) จากการพิจารณาค่า AIC..... | 25 |
| 3.14 ตัวอย่างการหารูปแบบ ARIMA(p,d,q) SARIMA (P,D,Q)T หรือ ARIMA(p,d,q)x(P,D,Q)s จากการพิจารณาค่า AIC..... | 26 |
| 3.15 กราฟเทียบข้อมูลพยากรณ์กับข้อมูลจริง (ข้อมูลทดสอบแบบจำลอง) ของ เทคนิค ARIMA..... | 26 |
| 3.16 ขั้นตอนการทำงานของระบบเพื่อพยากรณ์ข้อมูลใหม่..... | 28 |
| 4.1 กราฟแสดงผลการพยากรณ์ข้อมูลยอดขายรวมรายวัน (1-7 กันยายน 2564) ด้วยเทคนิค ARIMA..... | 34 |
| 4.2 ข้อมูลจริงของยอดขายรวมรายวัน วันที่ 1-7 กันยายน 2564..... | 35 |
| 4.3 รูปแบบของข้อมูลจริงก่อนนำเข้าพยากรณ์ด้วย Random Forest กรณีศึกษา ข้อมูลประเภทยอดขายของร้านกาแฟแบ่งตามรหัสสินค้า..... | 36 |
| 4.4 ผลการพยากรณ์ข้อมูลยอดขายรวมรายรหัสสินค้าของวันที่ 1 กันยายน พ.ศ. 2564 ด้วยเทคนิค Random Forest สำหรับ Export ไฟล์ในรูปแบบ.xlsx..... | 36 |
| 4.5 ข้อมูลจริงของยอดขายรวมรายรหัสสินค้าวันที่ 1 กันยายน 2564..... | 36 |
| 4.6 กราฟแสดงผลการพยากรณ์ข้อมูลจำนวนผู้ใช้บริการ (เดือน มกราคม ถึง เดือน กรกฎาคม พ.ศ. 2564) ด้วยเทคนิค SARIMA..... | 37 |
| 4.7 ข้อมูลจริงของจำนวนผู้ใช้บริการรายเดือน ตั้งแต่เดือน มกราคม ถึง กรกฎาคม 2564..... | 38 |
| 4.8 รูปแบบของข้อมูลจริงก่อนนำเข้าพยากรณ์ด้วย Random Forest กรณีศึกษา ข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3..... | 39 |
| 4.9 ผลการพยากรณ์ข้อมูลผู้เสียชีวิตรายใหม่จากสถานการณ์ COVID-19 ระลอก 3 ด้วยเทคนิค Random Forest สำหรับ Export ไฟล์ในรูปแบบ.xlsx..... | 40 |
| 4.10 ข้อมูลจริงของเคส 30 วันย้อนหลังของรายงานสถานการณ์ COVID-19..... | 40 |

บทที่ 1

บทนำ

1.1 ที่มาและความสำคัญของปัญหา

แนวคิดอนุกรมเวลาปรากฏเป็นพื้นฐานในแทบทุกกิจกรรมมีการประยุกต์ใช้การวิเคราะห์อนุกรมเวลาในทางปฏิบัติหลายอย่าง เช่น การคาดการณ์ทางเศรษฐกิจ การวิเคราะห์ตลาดหุ้น การพยากรณ์แผ่นดินไหว หรือการคาดการณ์ผลตอบแทนและอื่นๆ นอกจากนี้การวิเคราะห์อนุกรมเวลาช่วยให้เราเข้าใจว่าอะไรคือแรงผลักดันสำคัญที่นำไปสู่แนวโน้มเฉพาะในจุดข้อมูลอนุกรมเวลา

งานวิจัยนี้จึงมีวัตถุประสงค์เพื่อพัฒนาระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง โดยแบ่งเป็น 1. Autoregressive Integrated Moving Average (ARIMA) 2. Seasonal Autoregressive Integrated Moving Average (SARIMA) ประเภทแบบจำลองอนุกรมเวลา เพื่อจัดการกับข้อมูลที่มีลักษณะเกี่ยวข้องกันทางด้านเวลานั้นคือข้อมูลในอดีตสามารถส่งผลต่อข้อมูลในอนาคตได้ 3. เทคนิคการวิเคราะห์การถดถอยพหุคูณ (Multiple Linear Regression) เป็นการวิเคราะห์ความสัมพันธ์ในรูปแบบเชิงเส้นตรง (Linearity) ระหว่างตัวแปรอิสระ (Independent Variable) ที่มีมากกว่า 1 ตัวและตัวแปรตาม (Dependent Variable) 4. เทคนิคต้นไม้ตัดสินใจ (Decision Tree) และ 5. เทคนิคป่าสุ่ม (Random Forest) มาพยากรณ์ข้อมูลซึ่งหลักๆที่เลือกใช้หลายเทคนิคจุดประสงค์เพื่อเปรียบเทียบประสิทธิภาพของแต่ละเทคนิคว่าเทคนิคใดให้ความคาดเคลื่อนในการพยากรณ์ต่ำที่สุดเพื่อนำไปใช้จริงสำหรับข้อมูลชุดใหม่ๆที่ต้องการพยากรณ์ได้

ผลที่ได้จากงานวิจัยนี้หลังจากที่นำข้อมูลหลากหลายประเภท เช่น ข้อมูลยอดขาย ข้อมูลจำนวนผู้ใช้บริการรถไฟฟ้า รวมถึงข้อมูลยอดผู้ติดเชื้อ COVID-19 เข้าสู่ระบบเพื่อทำการประมวลผลสร้างแบบจำลองด้วยเทคนิคต่างๆข้างต้นจะแสดงผลเปรียบเทียบค่าความคาดเคลื่อนทำให้ผู้ใช้งานสามารถเลือกใช้แบบจำลองในการพยากรณ์ข้อมูลในอนาคตที่มีความแม่นยำ งานวิจัยนี้แบ่งข้อมูลเพื่อสร้างแบบจำลองที่แตกกันคือ

แบบที่ 1 คือ ชุดข้อมูลอนุกรมเวลาที่มีตัวแปรสร้างแบบจำลอง คือ เวลาและค่าที่ต้องการพยากรณ์เพื่อใช้กับเทคนิค ARIMA, SARIMA

แบบที่ 2 คือ ชุดข้อมูลอนุกรมเวลาที่มีตัวแปรสร้างแบบจำลอง คือ เวลา ตัวแปรอิสระรวม และค่าที่ต้องการพยากรณ์เพื่อใช้กับเทคนิคการวิเคราะห์การถดถอยพหุคูณ (Multiple Linear Regression) ต้นไม้ตัดสินใจ (Decision Tree) และ ป่าสุ่ม (Random Forest)

ใช้ตัวชี้วัดการประเมินผลด้วยค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน (Mean Absolute Percentage Error, MAPE)

1.2 วัตถุประสงค์ของการวิจัย

พัฒนาเว็บแอปพลิเคชันเพื่อนำเสนอรูปแบบเปรียบเทียบการพยากรณ์ข้อมูลด้วยเทคนิค Autoregressive Integrated Moving Average (ARIMA) Seasonal Autoregressive Integrated Moving Average (SARIMA) เทคนิคการวิเคราะห์การถดถอยพหุคูณ (Multiple Linear Regression) เทคนิคต้นไม้ตัดสินใจ (Decision Tree) และวิธีป่าสุ่ม (Random Forest) โดยมุ่งเน้นให้แบบจำลองจากเทคนิคต่างๆข้างต้นสามารถวิเคราะห์คุณลักษณะสำคัญจากอนุกรมข้อมูลที่นำมาใช้และพิจารณาความสัมพันธ์ของแต่ละคุณลักษณะของข้อมูลในอดีตมาเพิ่มประสิทธิภาพการพยากรณ์ให้มีความแม่นยำและความคลาดเคลื่อนที่ต่ำที่สุดเพื่อนำไปใช้ในการพยากรณ์ในอนาคตที่มีความน่าเชื่อถือ

1.3 ขอบเขตการวิจัย

1.3.1 ข้อมูลกรณีศึกษาที่นำมาใช้ในงานวิจัยคือ

1.3.1.1 ข้อมูลยอดขายจริงรายวันของร้านกาแฟตั้งแต่เดือน มีนาคม ถึง สิงหาคม ปี พ.ศ. 2564 รวมไปถึงข้อมูลเพิ่มเติมดังนี้

1.3.1.1.1 ข้อมูลอุณหภูมิ¹ ตั้งแต่เดือน มีนาคม ถึง สิงหาคม ปี พ.ศ. 2564

1.3.1.1.2 ข้อมูลวันหยุดตามประเพณีของสถาบันการเงิน² ปี พ.ศ. 2564

¹Bangkok, Bangkok, Thailand Weather History,

จาก <https://www.wunderground.com/history/monthly/th/bangkok>

²วันหยุดตามประเพณีของสถาบันการเงิน ประจำปี,

จาก <https://www.bot.or.th/Thai/FinancialInstitutions/FIholiday/Pages/2021.aspx>

1.3.1.2 ข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้ามหานคร³ สายฉลองรัชธรรมตั้งแต่เดือน สิงหาคม ปี 2559 ถึง สิงหาคม ปี พ.ศ. 2564

1.3.1.3 ข้อมูลจำนวนผู้ป่วย COVID-19 ระลอก 3⁴ ตั้งแต่เดือน มิถุนายน ถึง ตุลาคม ปี พ.ศ. 2564

1.3.2. ทำการเตรียมข้อมูลให้สมบูรณ์ให้พร้อมสำหรับการนำไปใช้กับแบบจำลองต่าง ๆ

1.3.3. พัฒนาเว็บแอปพลิเคชัน

1.3.3.1 สร้างแบบจำลองพยากรณ์ข้อมูลด้วยการเปรียบเทียบด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่องในรูปแบบต่างๆ

1.3.3.2 ทำการวัดประสิทธิภาพของแต่ละแบบจำลองโดยพิจารณาจากค่าความค่าเคลื่อนด้วยค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน (Mean Absolute Percentage Error, MAPE)

1.3.3.3 สร้างตัวเลือกการพยากรณ์ข้อมูลจากการเลือกแบบจำลองที่พิจารณาประสิทธิภาพเปรียบเทียบก่อนหน้าสำหรับข้อมูลใหม่

1.4 ประโยชน์ที่คาดว่าจะได้รับ

1.4.1 สามารถนำระบบพยากรณ์ข้อมูลที่มีประสิทธิภาพแม่นยำน่าเชื่อถือไปใช้ให้เกิดประสิทธิภาพในการวางแผนด้านการตลาด การจัดสรรทรัพยากร การคาดการณ์เศรษฐกิจ หรือ อื่นๆ ตามจุดประสงค์ของข้อมูลในการใช้งานได้

1.4.2 ระบบพยากรณ์ข้อมูลมีหลากหลายเทคนิคเปรียบเทียบสำหรับสร้างแบบจำลองทำให้ผู้ใช้งานสามารถเลือกใช้เทคนิคการพยากรณ์ให้เหมาะสมกับจุดประสงค์การใช้งานของตนเองได้

1.4.3 ระบบใช้งานง่ายไม่ซับซ้อนและไม่จำเป็นต้องมีความรู้เกี่ยวกับเทคนิคระดับเชี่ยวชาญ ก็สามารถนำระบบนี้ไปใช้งานให้เกิดประโยชน์ได้

³ข้อมูลผู้โดยสาร โครงการรถไฟฟ้ามหานคร สายเฉลิมรัชมงคล และสายฉลองรัชธรรม

จาก <https://data.go.th/dataset/mrta-crmk>

⁴รายงานสถานการณ์ COVID-19 ระลอก 3 (ตั้งแต่ 01/04/2021 –ปัจจุบัน)

จาก <https://covid19.ddc.moph.go.th/api/Cases/timeline-cases-all>

บทที่ 2

แนวคิด ทฤษฎี และผลงานวิจัยที่เกี่ยวข้อง

การศึกษาทฤษฎีที่เกี่ยวข้องกับงานวิจัยนี้ได้ศึกษาแนวคิด ทฤษฎีและผลงานวิจัยที่เกี่ยวข้องเพื่อนำมาประยุกต์ใช้เพื่อเป็นความรู้พื้นฐานและใช้กำหนดแนวทางในการศึกษาวิจัยประกอบไปด้วย 6 หัวข้อหลัก ดังนี้

2.1 ความหมาย ประโยชน์และความสำคัญของการพยากรณ์

2.1.1 ความหมายของการพยากรณ์

การคาดการณ์หรือทำนายสิ่งที่จะเกิดขึ้นในอนาคตโดยอาศัยการวิเคราะห์ข้อมูลในอดีตหรือข้อมูลปัจจุบันสามารถนำไปใช้เพื่อให้ทราบถึงแนวโน้มของการเปลี่ยนแปลงทิศทางการธุรกิจ เศรษฐกิจหรือสภาพแวดล้อมที่จะมีผลกระทบได้ในอนาคตทำให้สามารถที่จะวางแผนหรือกำหนดนโยบายเพื่อให้บรรลุวัตถุประสงค์ของการดำเนินการได้

2.1.2 ประโยชน์และความสำคัญของการพยากรณ์

1. ลดโอกาสที่จะเกิดความเสียหายหรือไม่แน่นอนที่อาจเกิดขึ้นในอนาคต
2. ช่วยในการวางแผนการดำเนินงานรวมทั้งกำหนดเป้าหมาย วัตถุประสงค์ วิสัยทัศน์ แผนกลยุทธ์ และแผนปฏิบัติการได้อย่างเหมาะสม
3. ทำให้สามารถสรรหาทรัพยากรอื่นๆ มาเพิ่มเติมจากพื้นฐานข้อมูลที่มีอยู่ปัจจุบัน

2.2 องค์ประกอบ ประเภทและคุณลักษณะของการพยากรณ์

2.2.1 องค์ประกอบของการพยากรณ์ที่ดี (หทัยชนก นานานอก, 2553)

1. ระบุวัตถุประสงค์ในการนำผลการพยากรณ์ไปใช้และกำหนดช่วงเวลาที่การพยากรณ์จะครอบคลุมถึงเพื่อที่จะสามารถเลือกใช้วิธีการพยากรณ์ได้อย่างถูกต้องเหมาะสม
2. รวบรวมข้อมูลอย่างมีระบบถูกต้องตามความเป็นจริง เพราะคุณภาพของข้อมูลมีผลสำคัญมากต่อการพยากรณ์
3. เมื่อมีสินค้าหลายชนิดควรจำแนกประเภทของสินค้าที่มีลักษณะตามความต้องการที่คล้ายคลึงกันไว้เป็นกลุ่มเดียวกัน โดยแบ่งเป็นพยากรณ์สำหรับกลุ่มแล้วจึงแยกพยากรณ์สำหรับแต่ละสินค้าในกลุ่มอีกครั้งโดยเลือกวิธีการพยากรณ์ที่เหมาะสมกับแต่ละกลุ่มและสินค้า

4. ควรระบุข้อจำกัดและสมมติฐานที่ตั้งไว้ในการพยากรณ์นั้นๆ เพื่อผู้นำผลการพยากรณ์ไปใช้จะได้ทราบถึงเงื่อนไขข้อจำกัดที่มีผลต่อการวิเคราะห์ผลของค่าพยากรณ์

5. หมั่นตรวจสอบความถูกต้องแม่นยำของค่าพยากรณ์ได้กับค่าจริงที่เกิดขึ้นเป็นอย่างไรสม่ำเสมอเพื่อปรับวิธีการที่ใช้ในการคำนวณให้เหมาะสม

2.2.2 ประเภทและคุณลักษณะของการพยากรณ์

แบ่งตามระยะเวลาที่ใช้ในการพยากรณ์

1. การพยากรณ์ระยะสั้น (Short-Range Forecasting)

ระยะเวลาปกติน้อยกว่า 1 ฤดูกาล 1 วัน ไปถึง 1 ปี ลักษณะเป็นการพยากรณ์ลงลึกถึงระดับเฉพาะเพื่อการวางแผนกิจกรรมและการปรับปรุงการจัดซื้อและการจัดการสินค้าคงเหลือเหมาะสมนำไปใช้สำหรับการควบคุมระยะสั้นวิธีพยากรณ์ส่วนใหญ่จะใช้การคาดการณ์แนวโน้มด้วยกราฟการปรับเรียบแบบเอ็กซ์โพเนนเชียลการใช้ดุลยพินิจ

2. การพยากรณ์ระยะปานกลาง (Medium-Range Forecasting)

ระยะเวลาปกติเป็นฤดูกาลจนไปถึง 3 ปี ลักษณะเป็นการพยากรณ์เป็นเชิงตัวเลขมักต้องการความถูกต้องน่าเชื่อถือเหมาะสมนำไปใช้สำหรับการวางแผนโดยรวมวิธีพยากรณ์ส่วนใหญ่คือการรวบรวมข้อคิดเห็นอนุกรมเวลาการวิเคราะห์การถดถอย ดัชนีทางเศรษฐกิจ การใช้ดุลยพินิจ

3. การพยากรณ์ระยะยาว (Long-Range Forecasting)

ระยะเวลาปกติ 5 ปีหรือมากกว่า 3 ปีขึ้นไป ลักษณะเป็นการพยากรณ์กว้างๆและมักเป็นเชิงคุณภาพเหมาะสมนำไปใช้สำหรับการวางแผนธุรกิจวิธีพยากรณ์ส่วนใหญ่จะนำไปใช้เกี่ยวกับด้านเทคโนโลยีภาวะเศรษฐกิจการศึกษาตลาดการใช้ดุลยพินิจ

2.2.3 วิธีการพยากรณ์

2.2.5.1 การพยากรณ์เชิงคุณภาพ (Qualitative Methods หรือ Objective Methods)

เป็นกลุ่มของวิธีการพยากรณ์ที่นำข้อมูลและวิธีการเชิงคุณภาพ พยากรณ์โดยผู้ที่มีประสบการณ์ความรู้ความสามารถโดยไม่ใช้ตัวแบบทางคณิตศาสตร์จึงตรวจสอบความแม่นยำของการพยากรณ์ได้ยากกว่าการพยากรณ์เชิงปริมาณตัวอย่างเทคนิคที่ใช้ในการพยากรณ์เชิงคุณภาพ เช่น การประมาณการ การระดมความคิด การสำรวจตลาด หรือเทคนิคเดลไฟ (Delphi) คือสามารถแสดงความคิดเห็นแบบไม่แสดงตัวตนเมื่อได้คำตอบและได้ข้อสรุปไม่ตรงกันก็จะทำใหม่เรื่อยๆ จนกว่าได้ข้อสรุปที่เป็นหนึ่งเดียว

2.2.5.2 การพยากรณ์เชิงปริมาณ (Quantitative Methods)

เป็นกลุ่มของวิธีการพยากรณ์ที่นำข้อมูล หรือตัวเลขจากอดีต โดยใช้รูปแบบทางคณิตศาสตร์ เพื่อใช้ในการสร้างแบบจำลองและพยากรณ์ไปในอนาคต โดยแบ่งตัวอย่างออกเป็น 2 เทคนิคย่อย

1. การพยากรณ์ความสัมพันธ์ (Casual Forecasting)

เป็นเทคนิคที่ใช้ปัจจัยที่คาดว่าจะมีความสัมพันธ์กับตัวแปรที่จะพยากรณ์ เช่น การพยากรณ์ยอดขายจะพิจารณาหาความสัมพันธ์ระหว่างยอดขายกับค่าโฆษณาจำนวนช่องทางการจัดจำหน่ายรายได้เฉลี่ยของประชากรคุณภาพของสินค้า เป็นต้น การหาความสัมพันธ์ดังกล่าวจะใช้เทคนิคที่เรียกว่าการวิเคราะห์ความถดถอยและสหสัมพันธ์

2. การพยากรณ์อนุกรมเวลา (Time series Forecasting)

เป็นเทคนิคที่ใช้ข้อมูลในอดีตของตัวแปรที่ต้องการพยากรณ์เพื่อพยากรณ์ค่าในอนาคต เช่น ใช้ข้อมูลจำนวนผู้เข้าใช้บริการของปี พ.ศ. 2559-2564 เพื่อพยากรณ์จำนวนผู้เข้าใช้บริการของปี พ.ศ. 2565

2.3 เทคนิคการวิเคราะห์การถดถอย

2.3.1 การวิเคราะห์ความถดถอยพหุคูณ (Multiple Regression Analysis)

เป็นการวิเคราะห์ความสัมพันธ์ระหว่างตัวแปรอิสระตั้งแต่ 2 ตัวขึ้นไปกับผลของตัวแปรตามที่เกิดขึ้น โดยมีรูปแบบสมการพยากรณ์ถดถอยพหุคูณดังนี้

$$\hat{y} = b_0 + b_1x_1 + b_2x_2 + \dots + b_k x_k$$

โดยที่สัญลักษณ์ที่ใช้มีความหมายดังนี้

x_i คือ ค่าของตัวแปรอิสระ

\hat{y} คือ ค่าของตัวแปรตาม

b_0 คือ ค่าคงที่ (Constant) ของสมการถดถอยโดยจะเป็นจุดตัดแกน Y ของสมการ

b_i คือ ค่าสัมประสิทธิ์การถดถอย (Regression Coefficient) ของตัวแปรอิสระ x_i แต่ละตัว คือ จำนวนของตัวแปรอิสระ

(ดร.สุทิน, 2560) กล่าวว่า ในการวิเคราะห์การถดถอยจำเป็นต้องมีข้อตกลง (Assumption) ในการวิเคราะห์การถดถอย ซึ่งมีจำนวนข้อตกลงเบื้องต้นที่สำคัญ เช่น

1. ตัวแปรอิสระและตัวแปรตามต้องเป็นตัวแปรเชิงปริมาณ (Quantitative Variable) หรือ ตัวแปรต่อเนื่อง (Continuous Variable) หรือมีระดับการวัดเป็น Interval หรือ Ratio Scale เช่น ส่วนสูงน้ำหนัก เป็นต้น ในขณะที่ตัวแปรอิสระมีระดับการวัดเป็น Nominal หรือ Ordinal Scale

จะต้องแปลงข้อมูลให้เป็นตัวเลขคือ มีค่า 0 กับ 1 ก่อนจึงจะสามารถนำไปวิเคราะห์แต่ไม่ควรจะมีการแปลงข้อมูลในรูปแบบดังกล่าวหลายตัวเพราะจะทำให้เกิดความคลาดเคลื่อนในการพยากรณ์ได้มากขึ้น

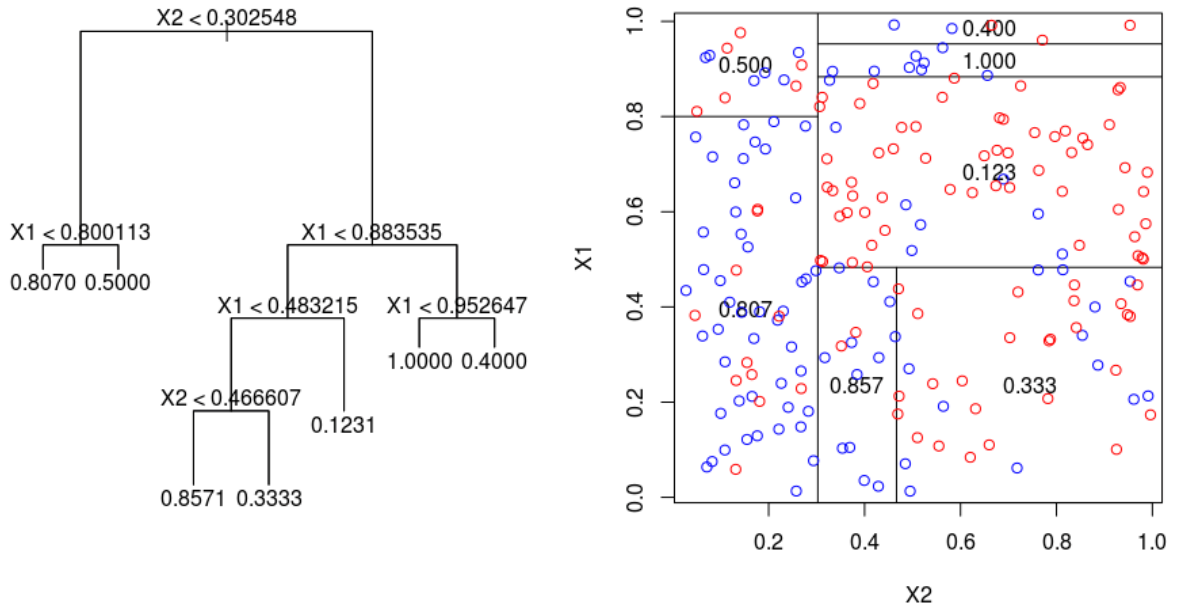
2. ตัวแปรอิสระและตัวแปรตามมีความสัมพันธ์เชิงเส้นตรง
3. ตัวแปรอิสระควรเป็นอิสระต่อกัน (ค่าสหสัมพันธ์ไม่ควรเกิน 0.7) เพราะจะทำให้เกิด Multicollinearity คือ การที่ตัวแปรอิสระบางตัวมีความสัมพันธ์กันเองซึ่งส่งผลกระทบต่อค่าสัมประสิทธิ์การตัดสินใจ (R^2) สูงเกินกว่าความเป็นจริง
4. มีตัวแปรตามมีการแจกแจงแบบปกติ (Normal Distribution)
5. ค่าความแปรปรวนของตัวแปรตามที่คงที่
6. ความแปรปรวนของค่าความคลาดเคลื่อนจากการพยากรณ์ (Residual) ทุกจุดบนเส้นถดถอยมีค่าเท่ากัน

2.3.2 เทคนิคต้นไม้ตัดสินใจและป่าสุ่ม รูปแบบการถดถอย (Decision Tree and Random Forest Regression)

2.3.2.1 ต้นไม้ตัดสินใจ (Decision Tree)

ต้นไม้ตัดสินใจเป็นวิธีการเรียนรู้ของเครื่อง (Machine Learning) รูปแบบหนึ่งที่ได้รับคามนิยมโดยเทคนิคการให้ผลลัพธ์จะออกมาในรูปแบบลักษณะของโครงสร้างต้นไม้เป้าหมายคือการสร้างแบบจำลองที่คาดการณ์ค่าของตัวแปรตามโดยการเรียนรู้กฎการตัดสินใจอย่างง่ายที่อนุมานจากคุณลักษณะข้อมูลโดยจำลองแผนผังการตัดสินใจโดยเลือกแอตทริบิวต์ที่มีอิทธิพลต่อการจำแนกข้อมูลสูงที่สุด (Root node) ไปยังโหนดถัดไปจนกว่าจะไม่มีเงื่อนไขให้ไปต่อ นั่นคือถึงใบ (Leaf node)

อัลกอริทึมที่ใช้สร้างแบบจำลองแผนผังต้นไม้ตัดสินใจมีดังนี้ ID3, C4.5, C5.0 และ CART โดยงานวิจัยนี้จะใช้ไลบรารีของ scikit-learn ที่ใช้ CART (Classification and Regression Trees) (Breiman, Leo; Friedman, J. H.; Olshen, R. A.; Stone, C. J.) (1984) ซึ่งสนับสนุนการพยากรณ์ในรูปแบบการถดถอย (Regression) อัลกอริทึม CART (Classification and Regression Trees) มีลักษณะของต้นไม้ตัดสินใจในรูปแบบของต้นไม้ไบนารี (Binary tree) ซึ่งแต่ละโหนดจะแตกออกเป็น 2 กิ่งเท่านั้นดังภาพที่ 2.1 โดยคำนวณดัชนีจินิ (Gini Index) สำหรับเลือกแอตทริบิวต์ที่จะเป็น โหนดราก (Root node) และใบ (Leaf node)



ภาพที่ 2.1 ลักษณะของต้นไม้ตัดสินใจในรูปแบบของต้นไม้ไบนารี (Binary tree)

ที่มา: <https://gdcoder.com/decision-tree-regressor-explained-in-depth/>

2.3.2.2 ป่าสุ่ม (Random Forest)

(ชนัท, 2561) ป่าสุ่มเป็นวิธีการเรียนรู้ของเครื่อง (Machine Learning) รูปแบบหนึ่งที่ใช้พื้นฐานพัฒนามาจากต้นไม้ตัดสินใจ (Decision Tree) ในทำนายในรูปแบบชุดของต้นไม้ตัดสินใจหลายๆ ต้น (Ensemble of Decision Trees) โดยสร้างจากการสุ่มข้อมูลตัวอย่างแบบเลือกแล้วใส่กลับ (random sampling with replacement) เพื่อนำมาสร้างเป็นแบบจำลองต้นไม้โดยแต่ละต้นไม้มีลักษณะที่ไม่ซ้ำกัน โดยแต่ละแบบจะจำลองจะมีการทำนายผลซึ่งผลจากการทำนายของต้นไม้แต่ละต้นจะทำการโหวตเลือกผลการทำนายที่ได้รับการโหวตมากที่สุดวิธีการนี้เรียกว่า Bagging หรือ Bootstrapping โดยงานวิจัยนี้ใช้ Random forest ของ scikit-learn เช่นเดียวกับเทคนิคต้นไม้ตัดสินใจ

2.4 การพยากรณ์ในรูปแบบอนุกรมเวลา

การวิเคราะห์อนุกรมเวลา (Time series analysis) ซึ่งเป็นวิธีการอาศัยข้อมูลจากอดีตเพื่อการพยากรณ์หรือคาดหมายสิ่งที่เกิดขึ้นในอนาคต

องค์ประกอบของอนุกรมเวลา (Time Series Components) ประกอบด้วย 4 ส่วนดังนี้

1. แนวโน้ม (Trend Component) มีรูปแบบการเพิ่มขึ้นหรือลดลงอย่างต่อเนื่องจะแสดงรูปแบบให้เห็นได้เมื่อมีระยะเวลาหลาย ๆ ปี
2. วัฏจักร (Cyclical Component) มีรูปแบบการเคลื่อนไหวซ้ำ ๆ ของการเพิ่มขึ้นหรือลดลงจะแสดงรูปแบบให้เห็นได้เมื่อมีระยะเวลามากกว่า 1 ปีอาจจะเกิดจากผลกระทบจากองค์ประกอบทางเศรษฐกิจ หรือวัฏจักรของธุรกิจ เป็นต้น
3. ฤดูกาล (Seasonal Component) มีรูปแบบความผันผวนขึ้นลงที่เกิดขึ้นเป็นประจำมีการเปลี่ยนแปลงตามฤดูกาล เป็นต้นจะแสดงรูปแบบให้เห็นเมื่อมีระยะเวลาประมาณ 1 ปี
4. การสุ่ม/ผิดปกติ (Random/Irregular Component) มีรูปแบบที่ไม่มีแบบแผนเป็นผลกระทบจากการเปลี่ยนแปลงแบบสุ่มหรือเหตุการณ์ที่ไม่คาดคิดมาก่อนเกิดขึ้นในช่วงเวลาสั้น ๆ และอาจจะไม่เกิดขึ้นซ้ำอีก

2.4.1 การพยากรณ์ด้วยเทคนิค Box-Jenkins

2.4.1.1 Autoregressive integrated moving average (ARIMA) เป็นเทคนิคพยากรณ์วิเคราะห์ข้อมูลอนุกรมเวลาอาศัยพฤติกรรมของข้อมูลในอดีตเพื่อกำหนดรูปแบบในปัจจุบันและอธิบายแนวโน้มหรือปรากฏการณ์ต่างๆ ของตัวข้อมูลเองในอนาคตโดยการหารูปแบบที่เหมาะสมให้กับข้อมูลเพื่อใช้ในการพยากรณ์โดยจะใช้เกณฑ์ AIC เลือกแบบจำลองที่ให้ค่า AIC น้อยที่สุดเป็นแบบจำลองที่ใช้ของข้อมูลอนุกรมเวลาที่พิจารณา

(ผศ.ดร. เฉลิมพล, 2562) แบบจำลอง ARIMA ประกอบด้วยองค์ประกอบ 3 ส่วน คือ

1. AR(p) คือ Y_t ที่ถูกกำหนดด้วย Autoregressive process

$$\text{โดย } Y_t = Y_{t-1}, Y_{t-2}, \dots, Y_p$$

2. I(d) คือ อันดับความหุคหนึ่งของ Y_t

$$\text{โดย } d = 0, 1, 2, \dots, d$$

3. MA(q) คือ Y_t ที่ถูกกำหนดด้วย Moving average process

$$\text{โดย } Y_t = \epsilon_{t-1}, \epsilon_{t-2}, \dots, \epsilon_q$$

p และ q หมายถึง ลำดับของคาบเวลาในอดีต (Lag length) ที่เหมาะสม

Autoregressive process หรือ AR(p) เป็นกระบวนการของค่าสังเกตค่าหนึ่ง Y_t ที่ถูกกำหนดขึ้นจากความสัมพันธ์ของตัวมันเองในอดีต

$$Y_t = \alpha_0 + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \epsilon_t$$

Y คือ ตัวแปรตาม, α คือ ค่าคงที่, ϕ คือ ค่าสัมประสิทธิ์ของ AR, t คือ เวลา, p คือ ลำดับของคาบเวลาในอดีต

Moving Average process หรือ MA(q) หมายถึง รูปแบบที่แสดงว่าค่า Y_t ถูกกำหนดมาจากค่าความคลาดเคลื่อน (\mathcal{E}) ในอดีต

$$Y_t = \alpha_0 + \mathcal{E}_t - \theta_1 \mathcal{E}_{t-1} - \theta_2 \mathcal{E}_{t-2} - \dots - \theta_q \mathcal{E}_{t-q}$$

Y คือ ตัวแปรตาม, α คือ ค่าคงที่, θ คือ ค่าสัมประสิทธิ์ของ MA, t คือ เวลา, q คือ ลำดับของคาบเวลาในอดีต

Autoregressive and Moving Average Process หรือ ARMA(p,q) เป็นการรวมกันระหว่าง AR กับ MA นั่นคือข้อมูลอนุกรมเวลาขึ้นอยู่กับทั้งค่าของข้อมูลอนุกรมเวลาในอดีตและค่าความคลาดเคลื่อนทั้งในปัจจุบันและในอดีตเป็นวิธีวิเคราะห์อนุกรมเวลาที่อาศัยขบวนการ Stochastic โดยถือว่าข้อมูลที่เกิดขึ้นตามเวลาที่เปลี่ยนแปลงไปมีลักษณะการเกิดที่เป็นไปตามกฎความน่าจะเป็นซึ่งการวิเคราะห์อนุกรมเวลาโดยวิธีนี้ลักษณะของอนุกรมเวลาต้องเป็นอนุกรมเวลาที่มีคุณสมบัติ Stationary

แบบจำลอง ARMA เมื่อตัวแปร Y มีความหยุดนิ่ง ณ ระดับปกติของข้อมูล หรือ $I(0)$, $d = 0$

$$\text{ARMA}(p,d) \text{ คือ } Y_t = \alpha_0 + \varphi_1 Y_{t-1} + \dots + \varphi_p Y_{t-p} + \mathcal{E}_t - \theta_1 \mathcal{E}_{t-1} - \dots - \theta_q \mathcal{E}_{t-q}$$

แบบจำลอง ARIMA เมื่อ ตัวแปร Y มีความหยุดนิ่ง ณ ผลต่างลำดับที่ d หรือ $I(d)$

$$\text{ARIMA}(p,d,q) \text{ คือ } \Delta^d Y_t = \alpha_0 + \varphi_1 \Delta^d Y_{t-1} + \dots + \varphi_p \Delta^d Y_{t-p} + \mathcal{E}_t - \theta_1 \mathcal{E}_{t-1} - \dots - \theta_q \mathcal{E}_{t-q}$$

2.4.1.2 Seasonal Autoregressive integrated moving average (SARIMA) เป็นเทคนิคพยากรณ์วิเคราะห์ข้อมูลอนุกรมเวลาถูกพัฒนาขึ้นเพื่อรองรับการทำงานที่คำนึงถึงความผันแปรตามฤดูกาลซึ่งเป็นส่วนประกอบที่มีความสำคัญดังนั้นรูปแบบจะแสดงเป็น SARIMA(p,d,q)x(P,D,Q) และขั้นตอนการสร้างแบบจำลองพยากรณ์แสดงรายละเอียดดังนี้ (วารสาร คณิตศาสตร์ และ วิทยาศาสตร์ มหิดล, 2556)

$$\varphi_p(B) \Phi_p(B^s)(1-B^s)^D Y_t = \delta + \theta_q(B) \Theta_q(B^s) \mathcal{E}_t$$

เมื่อ $\delta = \mu \varphi_p(B) \Phi_p(B^s)$ แทนค่าคงตัว โดยที่ μ คือค่าเฉลี่ยของอนุกรมเวลาที่คงที่ (Stationary)

$$\varphi_p(B) = 1 - \varphi_2 B^2 - \dots - \varphi_p B^p$$

คือตัวดำเนินการสหสัมพันธ์ในตัวอันดับที่ p กรณีไม่มีฤดูกาล (Non-Seasonal Moving Average Operator of Order q): AR(p))

$$\Phi_p(B^s) = 1 - \Phi_2 B^{2s} - \dots - \Phi_p B^{ps}$$

แทนตัวดำเนินการสหสัมพันธ์ในตัวอันดับที่ P กรณีมีฤดูกาล (Seasonal Moving Average Operator of Order Q): SAR(P))

$$\theta_p(B) = 1 - \theta_1 B - \theta_2 B^2 - \dots - \theta_p B^p$$

แทนตัวดำเนินการเฉลี่ยเคลื่อนที่อันดับที่ q กรณีไม่มีฤดูกาล (Non-Seasonal Moving Average Operator of Order q): $MA(q)$

$$\Theta_p(B^s) = 1 - \Theta_1 B^s - \Theta_2 B^{2s} - \dots - \Theta_Q B^{Qs}$$

แทนตัวดำเนินการเฉลี่ยเคลื่อนที่อันดับที่ Q กรณีมีฤดูกาล (Seasonal Moving Average Operator of Order Q): $SMA(Q)$

d และ D แทนลำดับที่ของการหาผลต่างและผลต่างฤดูกาลตามลำดับ

s แทนจำนวนคาบของฤดูกาล

B แทนตัวดำเนินการถอยหลัง (Backward Operator) โดยที่ $B^s Y_t = Y_{t-s}$

\mathcal{E}_t แทนอนุกรมเวลาของความคลาดเคลื่อนที่มีการแจกแจงปกติและเป็นอิสระกัน

2.4.1.3 การเลือกค่า p, d, q และ P, D, Q ให้เหมาะสมกับแบบจำลอง

หาค่า d ด้วยวิธีการตรวจสอบอันดับความหยุดนิ่ง Augmented Dickey-Fuller (ADF) เพื่อที่จะหลีกเลี่ยงข้อมูลที่มีค่าเฉลี่ย (mean) และความแปรปรวน (variances) ที่ไม่คงที่ในแต่ละช่วงเวลาที่แตกต่างกัน

หาค่า D ด้วยวิธีการตรวจสอบอันดับความหยุดนิ่ง (Canova-Hansen, 1995) ประมาณค่าของความแตกต่างตามฤดูกาล

ส่วนค่า p, q และ P, Q ใช้วิธีขยับตัวเลขไปที่ละขั้นเรียกว่า (Stepwise Hyndman and Khandakar, 2018) หาแบบจำลองโอกาสที่เป็นไปได้ ดังตัวอย่าง

ARIMA(2,d,2) SARIMA(2,d,2)(1,D,1)

ARIMA(0,d,0) SARIMA(0,d,0)(0,D,0)

ARIMA(1,d,0) SARIMA(1,d,0)(1,D,0)

ARIMA(0,d,1) SARIMA(0,d,1)(0,D,1)

พิจารณาแบบจำลองด้วยค่า Akaike information criterion (AIC)

(Akaike, 1973) ได้เสนอเกณฑ์การคัดเลือกแบบจำลองขึ้นมาที่ให้ค่าพยากรณ์แม่นยำที่สุดเป็นเกณฑ์ที่สร้างจากการประมาณค่าความคลาดเคลื่อนรวมเข้ากับข้อสนเทศ (Information) ของค่าสังเกตและใช้แนวคิดจากการหาค่าน้อยสุดของข้อสนเทศด้วยหลักการคูณแบบลิค-ไลเบอร์

(Kullback – Leibler, 1951)

$$AIC = n \cdot \ln\left(\frac{SSE}{n}\right) + 2p$$

n คือ ขนาดตัวอย่าง

\ln คือ ลอการิทึมฐาน e

SSE คือ ค่าความคลาดเคลื่อนกำลังสองของแบบจำลองการถดถอย

p คือ จำนวนพารามิเตอร์ของแบบจำลองการถดถอย

งานวิจัยนี้จะนำแบบจำลอง ARIMA(p,d,q) SARIMA(p,d,q)(P,D,Q) ที่ AIC มีค่าน้อยที่สุดไปใช้สร้างแบบจำลองต่อไป

2.5 เปรียบเทียบค่าความคลาดเคลื่อนในการประเมินผลการพยากรณ์

เทียบความแตกต่างระหว่างข้อมูลจริงกับค่าพยากรณ์จากค่าความคลาดเคลื่อนที่เกิดขึ้น สำหรับงานวิจัยนี้จะเลือกนำตัวชี้วัดในการประเมินคือค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน (Mean Absolute Percentage Error, MAPE)

$$MAPE = \frac{1}{n} \sum_{t=1}^n \left| \frac{Y_t - \hat{Y}_t}{Y_t} \right| \times 100$$

n คือจำนวนตัวอย่าง

Y_t คือข้อมูลจริง ณ เวลา t

\hat{Y}_t คือค่าพยากรณ์ของข้อมูล ณ เวลา t

2.6 งานวิจัยที่เกี่ยวข้อง

นิชา แก้วหาวงษ์ (2547) ในงานวิจัยนี้ได้นำเสนอการพยากรณ์ข้อมูลอนุกรมเวลาโดยวิธีการทำให้เรียบแบบเอ็กซ์โปเนนเชียลและวิธีการของบ็อกซ์-เจนกินส์ของการพยากรณ์มูลค่าการส่งออกข้าว ยางพารา และมันสำปะหลังพบว่าแบบจำลอง ARIMA โดยวิธีการของบ็อกซ์-เจนกินส์เป็นแบบจำลองที่ดีที่สุดโดยให้ความคลาดเคลื่อนต่ำสุด

วรางคณา เรียนสุทธิ์ (2562) ในงานวิจัยนี้ได้นำเสนอเทคนิคพยากรณ์ที่เหมาะสมกับอนุกรมเวลาราคามะพร้าวด้วยวิธีบ็อกซ์-เจนกินส์ของการพยากรณ์ราคามะพร้าวพบว่าวิธีบ็อกซ์-เจนกินส์ที่มีแบบจำลอง AR(2) I(2) MA(2, 13, 15) ไม่มีพจน์ค่าคงตัวเป็นวิธีที่มีความถูกต้องและเหมาะสมมากที่สุด

หทัยชนก นานานอก (2553) ในงานวิจัยนี้ได้นำเสนอการพยากรณ์ยอดขายสินค้าเพื่อการวางแผนผลิตจากการเปรียบเทียบด้วยเทคนิคการพยากรณ์ 5 แบบ คือวิธีการพยากรณ์การหาค่าเฉลี่ยเคลื่อนที่ วิธีเทคนิควินเตอร์ วิธีทำให้เรียบแบบเอ็กซ์โปเนนเชียลชั้นเดียว วิธีทำให้เรียบแบบ

เอกซ์โปเนนเชียลสองชั้นและวิธีวิเคราะห์แนวโน้มเชิงเส้นพบว่าวิธีการพยากรณ์ที่ให้ค่าความคลาดเคลื่อนน้อยที่สุดคือวิธีวิเคราะห์แนวโน้มเชิงเส้น

Abdulwahed Salam, Abdelaziz El Hibaoui (2018) ในงานวิจัยนี้ได้นำเสนอการเปรียบเทียบเทคนิคการเรียนรู้ของเครื่อง (Machine Learning) สำหรับทำนายการใช้พลังงานไฟฟ้าเมืองเตโตวอนด้วยเทคนิคการพยากรณ์ 5 วิธี Linear regression, Decision tree, Random forest, Feedforward Neural network และ Supper vector พบว่าวิธี Random forest เป็นวิธีที่เหมาะสมที่สุด

Balpreet Singh, Pawan Kumar, Nonita Sharma and K P Sharma (2020) ในงานวิจัยนี้ได้นำเสนอการพยากรณ์ยอดขายในอนาคตของ Amazon.com, Inc. ด้วยเทคนิคการพยากรณ์ 4 แบบคือ Holt-Winters Exponential Smoothing, Neural Network Autoregression Model, ARIMA (AutoRegressive Integrated Moving Average) และ SARIMA (Seasonal AutoRegressive Integrated Moving Average) พบว่า SARIMA ให้ผลลัพธ์ที่แม่นยำที่สุดเมื่อเปรียบเทียบกับวิธีอื่นๆ

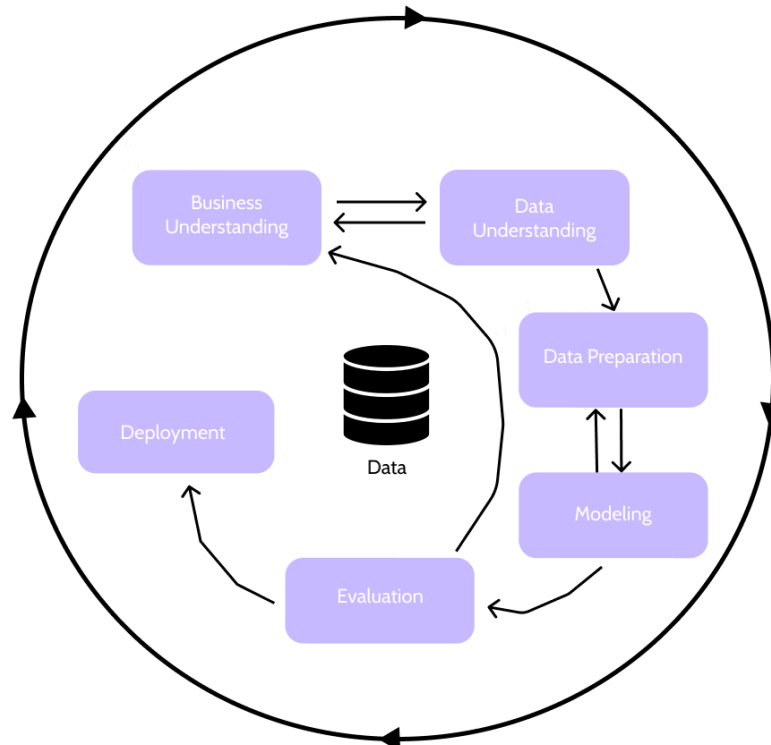
Takashi Tanizaki, Tomohiro Hoshino, Takeshi Shimmura and Takeshi Takenaka (2019) ในงานวิจัยนี้ได้นำเสนอการพยากรณ์จำนวนลูกค้าจะเข้ามาใช้บริการในร้านอาหารจากการเปรียบเทียบด้วยเทคนิคการพยากรณ์ 4 แบบคือ Bayesian Linear Regression, Boosted Decision Tree Regression, Decision Forest Regression and Stepwise method พบว่าแต่ละวิธีการพยากรณ์ให้อัตราการพยากรณ์ (Forecast rate) จากการใช้เครื่องมือสร้างแบบจำลองการเรียนรู้และการวิเคราะห์ทางสถิติ ที่ชื่อว่า Azure Machine Learning และ SPSS ให้ค่าเฉลี่ยเกิน 85% ขึ้นไป

บทที่ 3 ระเบียบวิธีวิจัย

การศึกษาวิจัยครั้งนี้เป็นการวิจัยเชิงพยากรณ์ (Predictive) เป็นการวิจัยเพื่อที่จะนำผลของการพยากรณ์ข้อมูลไปใช้เป็นข้อมูลสำหรับการวางแผนบริหารจัดการทางด้านต่างๆ เพื่อหลีกเลี่ยงความเสี่ยงที่อาจเกิดขึ้นในอนาคตโดยใช้ข้อมูลในรูปแบบของอนุกรมเวลา

3.1 แนวทางการวิจัย

แนวทางการวิจัยนำหลักการ Cross-Industry Standard Process For Data Mining (CRISP-DM) ที่เป็นที่ยอมรับในการทำเหมืองข้อมูลมาประยุกต์ใช้



ภาพที่ 3.1 กระบวนการวิเคราะห์ข้อมูลด้วย CRISP-DM

โดยนำข้อมูลอนุกรมเวลาที่แตกต่างกันมาทำการทดลองใช้กับระบบที่พัฒนาขึ้นมา เป็น 3 ประเภทของข้อมูลดังนี้

1. ข้อมูลยอดขายของร้านกาแฟ
2. ข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม
3. ข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3

3.1.1 การทำความเข้าใจกับธุรกิจ (Business Understanding)

3.1.1.1 กรณีศึกษาเป็นประเภทข้อมูลยอดขายของร้านกาแฟ

กรณีนี้เป็นธุรกิจประเภทร้านค้า SME เราได้ศึกษาและนำโมเดลธุรกิจพื้นฐานคือ Business Model Canvas ที่จะช่วยให้เรามองเห็นรายละเอียดต่างๆของธุรกิจตัวเองได้อย่างดีไม่ว่าจะเป็นทั้งจุดเด่นจุดด้อยข้อได้เปรียบข้อเสียเปรียบกิจกรรมหลักของธุรกิจ หรือแม้กระทั่งลูกค้าและลูกค้าของเราคือใครเพื่อที่จะมองเห็นจุดที่ต้องพัฒนาปรับปรุงให้ดีขึ้นดังภาพที่ 3.2

| | | | | |
|--|---|--|--|---|
| Key Partners (คู่ค้าที่เกี่ยวข้อง) <ul style="list-style-type: none"> - ผู้ขายวัตถุดิบ - โรงเรียนและผู้ฝึกสอนบาริสต้า - เจ้าของพื้นที่ - ผู้ขายเครื่องมือ อุปกรณ์ | Key Activities (กิจกรรมที่ต้องทำเพื่อให้โมเดลธุรกิจอยู่และไปได้) <ul style="list-style-type: none"> - การจัดหาวัตถุดิบ - การพัฒนาเมนูใหม่ๆและพัฒนาสูตรเครื่องดื่มให้ดีขึ้น - พัฒนาฝีมือบาริสต้า - พัฒนาการดูแลลูกค้า | Value Propositions (จุดขาย จุดเด่นที่เป็นคุณค่าของธุรกิจเรา) <ul style="list-style-type: none"> - เมนูที่หลากหลาย - คุณภาพของวัตถุดิบและเมล็ดกาแฟ - ราคาเครื่องดื่มไม่แพงเมื่อเทียบกับคุณภาพ | Customer Relationships (วิธีการรักษารฐานลูกค้าเป็นแบบไหน) <ul style="list-style-type: none"> - ให้อาหารชิมกาแฟเมนูใหม่ๆ - มีโปรแกรมสมาชิกส่งกาแฟให้ลูกค้า - สร้างความสัมพันธ์ที่ดีกับลูกค้าสามารถจดจำลูกค้าประจำได้ว่าลูกค้าชอบแบบไหนต้องการอะไร | Customer Segments (ลูกค้าเราคือใคร มีหน้าตาแบบไหน) <ul style="list-style-type: none"> - คนที่มีรสนิยมชื่นชอบกาแฟ - คนที่ชอบทานเครื่องดื่มและถ่ายรูป - คนชอบทานเครื่องดื่มมีคุณภาพ |
| Key Resources (ทรัพยากรที่จำเป็นที่ต้องใช้ในการดำเนิน) <ul style="list-style-type: none"> - สูตรกาแฟและเครื่องดื่มที่คิดขึ้นมาเอง - การฝึกอบรมพนักงานที่มีทักษะ | | Channels (ช่องทางที่เข้าถึงลูกค้า) <ul style="list-style-type: none"> - Social Media (Facebook, Instagram, Line) - Call - Food Delivery (Grab, Lineman, Robinhood) | | |
| Cost Structure <ul style="list-style-type: none"> - ค่าเช่าพื้นที่ - ค่าตกแต่งร้าน - วัตถุดิบและอุปกรณ์ - ค่าเงินเดือนพนักงาน - ค่าน้ำค่าไฟ - ค่าประชาสัมพันธ์ | | Revenue Streams (กระแสรายได้ของธุรกิจนี้) <ul style="list-style-type: none"> - เครื่องดื่ม | | |

ภาพที่ 3.2 Business Canvas ของร้านกาแฟที่ทำการวิจัย

นอกจากนั้นแล้วผู้วิจัยได้เข้าศึกษาเพิ่มเติมโดยการเข้าไปเป็นพนักงานในร้านเพื่อทำความเข้าใจและเพิ่มโอกาสเห็นปัญหาจริงที่สามารถกำหนดขอบเขตและนำข้อมูลมาวิเคราะห์เพื่อแก้ปัญหาได้พบว่าปัญหาหลักๆ ที่พบคือการประมาณการสั่งวัตถุดิบที่เก็บไว้ไม่ได้นาน เช่น นมสด

น้ำมะนาวสด หรือการผสมสูตรนมไม่เพียงพอต่อการขายต่อวันดังนั้นเราจึงนำขอบเขตของปัญหานี้ มากำหนดขอบเขตในการวิเคราะห์ในงานวิจัยนี้

3.1.1.2 กรณีศึกษาเป็นประเภทข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม

กรณีนี้เป็นองค์การรัฐวิสาหกิจในสังกัดกระทรวงคมนาคมที่จัดขึ้นเพื่อดำเนินการระบบขนส่งมวลชนในรูปแบบรถไฟฟ้าในพื้นที่กรุงเทพมหานครและปริมณฑลมีเป้าหมายเพื่อบรรเทา ความรุนแรงของปัญหาจราจรในประเทศเพื่ออำนวยความสะดวกให้ประชาชนเดินทางโดยไม่ต้อง ใช้รถยนต์ส่วนบุคคลดังนั้นเป้าหมายหลักของกรณีศึกษานี้ คือต้องการให้ประชาชนหันมาใช้ระบบขนส่งมวลชนให้มากขึ้นเพื่อลดปัญหาการจราจรลดค่าครองชีพให้กับประชาชน

3.1.1.3 กรณีศึกษาเป็นประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3

กรณีนี้เป็นองค์การของรัฐกระทรวงสาธารณสุขมีพันธกิจพัฒนาและอภิบาลระบบสุขภาพอย่างมีส่วนร่วมและยั่งยืน มีวิสัยทัศน์คือเป็นองค์กรหลักด้านสุขภาพที่รวมพลังสังคมเพื่อ ประชาชนสุขภาพดีจากสถานการณ์การแพร่ระบาดของโรคติดต่อเชื้อไวรัสโคโรนา 2019 ซึ่งมีการ แพร่ระบาดทั่วโลก รวมถึงประเทศไทยส่งผลกระทบในวงกว้างด้านสาธารณสุข เศรษฐกิจ สังคม ความมั่นคงตลอดการดำรงชีพของประชาชนดังนั้นพันธกิจหลักขององค์กรนี้คือให้คำแนะนำ วางแผนออกแนวทางการปฏิบัติ และมาตรการในการป้องกันความเสี่ยงจากโรคดังกล่าว

3.1.2 การทำความเข้าใจกับข้อมูล (Data Understanding)

3.1.2.1 กรณีศึกษาเป็นประเภทข้อมูลยอดขายของร้านกาแฟ

รวบรวมข้อมูลที่เกี่ยวข้องกับการสร้างแบบจำลองพยากรณ์ยอดขายรายวัน โดยเลือกใช้ ข้อมูลยอดขายรายวันจากเครื่อง POS (Point of sale system) ของร้านกาแฟเพราะมีความน่าเชื่อถือและ มีความเหมาะสมกับข้อมูล

| วันที่ | รหัสสินค้า | ชื่อสินค้า | กลุ่ม | หมวดสินค้า | ต้นทุนเฉลี่ย | ราคาขายเฉลี่ย | กำไรเฉลี่ย | จำนวนการขาย | ยอดโอน | ต้นทุน | ส่วนลด | ราคาสุทธิ |
|------------|------------|--------------------|-------|------------|--------------|---------------|------------|-------------|--------|--------|--------|-----------|
| 03/08/2021 | | สตรอบอรี่เมล็ดอ่อน | | ไทย ไรตา | 0 | 45 | 45 | 1 | 45 | 0 | 0 | 45 |
| 03/08/2021 | | สตรอบอรี่เมล็ดอ่อน | | ไทย ไรตา | 0 | 50 | 50 | 1 | 50 | 0 | 0 | 50 |
| 03/08/2021 | | บลูเบอร์รี่ | | ไทย ไรตา | 0 | 40 | 40 | 1 | 40 | 0 | 0 | 40 |
| 03/08/2021 | | ป่าละเมาะ | | ไทย ไรตา | 0 | 45 | 45 | 1 | 45 | 0 | 0 | 45 |
| 03/08/2021 | | แบบไรตา | | ไทย ไรตา | 0 | 35 | 35 | 1 | 35 | 0 | 0 | 35 |
| 02/08/2021 | | ธัญพืช | | ไทย เมม | 0 | 45 | 45 | 2 | 90 | 0 | 0 | 90 |
| 02/08/2021 | | โกโก้คั่ว | | ไทย เมม | 0 | 40 | 40 | 1 | 40 | 0 | 0 | 40 |
| 02/08/2021 | | เมล็ดทรายมา | | ไทย เมม | 0 | 35 | 35 | 1 | 35 | 0 | 0 | 35 |
| 02/08/2021 | | นมพาสเจอร์ | | ไทย เมม | 0 | 40 | 40 | 1 | 40 | 0 | 0 | 40 |
| 02/08/2021 | | โวลันโบลด์ | | ไทย เมม | 0 | 30 | 30 | 1 | 30 | 0 | 0 | 30 |
| 02/08/2021 | | โวลัน | | ไทย เมม | 0 | 30 | 30 | 1 | 30 | 0 | 0 | 30 |
| 03/08/2021 | | โกโก้ | | ไทย เมม | 0 | 38.75 | 38.75 | 8 | 310 | 0 | 0 | 310 |
| 03/08/2021 | | ธัญพืช | | ไทย เมม | 0 | 41.25 | 41.25 | 4 | 165 | 0 | 0 | 165 |
| 03/08/2021 | | ธัญพืชคั่ว | | ไทย เมม | 0 | 48.33 | 48.33 | 3 | 145 | 0 | 0 | 145 |
| 03/08/2021 | | นมผง | | ไทย เมม | 0 | 35 | 35 | 1 | 35 | 0 | 0 | 35 |
| 03/08/2021 | | นมผงโปรตีน | | ไทย เมม | 0 | 40 | 40 | 3 | 120 | 0 | 0 | 120 |
| 02/08/2021 | | โกโก้ | | ไทย เมม | 0 | 35 | 35 | 1 | 35 | 0 | 0 | 35 |
| 02/08/2021 | | ชานชา | | ไทย ยา | 0 | 42.5 | 42.5 | 4 | 170 | 0 | 0 | 170 |
| 02/08/2021 | | ชานชาเมล็ด | | ไทย ยา | 0 | 35 | 35 | 1 | 35 | 0 | 0 | 35 |
| 02/08/2021 | | ชาเขียวเมล็ด | | ไทย ยา | 0 | 40 | 40 | 16 | 640 | 0 | 0 | 640 |

ภาพที่ 3.3 ตัวอย่างข้อมูลยอดขายสินค้ารายวันที่จะนำมาวิเคราะห์

3.1.2.2 กรณีศึกษาเป็นประเภทข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม

รวบรวมข้อมูลที่เกี่ยวข้องกับการสร้างแบบจำลองพยากรณ์จำนวนผู้โดยสารรวมรายเดือน โดยเลือกใช้ข้อมูลจากองค์การไฟฟ้าขนส่งมวลชนแห่งประเทศไทยเพราะมีความน่าเชื่อถือ และมีการปรับปรุงข้อมูลให้อัปเดตอยู่เสมอ

| เดือน | ปี | โครงการ | จำนวนผู้โดยสารรวม | จำนวนโดยสารเฉลี่ยรายวัน | จำนวนผู้โดยสารเฉลี่ยรายวันธรรมดา | จำนวนผู้โดยสารเฉลี่ยรายวันหยุด |
|------------|------|-------------------------------------|-------------------|-------------------------|----------------------------------|--------------------------------|
| สิงหาคม | 2559 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 524,472 | 20,172 | 20,492 | 19,568 |
| กันยายน | 2559 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 629,963 | 20,999 | 22,658 | 16,437 |
| ตุลาคม | 2559 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 650,772 | 20,993 | 24,212 | 15,140 |
| พฤศจิกายน | 2559 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 647,337 | 21,578 | 24,238 | 14,264 |
| ธันวาคม | 2559 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 593,981 | 19,161 | 23,048 | 12,093 |
| มกราคม | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 634,552 | 20,469 | 24,282 | 13,538 |
| กุมภาพันธ์ | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 611,029 | 21,822 | 25,410 | 14,248 |
| มีนาคม | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 719,225 | 23,201 | 25,447 | 16,741 |
| เมษายน | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 631,343 | 21,045 | 25,678 | 15,750 |
| พฤษภาคม | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 748,456 | 24,144 | 28,127 | 16,901 |
| มิถุนายน | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 817,301 | 27,243 | 30,217 | 19,066 |
| กรกฎาคม | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 823,504 | 26,565 | 31,501 | 18,749 |
| สิงหาคม | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 1,331,334 | 42,946 | 47,655 | 31,436 |
| กันยายน | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 1,407,258 | 46,909 | 52,322 | 34,278 |
| ตุลาคม | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 1,376,904 | 45,897 | 53,861 | 32,140 |
| พฤศจิกายน | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 1,395,555 | 46,519 | 52,042 | 31,330 |
| ธันวาคม | 2560 | โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 1,249,626 | 40,311 | 48,421 | 27,470 |

ภาพที่ 3.4 ตัวอย่างข้อมูลจำนวนผู้โดยสารโครงการรถไฟฟ้ามหานครสายเฉลิมรัชธรรมที่นำมาวิเคราะห์

3.1.2.3 กรณีศึกษาเป็นประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3 รวบรวมข้อมูลที่เกี่ยวข้องกับการสร้างแบบจำลองพยากรณ์จำนวนผู้เสียชีวิตรายใหม่รายวันโดยเลือกใช้ข้อมูลจาก API ของกระทรวงสาธารณสุขเพราะมีความน่าเชื่อถือและมีการอัปเดตข้อมูลทุกวัน

```

[[{"tmn_date": "2021-04-01", "new_case": 16, "total_case": 2889, "new_case_excludeabroad": 13, "total_case_excludeabroad": 2890, "new_death": 0, "total_death": 0, "new_recovered": 122, "total_recovered": 12767, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-02", "new_case": 15, "total_case": 2894, "new_case_excludeabroad": 12, "total_case_excludeabroad": 2895, "new_death": 0, "total_death": 0, "new_recovered": 155, "total_recovered": 12922, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-03", "new_case": 84, "total_case": 2903, "new_case_excludeabroad": 69, "total_case_excludeabroad": 2904, "new_death": 1, "total_death": 1, "new_recovered": 162, "total_recovered": 13084, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-04", "new_case": 96, "total_case": 2912, "new_case_excludeabroad": 87, "total_case_excludeabroad": 2908, "new_death": 0, "total_death": 0, "new_recovered": 196, "total_recovered": 13280, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-05", "new_case": 194, "total_case": 2931, "new_case_excludeabroad": 163, "total_case_excludeabroad": 2913, "new_death": 0, "total_death": 0, "new_recovered": 166, "total_recovered": 13446, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-06", "new_case": 290, "total_case": 2951, "new_case_excludeabroad": 245, "total_case_excludeabroad": 2926, "new_death": 0, "total_death": 0, "new_recovered": 169, "total_recovered": 13615, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-07", "new_case": 334, "total_case": 2989, "new_case_excludeabroad": 327, "total_case_excludeabroad": 2975, "new_death": 0, "total_death": 0, "new_recovered": 131, "total_recovered": 13746, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-08", "new_case": 465, "total_case": 3039, "new_case_excludeabroad": 391, "total_case_excludeabroad": 3274, "new_death": 0, "total_death": 0, "new_recovered": 132, "total_recovered": 13878, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-09", "new_case": 559, "total_case": 3089, "new_case_excludeabroad": 549, "total_case_excludeabroad": 3269, "new_death": 1, "total_death": 1, "new_recovered": 127, "total_recovered": 14005, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-10", "new_case": 787, "total_case": 3168, "new_case_excludeabroad": 781, "total_case_excludeabroad": 3246, "new_death": 0, "total_death": 0, "new_recovered": 139, "total_recovered": 14144, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-11", "new_case": 987, "total_case": 3267, "new_case_excludeabroad": 984, "total_case_excludeabroad": 3249, "new_death": 0, "total_death": 0, "new_recovered": 143, "total_recovered": 14287, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-12", "new_case": 985, "total_case": 3366, "new_case_excludeabroad": 980, "total_case_excludeabroad": 3248, "new_death": 0, "total_death": 0, "new_recovered": 134, "total_recovered": 14421, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-13", "new_case": 965, "total_case": 3465, "new_case_excludeabroad": 937, "total_case_excludeabroad": 3276, "new_death": 0, "total_death": 0, "new_recovered": 140, "total_recovered": 14561, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-14", "new_case": 1138, "total_case": 3583, "new_case_excludeabroad": 1130, "total_case_excludeabroad": 3276, "new_death": 0, "total_death": 0, "new_recovered": 134, "total_recovered": 14695, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-15", "new_case": 1541, "total_case": 3700, "new_case_excludeabroad": 1540, "total_case_excludeabroad": 3245, "new_death": 0, "total_death": 0, "new_recovered": 161, "total_recovered": 14856, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-16", "new_case": 1552, "total_case": 3830, "new_case_excludeabroad": 1576, "total_case_excludeabroad": 3282, "new_death": 1, "total_death": 1, "new_recovered": 197, "total_recovered": 15053, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-17", "new_case": 1547, "total_case": 3985, "new_case_excludeabroad": 1544, "total_case_excludeabroad": 3736, "new_death": 0, "total_death": 0, "new_recovered": 190, "total_recovered": 15243, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-18", "new_case": 1767, "total_case": 4132, "new_case_excludeabroad": 1709, "total_case_excludeabroad": 3913, "new_death": 1, "total_death": 1, "new_recovered": 115, "total_recovered": 15358, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-19", "new_case": 1398, "total_case": 4374, "new_case_excludeabroad": 1380, "total_case_excludeabroad": 4057, "new_death": 1, "total_death": 1, "new_recovered": 104, "total_recovered": 15462, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-20", "new_case": 1442, "total_case": 4525, "new_case_excludeabroad": 1443, "total_case_excludeabroad": 4359, "new_death": 0, "total_death": 0, "new_recovered": 121, "total_recovered": 15583, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-21", "new_case": 1458, "total_case": 4663, "new_case_excludeabroad": 1454, "total_case_excludeabroad": 4342, "new_death": 1, "total_death": 1, "new_recovered": 110, "total_recovered": 15693, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-22", "new_case": 1476, "total_case": 4815, "new_case_excludeabroad": 1470, "total_case_excludeabroad": 4682, "new_death": 1, "total_death": 1, "new_recovered": 177, "total_recovered": 15870, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-23", "new_case": 2070, "total_case": 5028, "new_case_excludeabroad": 2061, "total_case_excludeabroad": 4684, "new_death": 1, "total_death": 1, "new_recovered": 144, "total_recovered": 16014, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-24", "new_case": 2839, "total_case": 5302, "new_case_excludeabroad": 2827, "total_case_excludeabroad": 4971, "new_death": 0, "total_death": 0, "new_recovered": 177, "total_recovered": 16191, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-25", "new_case": 2438, "total_case": 5560, "new_case_excludeabroad": 2431, "total_case_excludeabroad": 5206, "new_death": 1, "total_death": 1, "new_recovered": 147, "total_recovered": 16338, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-26", "new_case": 2048, "total_case": 5798, "new_case_excludeabroad": 2037, "total_case_excludeabroad": 5424, "new_death": 0, "total_death": 0, "new_recovered": 140, "total_recovered": 16478, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-27", "new_case": 2179, "total_case": 5987, "new_case_excludeabroad": 2174, "total_case_excludeabroad": 5443, "new_death": 1, "total_death": 1, "new_recovered": 195, "total_recovered": 16673, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-28", "new_case": 2012, "total_case": 6189, "new_case_excludeabroad": 2003, "total_case_excludeabroad": 5841, "new_death": 1, "total_death": 1, "new_recovered": 181, "total_recovered": 16854, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-29", "new_case": 1871, "total_case": 6370, "new_case_excludeabroad": 1864, "total_case_excludeabroad": 6020, "new_death": 1, "total_death": 1, "new_recovered": 192, "total_recovered": 17046, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-04-30", "new_case": 1583, "total_case": 6533, "new_case_excludeabroad": 1579, "total_case_excludeabroad": 6185, "new_death": 1, "total_death": 1, "new_recovered": 180, "total_recovered": 17226, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-01", "new_case": 1693, "total_case": 6704, "new_case_excludeabroad": 1694, "total_case_excludeabroad": 6374, "new_death": 1, "total_death": 1, "new_recovered": 182, "total_recovered": 17408, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-02", "new_case": 1940, "total_case": 6886, "new_case_excludeabroad": 1930, "total_case_excludeabroad": 6567, "new_death": 1, "total_death": 1, "new_recovered": 183, "total_recovered": 17591, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-03", "new_case": 2041, "total_case": 7102, "new_case_excludeabroad": 2040, "total_case_excludeabroad": 6773, "new_death": 1, "total_death": 1, "new_recovered": 170, "total_recovered": 17761, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-04", "new_case": 1753, "total_case": 7289, "new_case_excludeabroad": 1750, "total_case_excludeabroad": 6945, "new_death": 1, "total_death": 1, "new_recovered": 140, "total_recovered": 17901, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-05", "new_case": 2112, "total_case": 7480, "new_case_excludeabroad": 2107, "total_case_excludeabroad": 7157, "new_death": 1, "total_death": 1, "new_recovered": 188, "total_recovered": 18089, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-06", "new_case": 1911, "total_case": 7681, "new_case_excludeabroad": 1902, "total_case_excludeabroad": 7372, "new_death": 1, "total_death": 1, "new_recovered": 245, "total_recovered": 18334, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-07", "new_case": 2044, "total_case": 7893, "new_case_excludeabroad": 2040, "total_case_excludeabroad": 7512, "new_death": 1, "total_death": 1, "new_recovered": 237, "total_recovered": 18571, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-08", "new_case": 2439, "total_case": 8127, "new_case_excludeabroad": 2409, "total_case_excludeabroad": 7792, "new_death": 1, "total_death": 1, "new_recovered": 224, "total_recovered": 18795, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-09", "new_case": 2261, "total_case": 8375, "new_case_excludeabroad": 2080, "total_case_excludeabroad": 8000, "new_death": 1, "total_death": 1, "new_recovered": 199, "total_recovered": 19024, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-10", "new_case": 1630, "total_case": 8589, "new_case_excludeabroad": 1621, "total_case_excludeabroad": 8163, "new_death": 1, "total_death": 1, "new_recovered": 160, "total_recovered": 19184, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-11", "new_case": 1919, "total_case": 8824, "new_case_excludeabroad": 1903, "total_case_excludeabroad": 8353, "new_death": 1, "total_death": 1, "new_recovered": 182, "total_recovered": 19366, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-12", "new_case": 1983, "total_case": 8987, "new_case_excludeabroad": 1974, "total_case_excludeabroad": 8505, "new_death": 1, "total_death": 1, "new_recovered": 200, "total_recovered": 19566, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-13", "new_case": 4897, "total_case": 9374, "new_case_excludeabroad": 4875, "total_case_excludeabroad": 9075, "new_death": 1, "total_death": 1, "new_recovered": 1572, "total_recovered": 21138, "update_date": "2021-09-01 07:40:49"}, {"tmn_date": "2021-05-14", "new_case": 2286, "total_case": 9600, "new_case_excludeabroad": 2251, "total_case_excludeabroad": 9262, "new_death": 1, "total_death": 1, "new_recovered": 170, "total_recovered": 21308, "update_date": "2021-09-01 07:40:49"}]]

```

ภาพที่ 3.5 ตัวอย่างข้อมูลจำนวนป่วย COVID-19 ระลอก 3 ที่จะนำมาวิเคราะห์

3.1.3 การเตรียมข้อมูล (Data Preparation)

การเตรียมข้อมูลถือเป็นขั้นตอนที่ใช้เวลามากที่สุดในขั้นตอนทั้งหมดเพราะเป็นส่วนที่สำคัญที่ส่งผลต่อประสิทธิภาพในการพยากรณ์โดยการเตรียมข้อมูลประกอบด้วย

3.1.3.1 การคัดเลือกข้อมูล (Data Selection)

สำหรับงานวิจัยนี้จะเลือกข้อมูลที่เป็นตัวเลขมาใช้นี้เนื่องจากเป็นการวิเคราะห์อนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่องตามรายละเอียดดังภาพที่ 3.6

| ฟิลต์บังคับ | รูปแบบที่ 1: เทคนิค ARIMA, SARIMA | | รูปแบบที่ 2: ใช้กับเทคนิค Multiple Linear Regression, Decision Tree Random Forest | |
|---|---|-------------------|---|------------------------------------|
| | Attribute | Target | Attribute | Target |
| | Date | Total | Date, ID | Total |
| กรณีศึกษาเป็นประเภทข้อมูล ยอดขายของร้านกาแฟ | เวลา (Date), ยอดขาย (Total) | ยอดขาย (Total) | เวลา (Date), รหัสสินค้า (ID), อุณหภูมิ (Temperature), วันหยุด (Holiday), มีการจัดโปรโมชั่น (Promotion), | ยอดขาย (Total) |
| กรณีศึกษาเป็นประเภทข้อมูล จำนวนผู้ใช้บริการ โครงการ รถไฟฟ้าฟ้ามหานคร สายฉลองรัช ธรรม | เวลา (Date), จำนวน ผู้โดยสาร รวม (Total) | | | |
| กรณีศึกษาเป็นประเภทข้อมูล รายงานสถานการณ์ COVID-19 ระลอก 3 | | | วันแถลง (Date), รหัสประเทศ (ID), จำนวนผู้ป่วยรายใหม่ (new_case), จำนวนผู้ป่วยรักษาหายรายใหม่ (new_recovered) | จำนวนผู้ป่วย รายใหม่ (Total) |

ภาพที่ 3.6 การคัดเลือกข้อมูลสำหรับการสร้างแบบจำลองสำหรับกรณีศึกษาและเทคนิคต่าง ๆ

3.1.3.2 การกลั่นกรองข้อมูล (Data Cleaning)

ขั้นตอนนี้เนื่องจากข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม และข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3 ก่อนข้างสมบูรณ์แล้วดังนั้นเราจะทำขั้นตอนนี้เฉพาะข้อมูลยอดขายของร้านกาแฟที่ยังมีบางส่วนต้องกรองข้อมูลเพิ่มเติม

รูปแบบอนุกรมเวลาควรเป็นช่วงเวลาที่ต่อเนื่อง เช่น รายวัน รายเดือน รายปี ในกรณีศึกษานี้มีการพบข้อมูลที่หายไป (missing value) คือเป็นช่วงเวลาที่มีการปิดร้านค้าจึงได้มีการแทนที่ข้อมูลด้วยค่าเฉลี่ยของยอดขายภายในเดือนนั้นๆ มาทำการแทนที่และมีการกรองข้อมูลที่เป็นคำผิดปกติก่อนที่จะนำไปสร้างแบบจำลองการพยากรณ์ คือรายการที่ไม่มียอดขายจริงเช่น ชื่อสินค้าคือ “รายการไม่รับช้อนส้อมพลาสติก”, “ใช้สิทธิโครงการรัฐ”, “ลด 50% สูงสุด 100.- เพียงใส่โค้ด "KBANK" ในหน้าชำระเงิน” หรือ “ของแถม” เป็นต้น

3.1.3.3 การแปลงรูปแบบข้อมูล (Data Transformation)

ขั้นตอนนี้สร้างข้อมูลให้อยู่ในรูปแบบที่พร้อมนำไปใช้ในการวิเคราะห์คือ

1. กรณีศึกษาเป็นประเภทข้อมูลยอดขายของร้านกาแฟ

เพิ่มแอตทริบิวต์ที่อาจจะเพิ่มประสิทธิภาพในการวิเคราะห์เพิ่มเติม เช่น กรณีศึกษาเป็นประเภทข้อมูลยอดขายของร้านกาแฟ วันหยุด (อ้างอิงวันหยุดตามประเพณีของสถาบันการเงิน), มีการจัดโปรโมชั่นหรือไม่, อุณหภูมิ (องศาเซลเซียส)

แปลงข้อมูลให้เป็น Numerical (ตัวเลข) เช่น วันหยุด (Holiday), การจัดโปรโมชั่น (Promotion) ให้อยู่ในรูปแบบของ {0, 1} ดังตารางที่ 3.1 และแปลงข้อมูลวันที่ออกเป็นวันที่เท่าไรของสัปดาห์ (วันจันทร์-อาทิตย์ 0-6), เดือน, ปี, วันที่ ดังตารางที่ 3.2 ซึ่งขั้นตอนดังกล่าวจะถูกปรับให้หลังจากนำข้อมูลเข้าระบบแล้ว

ตารางที่ 3.1 ตัวอย่างการแปลงแอตทริบิวต์วันหยุด (Holiday), การจัดโปรโมชั่น (Promotion)

| IS Holiday | Promotion |
|------------|-----------|
| 0 | 0 |
| 0 | 0 |
| 0 | 1 |
| 0 | 0 |

ตารางที่ 3.2 ตัวอย่างการแปลงแอตทริบิวต์วันที่ (Date)

| Date_dayofweek | Date_month | Date_year | Date_day |
|----------------|------------|-----------|----------|
| 0 | 6 | 2020 | 22 |
| 3 | 3 | 2020 | 5 |

จัดกลุ่มข้อมูลสินค้าชนิดเดียวกันไว้เป็นข้อมูลชุดเดียวกัน เช่น ข้อมูลสินค้าที่นำมาใช้บางรายการคือประเภทสินค้าชนิดเดียวกันแต่ชื่อสินค้าไม่ตรงกันดังตารางด้านล่าง

ตารางที่ 3.3 ตัวอย่างการจัดกลุ่มชื่อสินค้าก่อนนำเข้าวิเคราะห์

| ID | ชื่อสินค้า |
|----|--------------------------------------|
| 1 | .อเมริกาโน่ น้ำผึ้ง |
| 1 | *พรีเมียม อเมริกาโน่เย็น (AMERICANO) |
| 1 | *อเมริกาโน่ |
| 1 | *อเมริกาโน่เย็น (AMERICANO) |
| 1 | *อเมริกาโน่เย็น* (AMERICANO) |
| 1 | พรีเมียม อเมริกาโน่เย็น (AMERICANO) |
| 1 | อเมริกาโน่ |

- กรณีศึกษาเป็นประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3
เพิ่มแอตทริบิวต์ ID เป็นรหัสประเทศเพื่อนำเข้าระบบ (ฟิลด์บังคับ) กรณีนี้จะใส่เป็นรหัส ID เดียวกัน และแปลงข้อมูลวันที่ตามตารางที่ 3.2
จากขั้นตอนในการเตรียมข้อมูลข้างต้น งานวิจัยนี้จะแบ่งรูปแบบไฟล์ข้อมูลเพื่อนำเข้าระบบเพื่อสร้างแบบจำลองและพยากรณ์ดังภาพที่ 3.7

| | รูปแบบที่ 1: เทคนิค | รูปแบบที่ 2: ใช้กับเทคนิค Multiple Linear Regression, Decision | |
|--|---|--|---|
| | ARIMA, SARIMA | Tree Random Forest | |
| การแบ่งไฟล์ก่อนนำเข้าแบบจำลองและพยากรณ์ | คอลัมภ์ข้อมูลในไฟล์ train | คอลัมภ์ข้อมูลในไฟล์ train | คอลัมภ์ข้อมูลในไฟล์ test |
| กรณีศึกษาเป็นประเภทข้อมูลยอดขายของร้านกาแฟ | เวลา (Date), ยอดขาย (Total) | เวลา (Date), รหัสสินค้า (ID), อุณหภูมิ (Temperature), วันหยุด (Holiday), มีการจัดโปรโมชั่น (Promotion), ยอดขาย (Total) | เวลา (Date), รหัสสินค้า (ID), อุณหภูมิ (Temperature), วันหยุด (Holiday), มีการจัดโปรโมชั่น (Promotion), |
| กรณีศึกษาเป็นประเภทข้อมูลจำนวนผู้ใช้บริการโครงการรถไฟฟ้ามหานครสายฉลองรัชธรรม | เวลา (Date), จำนวนผู้โดยสารรวม (Total) | | |
| กรณีศึกษาเป็นประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3 | | วันแถลง (Date), รหัสประเทศ (ID), จำนวนผู้ป่วยรายใหม่ (new_case), จำนวนผู้ป่วยรักษาหายรายใหม่ (new_recovered), จำนวนผู้ป่วยตายรายใหม่ (Total) | วันแถลง (Date), รหัสประเทศ (ID), จำนวนผู้ป่วยรายใหม่ (new_case), จำนวนผู้ป่วยรักษาหายรายใหม่ (new_recovered) |

ภาพที่ 3.7 คอลัมภ์ในข้อมูลไฟล์ train, test สำหรับการสร้างแบบจำลองสำหรับกรณีศึกษาและเทคนิคต่าง ๆ

และมีรูปแบบของข้อมูลก่อนนำเข้าแบบจำลองจริงหลังการแปลงข้อมูลผ่านระบบอัตโนมัติในแต่ละกรณีดังภาพที่ 3.8 - 3.11

| Date | Total |
|----------|------------|
| 3/1/2021 | 7,100.0000 |
| 3/2/2021 | 6,940.0000 |
| 3/3/2021 | 7,355.0000 |
| 3/4/2021 | 7,395.0000 |
| 3/5/2021 | 6,890.0000 |
| 3/6/2021 | 8,060.0000 |

ภาพที่ 3.8 รูปแบบของข้อมูลจริงก่อนนำเข้าสร้างแบบจำลองด้วยเทคนิค ARIMA, SARIMA
กรณีศึกษาข้อมูลประเภทยอดขายของร้านกาแฟ

| Date | Total |
|----------|---------|
| 1/1/2562 | 1591763 |
| 2/1/2562 | 1443333 |
| 3/1/2562 | 1664735 |
| 4/1/2562 | 1458908 |
| 5/1/2562 | 1530913 |
| 6/1/2562 | 1675410 |

ภาพที่ 3.9 รูปแบบของข้อมูลจริงก่อนนำเข้าสร้างแบบจำลองด้วยเทคนิค ARIMA, SARIMA
กรณีศึกษาข้อมูลจำนวนผู้ใช้บริการโครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม

| Temperature | IS Holiday | Promotion | ID | UnitPrice | Date_dayofweek | Date_month | Date_year | Date_day |
|-------------|------------|-----------|-----|-----------|----------------|------------|-----------|----------|
| 29.1100 | 0 | 1 | 181 | 210 | 3 | 8 | 2021 | 5 |
| 30.3900 | 0 | 1 | 84 | 65 | 6 | 5 | 2021 | 30 |
| 28.0000 | 0 | 1 | 137 | 50 | 2 | 4 | 2021 | 28 |
| 32.1100 | 0 | 1 | 155 | 50 | 0 | 7 | 2021 | 5 |
| 31.0000 | 0 | 1 | 84 | 65 | 2 | 3 | 2021 | 24 |
| 29.1700 | 0 | 1 | 160 | 65 | 2 | 3 | 2021 | 3 |
| 32.2200 | 0 | 1 | 62 | 50 | 3 | 5 | 2021 | 6 |
| 28.9400 | 0 | 1 | 100 | 50 | 0 | 3 | 2021 | 22 |
| 31.0000 | 0 | 1 | 105 | 60 | 3 | 4 | 2021 | 22 |
| 32.1100 | 0 | 1 | 1 | 100 | 0 | 7 | 2021 | 5 |

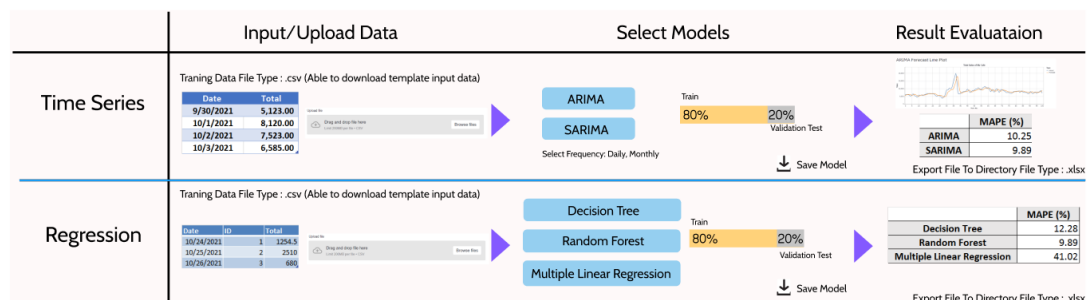
ภาพที่ 3.10 รูปแบบของข้อมูลจริงก่อนนำเข้าสร้างแบบจำลองด้วยเทคนิค Multiple Linear
Regression, Decision Tree, Random Forest กรณีศึกษาข้อมูลประเภทยอดขายของร้าน
กาแฟแบ่งตามรหัสสินค้า

| ID | new_case | new_recovered | Date_dayofweek | Date_month | Date_year | Date_day |
|----|----------|---------------|----------------|------------|-----------|----------|
| 1 | 8685 | 3797 | 1 | 7 | 2021 | 13 |
| 1 | 11786 | 14738 | 1 | 9 | 2021 | 14 |
| 1 | 14956 | 17936 | 3 | 9 | 2021 | 2 |
| 1 | 6519 | 4148 | 2 | 7 | 2021 | 7 |
| 1 | 16031 | 15417 | 3 | 9 | 2021 | 9 |
| 1 | 2804 | 4143 | 6 | 6 | 2021 | 13 |
| 1 | 2331 | 4947 | 2 | 6 | 2021 | 16 |

ภาพที่ 3.11 รูปแบบของข้อมูลจริงก่อนนำเข้าสู่สร้างแบบจำลองด้วยเทคนิค Multiple Linear Regression, Decision Tree, Random Forest กรณีศึกษาข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3

3.1.4 การสร้างแบบจำลอง (Modeling)

Build A New Predictive Model



ภาพที่ 3.12 ขั้นตอนการทำงานของระบบเพื่อสร้างแบบจำลอง

ขั้นตอนนี้จะใช้ข้อมูลจากข้อ 3.1.3 มาทำการสร้างแบบจำลองโดยการแบ่งออกเป็น 2 กลุ่มหลักๆ คือ

3.1.4.1 อนุกรมเวลา (Time series)

ข้อมูลจะถูกนำเข้าเพื่อสร้างแบบจำลองด้วยวิธี ARIMA หรือ SARIMA โดยแบ่งข้อมูลจำนวน 80% เป็นข้อมูลสำหรับสร้างแบบจำลองและข้อมูลจำนวน 20% เป็นข้อมูลใช้เป็นข้อมูลทดสอบแบบจำลองที่สร้างขึ้นได้ใช้วิธีการ Split Test

การพิจารณารูปแบบของ ARIMA

กรณีที่เป็น Non Stationary series ที่มีแนวโน้มและมีความจำเป็นต้องทำ Regular Differencing หากค่า (d) อันดับความหยาบหนึ่ง (จำนวนครั้งที่หาผลต่างเพื่อปรับอนุกรมเวลาให้เป็นอนุกรมเวลาที่ Stationary) ด้วยวิธี ADF test หาก p, q ใช้วิธีขยับตัวเลขไปที่ละขั้นหาแบบจำลองโอกาสที่เป็นไปได้และพิจารณาค่า AIC น้อยที่สุดเป็นแบบจำลองในการพยากรณ์

```

Performing stepwise search to minimize aic
ARIMA(0,1,0)(0,0,0)[0] intercept : AIC=6710.951, Time=0.05 sec
ARIMA(1,1,0)(0,0,0)[0] intercept : AIC=6675.320, Time=0.07 sec
ARIMA(0,1,1)(0,0,0)[0] intercept : AIC=6662.405, Time=0.21 sec
ARIMA(0,1,0)(0,0,0)[0] intercept : AIC=6708.995, Time=0.01 sec
ARIMA(1,1,1)(0,0,0)[0] intercept : AIC=6657.294, Time=0.36 sec
ARIMA(2,1,1)(0,0,0)[0] intercept : AIC=6653.920, Time=0.26 sec
ARIMA(2,1,0)(0,0,0)[0] intercept : AIC=6670.665, Time=0.06 sec
ARIMA(3,1,1)(0,0,0)[0] intercept : AIC=6654.851, Time=0.39 sec
ARIMA(2,1,2)(0,0,0)[0] intercept : AIC=inf, Time=0.61 sec
ARIMA(1,1,2)(0,0,0)[0] intercept : AIC=6651.920, Time=0.49 sec
ARIMA(0,1,2)(0,0,0)[0] intercept : AIC=6660.093, Time=0.18 sec
ARIMA(1,1,3)(0,0,0)[0] intercept : AIC=inf, Time=0.56 sec
ARIMA(0,1,3)(0,0,0)[0] intercept : AIC=6659.099, Time=0.34 sec
ARIMA(2,1,3)(0,0,0)[0] intercept : AIC=inf, Time=0.67 sec
ARIMA(1,1,2)(0,0,0)[0] intercept : AIC=6651.986, Time=0.12 sec

Best model: ARIMA(1,1,2)(0,0,0)[0] intercept
Total fit time: 4.385 seconds
    
```

ภาพที่ 3.13 ตัวอย่างการหารูปแบบ ARIMA (p,d,q) จากการพิจารณาค่า AIC

การพิจารณารูปแบบของ SARIMA

กรณีที่เป็น Non Stationary series ที่มีแนวโน้มและมีอิทธิพลของฤดูกาลมาเกี่ยวข้องมีความจำเป็นต้องทำ Seasonal Differencing หากค่า (D) อันดับความหยาบหนึ่ง (จำนวนครั้งที่หาผลต่างเพื่อปรับอนุกรมเวลาให้เป็นอนุกรมเวลาที่ Stationary และไม่มีความเป็นฤดูกาล) ด้วยวิธี CH test หาก p, q และ P,Q ใช้วิธีขยับตัวเลขไปที่ละขั้นหาแบบจำลองโอกาสที่เป็นไปได้และแทนค่า s ด้วยคาบของฤดูกาล⁵ รายวันแทนด้วย 7 รายเดือนแทนด้วย 12 งานวิจัยนี้จะมีข้อมูลที่น่ามาทดลองเป็นลักษณะ รายวันและรายเดือน ซึ่งค่า s จะถูกปรับตามคาบของฤดูกาลที่ใช้และพิจารณาค่า AIC น้อยที่สุดเป็นแบบจำลองในการพยากรณ์

⁵ pmdarima: ARIMA estimators for Python, © Copyright 2017-2021, Taylor G Smith

จาก https://alkaline-ml.com/pmdarima/tips_and_tricks.html

```

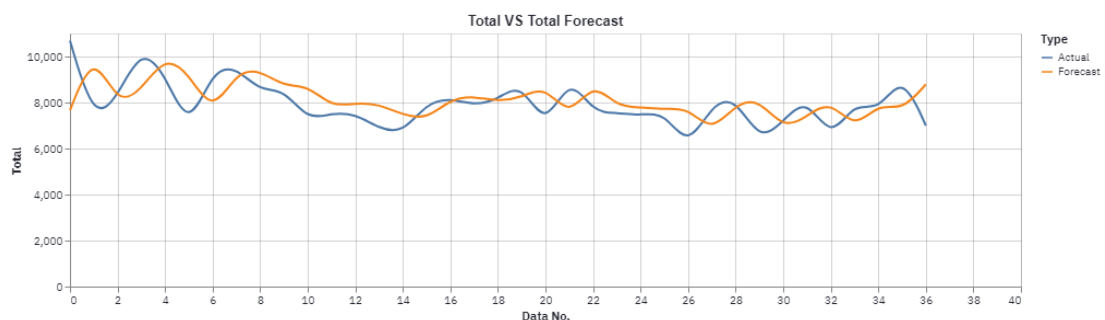
Performing stepwise search to minimize aic
ARIMA(0,1,0)(0,1,1)[12] : AIC=734.923, Time=0.04 sec
ARIMA(0,1,0)(0,1,0)[12] : AIC=747.992, Time=0.01 sec
ARIMA(1,1,0)(1,1,0)[12] : AIC=745.017, Time=0.03 sec
ARIMA(0,1,1)(0,1,1)[12] : AIC=738.959, Time=0.05 sec
ARIMA(0,1,0)(1,1,1)[12] : AIC=736.904, Time=0.05 sec
ARIMA(0,1,0)(0,1,2)[12] : AIC=736.450, Time=0.12 sec
ARIMA(0,1,0)(1,1,0)[12] : AIC=747.011, Time=0.03 sec
ARIMA(0,1,0)(1,1,2)[12] : AIC=737.523, Time=0.11 sec
ARIMA(1,1,0)(0,1,1)[12] : AIC=738.079, Time=0.05 sec
ARIMA(1,1,1)(0,1,1)[12] : AIC=740.918, Time=0.08 sec
ARIMA(0,1,0)(0,1,1)[12] intercept : AIC=737.099, Time=0.04 sec

Best model: ARIMA(0,1,0)(0,1,1)[12]
Total fit time: 0.616 seconds
    
```

ภาพที่ 3.14 ตัวอย่างการหารูปแบบ ARIMA(p,d,q) SARIMA (P,D,Q)_T หรือ ARIMA(p,d,q)x(P,D,Q)_s จากการพิจารณาค่า AIC

หลังจากได้แบบจำลองที่เหมาะสมแล้วจะใช้ข้อมูล train ทั้งหมดเพื่อเข้าเรียนรู้สร้างแบบจำลองนำค่าพยากรณ์เทียบกับข้อมูล test ของแต่ละวิธี สามารถร่างกราฟเปรียบเทียบได้ดังตัวอย่างภาพที่ 3.15

ARIMA Forecast Line Plot



ภาพที่ 3.15 กราฟเทียบข้อมูลพยากรณ์กับข้อมูลจริง (ข้อมูลทดสอบแบบจำลอง) ของเทคนิค ARIMA

3.1.4.2 สถิติข้อมูล/การวิเคราะห์เชิงถดถอย (Statistic/Regression Analysis)

ข้อมูลจะถูกนำเข้าไปเพื่อสร้างแบบจำลองด้วยวิธี Decision Tree, Random Forest และ Multiple Linear Regression โดยแบ่งข้อมูลจำนวน 80% เป็นข้อมูลสำหรับสร้างแบบจำลองและข้อมูลจำนวน 20% ใช้เป็นข้อมูลทดสอบโมเดลที่สร้างขึ้นได้ใช้วิธีการ Split Test แบบสุ่มข้อมูลก่อนแบ่ง

หลังจากได้แบบจำลองที่เหมาะสมแล้วจะใช้ข้อมูล train ทั้งหมดเพื่อเข้าเรียนรู้สร้างแบบจำลองและพยากรณ์ค่าออกมาดังตัวอย่างตารางที่ 3.4 และวัดค่าความคลาดเคลื่อนจากข้อมูล test เทียบกับค่าที่พยากรณ์ได้ของแต่ละวิธีเพื่อเปรียบเทียบในลำดับถัดไป

ตารางที่ 3.4 ตัวอย่างค่าพยากรณ์ที่ได้จากเทคนิค Random Forest (ข้อมูลทดสอบแบบจำลอง)

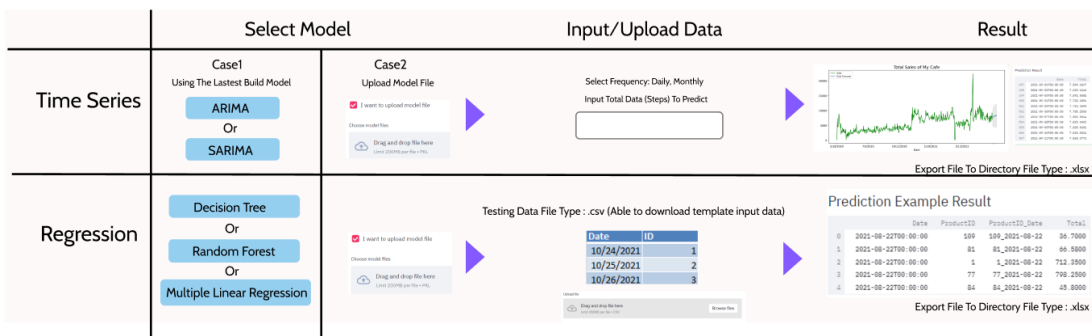
| ID_Date | Actual | Forecast |
|--------------|--------|----------|
| 1_2021-06-19 | 32 | 29.59 |
| 1_2021-07-16 | 67 | 94.02 |
| 1_2021-07-18 | 101 | 98.33 |
| 1_2021-08-29 | 264 | 274.96 |

3.1.5 การประเมินผล (Evaluation)

หลังจากที่แบ่งข้อมูลออกเป็น 2 ส่วนและใช้แบบจำลองที่สร้างได้จากข้อมูลจำนวน 80% มาทำการพยากรณ์ข้อมูลที่อยู่ในข้อมูลจำนวน 20% ที่เหลือนั้น และใช้ค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน (Mean Absolute Percentage Error, MAPE) แสดงค่าความคลาดเคลื่อนตามเปรียบเทียบด้วยเทคนิควิธีต่าง ๆ สามารถเลือกจัดเก็บแบบจำลองที่สนใจผ่านการ download ในรูปแบบของ .pkl ได้หรือ export ไฟล์ข้อมูลเปรียบเทียบดังตารางที่ 3.4 ในรูปแบบของ .xlsx

3.1.6 การปฏิบัติตามผลเสนอแนะ (Deployment)

Apply Model For New Dataset



ภาพที่ 3.16 ขั้นตอนการทำงานของระบบเพื่อพยากรณ์ข้อมูลใหม่

ขั้นตอนนี้หลังจากที่วิจัยและเปรียบเทียบค่าความคาดเคลื่อนของแต่ละเทคนิคแล้วแบบจำลองใดเหมาะสมกับความต้องการพยากรณ์ข้อมูลในอนาคตสามารถเลือกใช้แบบจำลองได้ออกเป็น 2 รูปแบบคือ

3.1.6.1 เลือกใช้แบบจำลองที่ทำการสร้างล่าสุดในระบบ

รูปแบบการพยากรณ์นี้ระบบจะทำการจัดเก็บแบบจำลองล่าสุดอัตโนมัติผู้ใช้งานสามารถเลือกวิธีการพยากรณ์และทำการระบุคาบเวลาในอนาคตที่ต้องการพยากรณ์สำหรับเทคนิค ARIMA, SARIMA หรือ upload ข้อมูลสำหรับทดสอบสำหรับเทคนิค Decision Tree, Random Forest, Multiple Linear Regression ที่ต้องการพยากรณ์ข้อมูลในอนาคตระบบจะแสดงผลลัพธ์ออกมาในรูปแบบของกราฟและ/หรือตาราง สามารถ export ให้อยู่ในรูปแบบของ .xlsx ไฟล์ ให้สำหรับนำไปใช้งาน

3.1.6.2 เลือกใช้แบบจำลองที่ได้จัดเก็บไว้มาใช้พยากรณ์ข้อมูลใหม่

รูปแบบการพยากรณ์นี้ผู้ใช้งานสามารถ upload แบบจำลองที่เคยจัดเก็บไว้นำมาพยากรณ์ข้อมูลได้โดยทำการระบุคาบเวลาในอนาคตที่ต้องการพยากรณ์สำหรับเทคนิค ARIMA, SARIMA หรือ upload ข้อมูลสำหรับทดสอบสำหรับเทคนิค Decision Tree, Random Forest, Multiple Linear Regression ที่ต้องการพยากรณ์ข้อมูลในอนาคตระบบจะแสดงผลลัพธ์ออกมาในรูปแบบของกราฟและ/หรือตารางสามารถ export ให้อยู่ในรูปแบบของ .xlsx ไฟล์ ให้สำหรับนำไปใช้งาน

3.2 เครื่องมือที่ใช้ในการวิจัย

3.2.1 Python

เป็นภาษาที่ได้รับความนิยมมากที่สุดสำหรับการทำ Machine Learning, Deep Learning และ Data Science ในตอนนี้เพราะมีไวยากรณ์ที่เข้าใจง่ายไม่ซับซ้อน มี library มากมายให้ใช้ และมี Community ขนาดใหญ่สำหรับแลกเปลี่ยนความรู้และสอบถามได้เป็นอย่างดี

3.2.2 Library ที่ใช้สำหรับงานวิจัยนี้

Streamlit เป็น library ที่นำมาใช้พัฒนา Web Application เนื่องจากเป็นการใช้ภาษา Python ทำให้สามารถทำงานกับ Object อื่นๆ ได้เช่น Pandas, Matplotlib, NumPy เป็นต้นไม่จำเป็นต้องใช้ Library API อื่น (เช่น Flask, Django) เป็นตัวช่วยในการติดต่อ กับ HTML

Pandas เป็น library สำหรับจัดการข้อมูลต่างๆ เช่น การโหลดข้อมูลไฟล์ CSV เข้ามาแล้วแสดงข้อมูลให้ออกมาในรูปแบบ Row กับ Column คล้าย Table เราเรียกสิ่งนี้ว่า Data Frame ทำให้สามารถใช้งานข้อมูลเหล่านั้นได้ง่ายและนำไปใช้งานต่อได้อย่างหลากหลาย

NumPy เป็น library ที่ใช้ในการคำนวณทางคณิตศาสตร์ในภาษา Python ซึ่งภายในถูกเขียนด้วยภาษา C จึงทำงานได้เร็วและมีประสิทธิภาพโดย NumPy โดยทั่วไปจะเกี่ยวกับการจัดการข้อมูล Array ขนาดใหญ่และเมทริกซ์ขนาดใหญ่และหลายมิติ

Matplotlib คือ library ที่ใช้สำหรับวาดกราฟที่ใช้งานได้หลากหลายยืดหยุ่น แสดงผลในรูปแบบ visualization เช่น กราฟเส้น, แผนภูมิแท่ง, แผนภูมิวงกลม เป็นต้น

Statsmodels คือ library ที่ช่วยให้สามารถใช้ Algorithm ที่สำคัญ ๆ อย่าง ANOVA และ ARIMA ที่ Machine Learning Library มาตรฐาน อย่าง Sci-kit Learn ไม่มีแล้วยังมีการแสดงรายละเอียด เช่น Features และ Metrics ของ Data ที่เป็นประโยชน์อีกมากมาย เช่น การใช้ Seasonal-Trend Decomposition

Scikit-learn คือ library สำหรับ Machine Learning ที่ได้รับความนิยมสุดๆ โดยมี algorithm ต่างๆ ทาง Machine Learning ให้ใช้งานง่ายและครบครัน เช่น Support Vector Machines (SVMs), Decision Tree, Random Forests, Gradient Boosting, Linear Regression เป็นต้น

Pmdarima เป็นโปรแกรมเสริม Python & Cython ของ Library ทางสถิติและการเรียนรู้ของเครื่อง (statsmodels และ scikit-learn) และดำเนินการการสร้างแบบจำลองกระบวนการ ARIMA ข้อดีของ Library นี้คือสามารถหาพารามิเตอร์ ARIMA (p, d, q) ที่ดีที่สุดสำหรับให้ได้เลย ในขณะที่ statsmodels จำเป็นจะต้องค้นหาพารามิเตอร์ที่เหมาะสมด้วยตนเอง

3.2.3 Anaconda

เป็น open-source ที่จะช่วยให้เราพัฒนาโปรแกรมทางด้าน data science และ machine learning ได้ดียิ่งขึ้นมีการสร้าง Environment สำหรับการใส่ library ที่แตกต่างกันได้ง่าย รวมถึงมีการรวบรวมเครื่องมือต่างๆ ไว้ เช่น Jupyter Notebook, Code editor ที่ใช้สำหรับการเขียน python ที่นิยมก็จะมี VS Code, Spyder หรือ R Studio สำหรับการเขียนภาษา R เป็นต้น

บทที่ 4

ผลการศึกษา

การวิจัยครั้งนี้จะแบ่งผลการวัดประสิทธิภาพและความถูกต้องออกเป็น 2 ส่วนคือ การวัดผลประสิทธิภาพความถูกต้องของแบบจำลองและการวัดผลประสิทธิภาพความถูกต้องจากการนำแบบจำลองพยากรณ์ในอนาคตเทียบกับข้อมูลจริงที่เกิดขึ้นในอนาคต

4.1 ผลการวัดประสิทธิภาพความถูกต้องของแบบจำลอง

งานวิจัยนี้ใช้ข้อมูลสร้างแบบจำลองแบ่งตามกรณีศึกษา 3 กรณีคือ

1. กรณีศึกษาเป็นประเภทข้อมูลยอดขายของร้านกาแฟ ราชวันตั้งแต่เดือน มีนาคม ถึง สิงหาคม พ.ศ. 2564
2. กรณีศึกษาเป็นประเภทข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม ราชเดือนตั้งแต่เดือน สิงหาคม พ.ศ. 2559 ถึง ธันวาคม พ.ศ. 2563
3. กรณีศึกษาเป็นประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3 ราชวัน ตั้งแต่เดือน มิถุนายน ถึง กันยายน พ.ศ. 2564

นำข้อมูลแบบจำลองที่สร้างได้จากข้อมูลข้างต้นจำนวน 80% มาทำการพยากรณ์ ส่วนที่อยู่ในข้อมูลจำนวน 20% ที่เหลือนั้นใช้ค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน (Mean Absolute Percentage Error, MAPE) วัดผลค่าความคลาดเคลื่อนเพื่อเปรียบเทียบกับเทคนิคต่างๆ ได้ดังนี้

ตารางที่ 4.1 เปรียบเทียบค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของการพยากรณ์ด้วย
กรณีศึกษาและเทคนิคในรูปแบบต่าง ๆ

| ค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน (Mean Absolute Percentage Error, MAPE) | | | | | |
|---|--------------------------------------|--------|--|------------------|----------------------------------|
| เทคนิค | รูปแบบที่ 1: เทคนิค ARIMA, SARIMA | | รูปแบบที่ 2: ใช้กับเทคนิค Multiple Linear Regression, Decision Tree Random Forest | | |
| | ARIMA | SARIMA | Decision Tree | Random Forest | Multiple Linear Regression |
| กรณีศึกษาเป็นประเภท ข้อมูลยอดขายของร้าน กาแฟ | 14.22 | 14.22 | 57.54 | 54.51 | 156.86 |
| กรณีศึกษาเป็นประเภท ข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้ามหานคร สายฉลองรัชธรรม | 35.15 | 25.90 | | | |
| กรณีศึกษาเป็นประเภท ข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3 | | | 22.45 | 18.78 | 28.08 |

วิเคราะห์เปรียบเทียบค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของแบบจำลองทาง
อนุกรมเวลาพบว่าข้อมูลที่ใช้สร้างแบบจำลองด้วยเทคนิค ARIMA และ SARIMA ในครั้งนี้หากไม่
มีแนวโน้มและอิทธิพลของฤดูกาลจะทำให้ผลพยากรณ์ที่ได้ มีค่าคลาดเคลื่อนไม่ต่างกันซึ่งต่างจาก
กรณีข้อมูลที่มีอิทธิพลของฤดูกาลแบบจำลองด้วยเทคนิค SARIMA จะให้ผลค่าความคลาดเคลื่อน
น้อยกว่า

กรณีศึกษาประเภทข้อมูลยอดขายของร้านกาแฟมีจุดประสงค์เพื่อพยากรณ์ค่ายอดขายรวมรายวันจากค่าความคลาดเคลื่อนดังตารางที่ 4.1 แสดงให้เห็นว่าสามารถเลือกใช้เพียงเทคนิค ARIMA ก็เพียงพอสำหรับการพยากรณ์ข้อมูลชุดใหม่

กรณีศึกษาเป็นประเภทข้อมูลจำนวนผู้ใช้บริการโครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม มีจุดประสงค์เพื่อพยากรณ์จำนวนผู้โดยสารรวมรายเดือนจากค่าความคลาดเคลื่อนดังตารางที่ 4.1 แสดงให้เห็นว่าการใช้เทคนิค SARIMA ในการสร้างแบบจำลอง มีค่าความคลาดเคลื่อนน้อยกว่าซึ่งเหมาะสำหรับนำไปใช้กับการพยากรณ์ข้อมูลชุดใหม่

ส่วนวิเคราะห์เปรียบเทียบค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของแบบจำลองทางอนุกรมเวลาพบว่าข้อมูลที่ใช้สร้างแบบจำลองทางสถิติข้อมูล/การวิเคราะห์เชิงถดถอย (Decision Tree, Random Forest, Multiple Linear Regression) สำหรับกรณีศึกษาประเภทข้อมูลยอดขายของร้านกาแฟและกรณีศึกษาประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3 พบว่าจากค่าความคลาดเคลื่อนมีลักษณะไปในทิศทางเดียวกันคือเทคนิค Random Forest ให้ค่าความคลาดเคลื่อนน้อยที่สุดรองลงมาคือ Decision Tree และ Multiple Linear Regression ตามลำดับดังตารางที่ 4.1

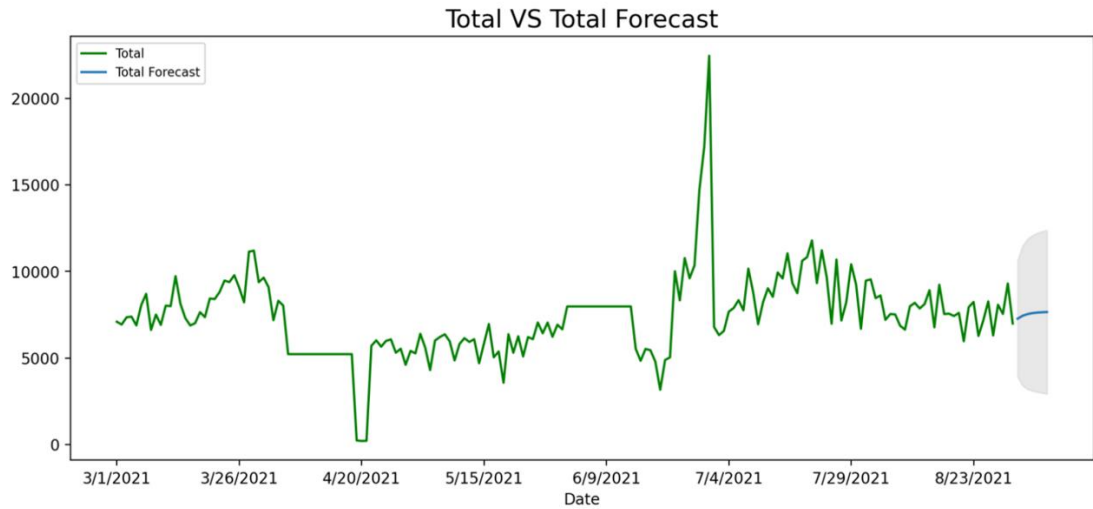
4.2 การวัดผลประสิทธิภาพความถูกต้องจากการนำแบบจำลองพยากรณ์ในอนาคตเทียบกับข้อมูลจริง

จากที่ได้ทำการวัดผลประสิทธิภาพความถูกต้องของแบบจำลองแล้วจากข้อ 4.1 ขั้นตอนนี้จะเป็นการเลือกแบบจำลองที่เหมาะสมมาใช้พยากรณ์ข้อมูลอนาคต

4.2.1 กรณีศึกษาประเภทข้อมูลยอดขายของร้านกาแฟ

4.2.1.1 ตัวอย่างการนำแบบจำลองพยากรณ์ไปใช้พยากรณ์กับข้อมูลยอดขายรวมรายวัน

เลือกใช้เทคนิค ARIMA จากเทคนิค ARIMA, SARIMA สำหรับการสร้างแบบจำลองและพยากรณ์ข้อมูลอนาคต โดยการใส่จำนวนคาบเวลาของข้อมูลที่ต้องการพยากรณ์ในอนาคตระบบ เช่น ในกรณีนี้แทนค่าด้วยคาบเวลาเท่ากับ 7 วัน (วันที่ 1-7 กันยายน พ.ศ. 2564) โดยมีระบบแสดงผลการพยากรณ์ปรากฏดังภาพที่ 4.1 และ Export ข้อมูลพยากรณ์ให้ในรูปแบบของไฟล์ .xlsx ดังตัวอย่างตารางที่ 4.2



ภาพที่ 4.1 กราฟแสดงผลการพยากรณ์ข้อมูลยอดขายรวมรายวัน (1-7 กันยายน 2564) ด้วยเทคนิค ARIMA

ตารางที่ 4.2 ตัวอย่างข้อมูลพยากรณ์ยอดขายรวมรายวันด้วยเทคนิค ARIMA เมื่อกำหนดคาบเวลาเท่ากับ 7 สำหรับ Export ไฟล์ในรูปแบบ.xlsx

| Future Step Forecast. | Total |
|-----------------------|-----------|
| 1 | 7,276.209 |
| 2 | 7,441.324 |
| 3 | 7,538.273 |
| 4 | 7,595.198 |
| 5 | 7,628.622 |
| 6 | 7,648.247 |
| 7 | 7,659.77 |

เทียบกับข้อมูลจริงของยอดขายรวมรายวันของดังภาพที่ 4.2

| วันที่ | ยอดก่อนลด | ส่วนลด | ส่วนลดภาษี | ยอดรวม ๑ | ค่าบริการ | ยอดก่อนภาษี ๒ | ภาษี | ลดภาษีขา | ยอดปดเศษ | รวมสุทธิ ๒ |
|------------|-----------|--------|------------|-----------|-----------|---------------|----------|----------|----------|------------|
| 01/09/2021 | 6,075 | 0 | 0 | 6,075 | 0 | 5,677.53 | 397.47 | 0 | 0 | 6,075 |
| 02/09/2021 | 6,385 | 0 | 0 | 6,385 | 0 | 5,967.28 | 417.72 | 0 | 0 | 6,385 |
| 03/09/2021 | 6,665 | 0 | 0 | 6,665 | 0 | 6,228.95 | 436.05 | 0 | 0 | 6,665 |
| 04/09/2021 | 5,530 | 0 | 0 | 5,530 | 0 | 5,168.24 | 361.76 | 114.75 | 2.75 | 5,418 |
| 05/09/2021 | 5,990 | 39 | 0 | 5,951 | 0 | 5,561.66 | 389.34 | 219.75 | 5.75 | 5,737 |
| 06/09/2021 | 6,060 | 0 | 0 | 6,060 | 0 | 5,663.57 | 396.43 | 274.5 | 6.5 | 5,792 |
| 07/09/2021 | 4,730 | 0 | 0 | 4,730 | 0 | 4,420.52 | 309.48 | 191.25 | 6.25 | 4,545 |
| Total | 41,435.00 | 39.00 | 0.00 | 41,396.00 | 0.00 | 38,687.75 | 2,708.25 | 800.25 | 21.25 | 40,617.00 |

ภาพที่ 4.2 ข้อมูลจริงของยอดขายรวมรายวัน วันที่ 1-7 กันยายน 2564

เมื่อหาค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของข้อมูลพยากรณ์เทียบกับข้อมูลจริง 7 วันได้ (อ้างอิงตารางที่ 4.2 และภาพที่ 4.2) เท่ากับ 39.23%

4.2.1.2 ตัวอย่างการนำแบบจำลองพยากรณ์ไปใช้พยากรณ์กับข้อมูลยอดขายรวมรายวัน แบ่งตามรหัสสินค้า

เลือกใช้เทคนิค Random Forest จากเทคนิค (Decision Tree, Random Forest, Multiple Regression) สำหรับการสร้างแบบจำลองและพยากรณ์ข้อมูลอนาคตโดยการ upload ข้อมูลที่ต้องการพยากรณ์ดังตัวอย่างข้อมูลตามตารางที่ 4.3 ระบบแสดงข้อมูลจริงที่ทำการแปลงข้อมูลผ่านระบบดังภาพที่ 4.3 และแสดงผลการพยากรณ์และ Export ข้อมูลพยากรณ์ไว้ในรูปแบบของไฟล์ .xlsx ปรากฏดังภาพที่ 4.4

ตารางที่ 4.3 ตัวอย่างข้อมูล Test สำหรับการพยากรณ์ข้อมูลใหม่ ด้วยเทคนิค Random Forest

| Date | Temperature | IS Holiday | Promotion | ID | UnitPrice |
|----------|-------------|------------|-----------|-----|-----------|
| 9/1/2021 | 30 | FALSE | TRUE | 77 | 40 |
| 9/1/2021 | 30 | FALSE | TRUE | 141 | 50 |
| 9/1/2021 | 30 | FALSE | TRUE | 152 | 60 |
| 9/1/2021 | 30 | FALSE | TRUE | 158 | 55 |
| 9/1/2021 | 30 | FALSE | TRUE | 160 | 65 |
| 9/1/2021 | 30 | FALSE | TRUE | 111 | 45 |
| 9/1/2021 | 30 | FALSE | TRUE | 102 | 60 |

| Temperature | IS Holiday | Promotion | ID | UnitPrice | Date_dayofweek | Date_month | Date_year | Date_day |
|-------------|------------|-----------|-----|-----------|----------------|------------|-----------|----------|
| 30 | 0 | 1 | 77 | 40 | 2 | 9 | 2021 | 1 |
| 30 | 0 | 1 | 141 | 50 | 2 | 9 | 2021 | 1 |
| 30 | 0 | 1 | 152 | 60 | 2 | 9 | 2021 | 1 |
| 30 | 0 | 1 | 158 | 55 | 2 | 9 | 2021 | 1 |
| 30 | 0 | 1 | 160 | 65 | 2 | 9 | 2021 | 1 |
| 30 | 0 | 1 | 111 | 45 | 2 | 9 | 2021 | 1 |
| 30 | 0 | 1 | 102 | 60 | 2 | 9 | 2021 | 1 |

ภาพที่ 4.3 รูปแบบของข้อมูลจริงก่อนนำเข้าพยากรณ์ด้วย Random Forest กรณีศึกษาข้อมูลประเภทยอดขายของร้านกาแฟแบ่งตามรหัสสินค้า

Prediction Example Result

| | Date | ID | ID_Date | Total |
|---|---------------------|-----|----------------|----------|
| 0 | 2021-09-01T00:00:00 | 77 | 77_2021-09-01 | 675.0000 |
| 1 | 2021-09-01T00:00:00 | 141 | 141_2021-09-01 | 506.8333 |
| 2 | 2021-09-01T00:00:00 | 152 | 152_2021-09-01 | 478.2500 |
| 3 | 2021-09-01T00:00:00 | 158 | 158_2021-09-01 | 491.1500 |
| 4 | 2021-09-01T00:00:00 | 160 | 160_2021-09-01 | 285.8000 |
| 5 | 2021-09-01T00:00:00 | 111 | 111_2021-09-01 | 435.5500 |
| 6 | 2021-09-01T00:00:00 | 102 | 102_2021-09-01 | 636.3000 |

ภาพที่ 4.4 ผลการพยากรณ์ข้อมูลยอดขายรวมรายรหัสสินค้าของวันที่ 1 กันยายน พ.ศ. 2564 ด้วยเทคนิค Random Forest สำหรับ Export ไฟล์ในรูปแบบ.xlsx

เทียบกับข้อมูลจริงของยอดขายรวมรายรหัสสินค้าของวันที่ 1 กันยายน พ.ศ. 2564 ดัง

ภาพที่ 4.5

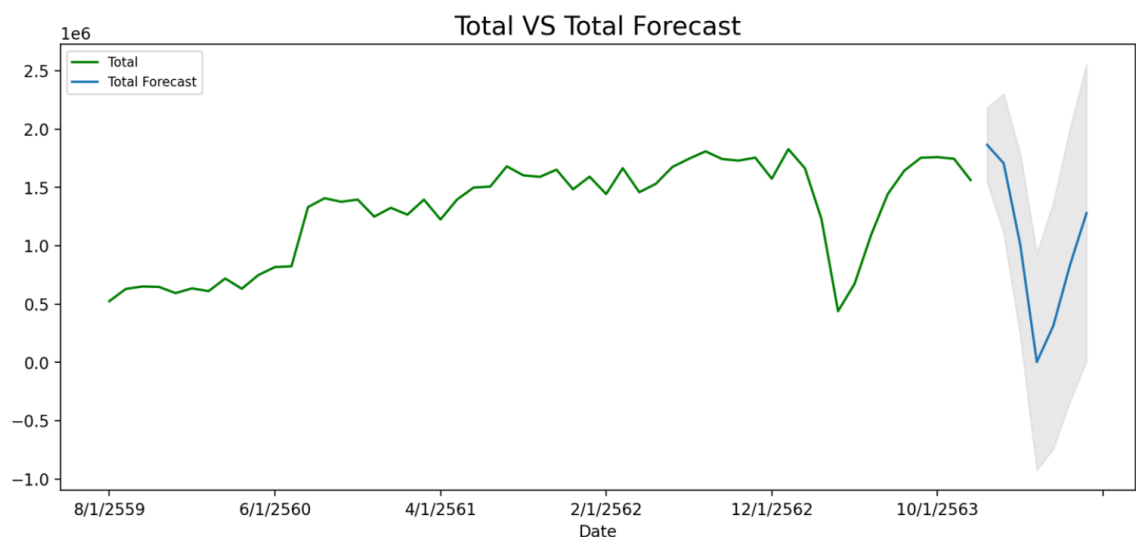
| รหัสสินค้า | ชื่อสินค้า | กลุ่ม | หมวดสินค้า | ต้นทุนเฉลี่ย | ราคาขายเฉลี่ย | กำไรเฉลี่ย | จำนวนการขาย | ยอดก่อนลด | ต้นทุน | ส่วนลด | ราคาสุทธิ |
|------------|-------------------------|-------|---------------|--------------|---------------|------------|-------------|-----------|--------|--------|-----------|
| 77 | เอส-เย็น | | ไทย กาแฟ | 0 | 41.43 | 41.43 | 14 | 580 | 0 | 0 | 580 |
| 141 | พรีเมียม อเมริกาโน่เย็น | | พรีเมียม เย็น | 0 | 60 | 60 | 7 | 420 | 0 | 0 | 420 |
| 152 | ลาเต้ | | ไทย กาแฟ | 0 | 41 | 41 | 10 | 410 | 0 | 0 | 410 |
| 158 | อเมริกาโน่ | | ไทย กาแฟ | 0 | 40 | 40 | 10 | 400 | 0 | 0 | 400 |
| 160 | อเมริกาโน่มีฟอง | | ไทย กาแฟ | 0 | 50 | 50 | 7 | 350 | 0 | 0 | 350 |
| 111 | ชาไทยนมสด | | ไทย ชา | 0 | 36.67 | 36.67 | 9 | 330 | 0 | 0 | 330 |
| 102 | คาปูชิโน่ | | ไทย กาแฟ | 0 | 41.43 | 41.43 | 7 | 290 | 0 | 0 | 290 |

ภาพที่ 4.5 ข้อมูลจริงของยอดขายรวมรายรหัสสินค้าวันที่ 1 กันยายน 2564

เมื่อหาค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของข้อมูลพยากรณ์เทียบกับข้อมูลจริงของวันที่ 1 กันยายน 2564 จากตัวอย่างสินค้า 7 ตัวอย่างพบว่าได้ค่าความคลาดเคลื่อน (อ้างอิงภาพที่ 4.4 และภาพที่ 4.5) เท่ากับ 35.18%

4.2.2 กรณีศึกษาประเภทข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม

เลือกใช้เทคนิค SARIMA จากเทคนิค ARIMA, SARIMA สำหรับการสร้างแบบจำลองและพยากรณ์ข้อมูลอนาคตโดยการใส่จำนวนคาบเวลาของข้อมูลที่ต้องการพยากรณ์ในอนาคต ระบบ เช่น ในกรณีนี้แทนค่าด้วยคาบเวลาเท่ากับ 7 เดือน (เดือน มกราคม ถึง กรกฎาคม พ.ศ. 2564) โดยมีระบบแสดงผลการพยากรณ์ปรากฏดังภาพที่ 4.6 และ Export ข้อมูลพยากรณ์ให้ในรูปแบบของไฟล์ .xlsx ดังตัวอย่างตารางที่ 4.4



ภาพที่ 4.6 กราฟแสดงผลการพยากรณ์ข้อมูลจำนวนผู้ใช้บริการ (เดือน มกราคม ถึง เดือน กรกฎาคม พ.ศ. 2564) ด้วยเทคนิค SARIMA

ตารางที่ 4.4 ตัวอย่างข้อมูลพยากรณ์จำนวนผู้ใช้บริการรายเดือนด้วยเทคนิค SARIMA เมื่อกำหนดคาบเวลาเท่ากับ 7 สำหรับ Export ไฟล์ในรูปแบบ.xlsx

| Future Step Forecast. | Total |
|-----------------------|--------------|
| 1 | 1,864,581.63 |
| 2 | 1,706,402.79 |
| 3 | 1,008,853.91 |
| 4 | 3,758.38 |
| 5 | 315,255.02 |
| 6 | 833,198.08 |
| 7 | 1,279,183.94 |

กรณีนี้ผลลัพธ์ของจำนวนผู้ใช้บริการควรเป็นจำนวนเต็มก่อนนำมาใช้ให้ทำการปัดเศษเป็นจำนวนเต็มเพื่อนำไปเทียบกับข้อมูลจริงของจำนวนผู้ใช้บริการ ดังภาพที่ 4.7

| เดือน | ปี | โครงการ | จำนวนผู้โดยสารรวม |
|------------|------|--|-------------------|
| มกราคม | 2564 | โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม | 874,003 |
| กุมภาพันธ์ | 2564 | โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม | 1,081,904 |
| มีนาคม | 2564 | โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม | 1,519,467 |
| เมษายน | 2564 | โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม | 793,625 |
| พฤษภาคม | 2564 | โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม | 521,799 |
| มิถุนายน | 2564 | โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม | 677,959 |
| กรกฎาคม | 2564 | โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม | 466,159 |

ภาพที่ 4.7 ข้อมูลจริงของจำนวนผู้ใช้บริการรายเดือน ตั้งแต่เดือน มกราคม ถึง กรกฎาคม 2564

เมื่อหาค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของข้อมูลพยากรณ์เทียบกับข้อมูลจริง 7 เดือนได้ (อ้างอิงตารางที่ 4.4 และภาพที่ 4.7) เท่ากับ 77.30%

4.2.3 กรณีศึกษาเป็นประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3

เลือกใช้เทคนิค Random Forest จากเทคนิค (Decision Tree, Random Forest, Multiple Regression) สำหรับการสร้างแบบจำลองและพยากรณ์ข้อมูลอนาคตโดยการ upload ข้อมูลที่ต้องการพยากรณ์ดังตัวอย่างข้อมูลตามตารางที่ 4.5 ระบบจะทำการแปลงข้อมูลให้ดังภาพที่ 4.8 หลังจากนำข้อมูลเข้าระบบจะแสดงผลการพยากรณ์ประมวลผลด้วยแบบจำลองข้างต้นและสามารถ Export ข้อมูลพยากรณ์ให้ในรูปแบบของไฟล์ .xlsx ปรากฏดังภาพที่ 4.9

ตารางที่ 4.5 ตัวอย่างข้อมูล Test สำหรับการพยากรณ์ข้อมูลใหม่ ด้วยเทคนิค Random Forest

| ID | Date | new_case | new_recovered |
|----|-----------|----------|---------------|
| 1 | 10/1/2021 | 11754 | 12473 |
| 1 | 10/2/2021 | 11375 | 13127 |
| 1 | 10/3/2021 | 10828 | 11894 |
| 1 | 10/4/2021 | 9930 | 12336 |
| 1 | 10/5/2021 | 9869 | 11152 |
| 1 | 10/6/2021 | 9866 | 10115 |

| ID | new_case | new_recovered | Date_dayofweek | Date_month | Date_year | Date_day |
|----|----------|---------------|----------------|------------|-----------|----------|
| 1 | 11754 | 12473 | 4 | 10 | 2021 | 1 |
| 1 | 11375 | 13127 | 5 | 10 | 2021 | 2 |
| 1 | 10828 | 11894 | 6 | 10 | 2021 | 3 |
| 1 | 9930 | 12336 | 0 | 10 | 2021 | 4 |
| 1 | 9869 | 11152 | 1 | 10 | 2021 | 5 |
| 1 | 9866 | 10115 | 2 | 10 | 2021 | 6 |

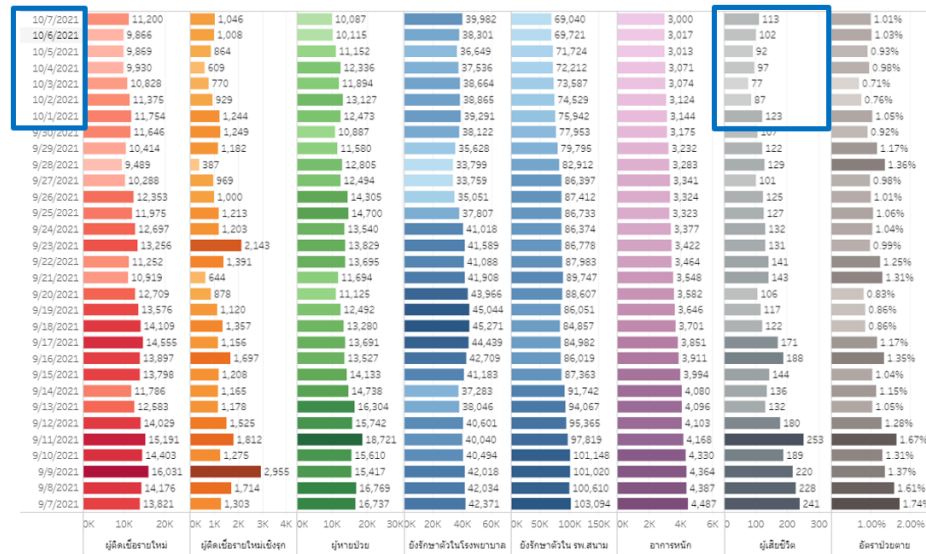
ภาพที่ 4.8 รูปแบบของข้อมูลจริงก่อนนำเข้าพยากรณ์ด้วย Random Forest กรณีศึกษาข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3

Prediction Example Result

| | Date | ID | ID_Date | Total |
|---|---------------------|----|--------------|----------|
| 0 | 2021-10-01T00:00:00 | 1 | 1_2021-10-01 | 132.0900 |
| 1 | 2021-10-02T00:00:00 | 1 | 1_2021-10-02 | 136.2400 |
| 2 | 2021-10-03T00:00:00 | 1 | 1_2021-10-03 | 131.8100 |
| 3 | 2021-10-04T00:00:00 | 1 | 1_2021-10-04 | 124.8400 |
| 4 | 2021-10-05T00:00:00 | 1 | 1_2021-10-05 | 127.0400 |
| 5 | 2021-10-06T00:00:00 | 1 | 1_2021-10-06 | 129.2400 |

ภาพที่ 4.9 ผลการพยากรณ์ข้อมูลผู้เสียชีวิตรายใหม่จากสถานการณ์ COVID-19 ระลอก 3 ด้วยเทคนิค Random Forest สำหรับ Export ไฟล์ในรูปแบบ.xlsx

เทียบกับข้อมูลจริงของเคส 30 วันย้อนหลังของรายงานสถานการณ์ COVID-19 ของวันที่ 1-6 ตุลาคม พ.ศ. 2564 ดังภาพที่ 4.10



ภาพที่ 4.10 ข้อมูลจริงของเคส 30 วันย้อนหลังของรายงานสถานการณ์ COVID-19

ที่มา: สรุปข้อมูลเคส 30 วันย้อนหลัง กรมควบคุมโรค กระทรวงสาธารณสุข

เมื่อหาค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อนของข้อมูลพยากรณ์เทียบกับข้อมูลจริง
ของจำนวนผู้เสียชีวิตของรายงานสถานการณ์ COVID-19 วันที่ 1-6 ตุลาคม 2564 พบว่าได้ค่าความ
คลาดเคลื่อน (อ้างอิงภาพที่ 4.9 และภาพที่ 4.10) เท่ากับ 37.69%

บทที่ 5

บทสรุปและข้อเสนอแนะ

งานวิจัยนี้ได้นำเสนอเกี่ยวกับการพัฒนาระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง โดยผู้วิจัยได้วัดประสิทธิภาพของแต่ละตัวแบบจำลองแต่ละเทคนิค เพื่อนำผลลัพธ์มาเปรียบเทียบและสามารถสรุปผลการวิจัยได้ดังนี้

5.1 สรุปผลการวิจัย

การพัฒนาระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่องสามารถประยุกต์ใช้การวิเคราะห์อนุกรมเวลาในทางปฏิบัติหลายอย่างเช่นกรณีศึกษาที่งานวิจัยนี้นำมาประยุกต์ใช้ คือ

1. กรณีศึกษาเป็นประเภทข้อมูลยอดขายของร้านกาแฟรายวันตั้งแต่เดือน มีนาคม ถึง สิงหาคม พ.ศ. 2564
2. กรณีศึกษาเป็นประเภทข้อมูลจำนวนผู้ใช้บริการ โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรมรายเดือนตั้งแต่เดือน สิงหาคม พ.ศ. 2559 ถึง ธันวาคม พ.ศ. 2563
3. กรณีศึกษาเป็นประเภทข้อมูลรายงานสถานการณ์ COVID-19 ระลอก 3 รายวันตั้งแต่เดือน มิถุนายน ถึง กันยายน พ.ศ. 2564

โดยใช้เทคนิคที่หลากหลายได้แก่ เทคนิค ARIMA, SARIMA, Decision Tree, Random Forest และ Multiple Linear Regression

จากผลการทดลองพิจารณาประสิทธิภาพของแบบจำลองเปรียบเทียบเทคนิค ARIMA กับ SARIMA เมื่อนำข้อมูลกรณีศึกษาของยอดขายรวม (รายวัน) ของร้านกาแฟและจำนวนผู้ใช้บริการ (รายเดือน) โครงการรถไฟฟ้าฟ้ามหานคร สายฉลองรัชธรรม พบว่าหากข้อมูลมีการเคลื่อนไหวแนวโน้มที่ไม่เสถียรและไม่มีอิทธิพลตามฤดูกาลผลของการพยากรณ์จะให้ค่า MAPE ใกล้เคียงกันหรือแทบไม่ต่างกันเลย เนื่องจากพารามิเตอร์ที่ทำผลของเกณฑ์ในการคัดเลือกแบบจำลอง (AIC ต่ำที่สุด) สำหรับเทคนิคทั้งสองข้างต้นมีค่าเดียวกัน ในทางตรงกันข้ามหากข้อมูลมีแนวโน้มอิทธิพลตามฤดูกาลผลของการพยากรณ์ของเทคนิค SARIMA จะให้ผลดีกว่าโดยให้ค่า MAPE ต่ำกว่าเทคนิค ARIMA และเมื่อพิจารณาประสิทธิภาพของแบบจำลองเปรียบเทียบเทคนิค

Decision Tree, Random Forest และ Multiple Linear Regression เมื่อนำข้อมูลกรณีศึกษาของ ยอดขายรวมของแต่ละรายสินค้า (รายวัน) ของร้านกาแฟและจำนวนผู้ป่วยและผู้เสียชีวิตจาก สถานการณ์ COVID-19 ระลอก 3 (รายวัน) พบว่าเทคนิค Random Forest พยากรณ์ยอดขายตาม รายการสินค้าและจำนวนผู้เสียชีวิตรายใหม่ของสถานการณ์ COVID-19 ให้ผลค่าความคลื่อนไปใน ทิศทางเดียวกันโดยให้ค่า MAPE ต่ำที่สุดรองลงมาคือเทคนิค Decision Tree และ Multiple Linear Regression ตามลำดับ สอดคล้องกับผลการศึกษาของ Abdulwahed Salam and Abdelaziz El Hibaoui (2018) ซึ่งศึกษาเปรียบเทียบวิธีพยากรณ์ 5 วิธี Linear regression, Decision tree, Random forest, Feedforward Neural network และ Supper vector ของการพยากรณ์ปริมาณการใช้ไฟฟ้าของ เมืองเตโตวอนประเทศโมร็อกโกพบว่าวิธี Random Forest เป็นวิธีที่เหมาะสมที่สุด

5.2 อภิปรายผลการวิจัย

จากการตั้งสมมติฐานที่ว่าค่าการพยากรณ์ควรมีความใกล้เคียงและมีค่า MAPE ไม่น้อย กว่า 50% แต่เมื่อนำข้อมูลสร้างแบบจำลองและวัดประสิทธิภาพจากค่าความคลื่อนพบว่าค่า MAPE ที่ได้ส่วนใหญ่ให้ผลตามสมมติฐานที่ตั้งไว้ แต่ก็มีบางกรณีที่มีค่า MAPE สูงกว่าสมมติฐาน ผู้วิจัยได้ศึกษาเพิ่มเติมว่าอะไรเป็นอาจเป็นเหตุที่ทำให้ประสิทธิภาพของการพยากรณ์ไม่ดีเท่าที่ควร หลังจากได้ไปศึกษารายละเอียดข้อมูลและสภาพแวดล้อมหรือปัจจัยที่ส่งผลกระทบต่อของข้อมูล ได้ พบข้อสังเกตที่น่าสนใจ คือช่วงเวลาของข้อมูลที่นำมาวิจัยและเทียบผลการพยากรณ์นั้นมีความ ทิศทางที่ไม่เสถียรจากสถานการณ์ COVID-19 เช่น การขอความร่วมมือประชาชนงดออกจาก เลหสถาน หรือ การออกแนวทางปฏิบัติด้านสาธารณสุขเพื่อป้องกันการแพร่ระบาดของเชื้อ COVID-19 สำหรับห้างสรรพสินค้า ศูนย์การค้า คอมมูนิตี้มอลล์ ร้านอาหาร แหล่งท่องเที่ยว หรือ สถานที่ท่องเที่ยว เป็นต้น หรือแม้กระทั่งปัจจัยทางด้านสิ่งแวดล้อม สภาพภูมิอากาศ ตัวอย่างเช่น มี เหตุการณ์พายุเข้าไทยทำให้ฝนตกหนักช่วงเวลานั้นๆ ส่งผลให้คนงดออกเดินทางหรือออกมาจับจ่าย ใช้สอยซึ่งส่งผลกระทบต่อโดยตรงกับข้อมูลที่นำมาศึกษาวิจัย

5.3 ข้อเสนอแนะ

5.3.1 ข้อมูลที่นำมาศึกษาวิจัยควรมีความสมบูรณ์ถูกต้องและมีความน่าเชื่อถือ เช่น การเตรียม สร้างรูปแบบการบันทึกของข้อมูล เช่น การกำหนดรหัสสินค้าให้ตรงกับสินค้าและราคาขายเพื่อที่ ได้ข้อมูลที่สมบูรณ์ลดการเตรียมข้อมูลก่อนการนำเข้าสร้างแบบจำลอง

5.3.2 ควรแบ่งข้อมูลของคาบเวลามาสร้างแบบจำลองให้เหมาะสมตามจุดประสงค์ของการพยากรณ์ เช่น ใช้ปัจจัยพิจารณาาร่วมคือหากมีสถานการณ์ลือกคาวน หรือการขอความร่วมมือประชาชนงคออกจากเคหสถาน ที่ทำให้ยอคขาย หรือจำนวนผู้ให้บริการในช่วงเวลานั้นๆ ไม่ปกติ หากเรามีจุดประสงค์ว่าอยากทราบว่ถ้าช่วงเวลาปกติจะมียอคขายหรือจำนวนผู้ให้บริการประมาณเท่าไร ลองปรับเลือกข้อมูลยอคขายหรือจำนวนผู้ให้บริการช่วงหลังจากคลายลือกคาวน มาทำการพยากรณ์ หรืออาจจะใช้การวิเคราะห์เพิ่มเติมว่าคาบเวลาใดที่เหมาะสมกับการสร้างแบบจำลองเช่น 1 หรือ 3 เดือนก็อาจเพียงพอกับการพยากรณ์ข้อมูลรายวัน เป็นต้น

5.3.3 ควรมีการพิจารณาประเมินผลที่ได้จากการพิจารณาอย่างสม่ำเสมอ ปรับเพิ่มลดตัวแปรอิสระเพิ่มเติมที่อาจช่วยเพิ่มประสิทธิภาพของการสร้างแบบจำลองมากขึ้น หรือลองนำเทคนิคใหม่ๆ ที่ถูกพัฒนามาปรับใช้ให้มีประสิทธิภาพมากขึ้น

5.3.4 นอกจากตัวระบบที่ช่วยพยากรณ์เบื้องต้นแล้ว ควรมีการพิจารณาปัจจัยอื่นๆที่อาจส่งผลกับค่าพยากรณ์เพิ่มเติมร่วมด้วยเช่น ข้อมูลทางเศรษฐกิจ, แผนการพัฒนาผลิตภัณฑ์และรูปแบบทางการตลาดร่วมด้วย

บรรณานุกรม

บรรณานุกรม

ภาษาไทย

การพยากรณ์.(2564). สืบค้น 27 กันยายน 2564, จาก

http://119.46.166.126/digitalschool/p5/ma5_1/lesson5/content1/more/page2.php

ดร.สุทิน ชนะบุญ. (2560). สถิติการวิเคราะห์ข้อมูลในงานวิจัยเบื้องต้น บทที่ 6 การวิเคราะห์เชิง

อนุमान. R2R สำนักงานสาธารณสุขจังหวัดขอนแก่น ปี 2560. น.148-160

ชนัท จระณะสมบุรณ์. (2561). การทำนายการซื้อซ้ำของผู้ซื้อโดยใช้เทคนิคการเรียนรู้ของเครื่องจักร.

ปริญญาณิพนธ์ วิทยาศาสตร์มหาบัณฑิต สาขาวิชาเทคโนโลยีสารสนเทศ

มหาวิทยาลัยศรีนครินทรวิโรฒ.

นิฉา แก้วหาวงษ์. (2547). การศึกษาเปรียบเทียบการพยากรณ์ข้อมูลอนุกรมเวลาโดยวิธีการทำให้

เรียบแบบเอ็กซ์โปเนนเชียลและวิธีการของบ็อกซ์-เจนกินส์: กรณีศึกษาการพยากรณ์มูลค่า

การส่งออกข้าว ยางพารา และมันสำปะหลัง. วารสารวิทยาศาสตร์ มหาวิทยาลัยนเรศวร ปีที่

1 ฉบับที่ 2. (น.73-90).

ปณิธิ ชัมมวิจยะ. (2560) การพยากรณ์โรคโดยการวิเคราะห์ข้อมูลอนุกรมเวลา. อมรมพัฒนา

ศักยภาพบุคลากรด้านการพยากรณ์โรค กองยุทธศาสตร์และแผนงาน สำนักงาน

ปลัดกระทรวงสาธารณสุข. สืบค้น 30 กันยายน 2564, จาก

https://bps.moph.go.th/new_bps/sites/default/files/ime%20Series%20Analysis_2017.pdf

ปริมาณผู้โดยสารรถไฟฟ้ามหานคร. (2564). สืบค้น 5 ตุลาคม 2564, จาก

<https://data.go.th/dataset/mrta-crmk>

ปิยมาส กล้าแข็ง. (2561). การประยุกต์ใช้เทคนิคการพยากรณ์ เพื่อการจัดการสินค้าคงคลัง.

วิทยานิพนธ์ ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาการจัดการซัพพลายเชนธุรกิจ

มหาวิทยาลัยราชภัฏสวนสุนันทา.

ศศ.ดร. เณลิมพล. (2562). การพยากรณ์ทางอนุกรมเวลา (Time Series Forecasting). สืบค้น 30

กันยายน 2564, จาก [https://cj007blog.files.wordpress.com/2020/04/07-time-series-](https://cj007blog.files.wordpress.com/2020/04/07-time-series-forecasting.pdf)

[forecasting.pdf](https://cj007blog.files.wordpress.com/2020/04/07-time-series-forecasting.pdf)

- วรางคณา เรียนสุทธิ. (2562). การพยากรณ์ราคามะพร้าวด้วยวิธีบ็อกซ์-เจนกินส์.วารสารวิชาการ มทร.สุวรรณภูมิ (RMUTSB Acad. J.) ปีที่ 7. (น. 87 – 100). มหาวิทยาลัยทักษิณ.
- วันหยุดตามประเพณีของสถาบันการเงิน ประจำปี. (2564). สืบค้น 14 มิถุนายน 2564, จาก <https://www.bot.or.th/Thai/FinancialInstitutions/FIholiday/Pages/2021.aspx>
- สถานการณ์ผู้ติดเชื้อ COVID-19 อัปเดตรายวัน รูปแบบ API(Json/CSV Data Format). (2564). สืบค้น 28 กันยายน 2564, จาก <https://covid19.ddc.moph.go.th/>
- หทัยชนก นานานอก. (2553). ศึกษาการพยากรณ์ยอดขายเพื่อวางแผนการผลิตของสินค้าในอุตสาหกรรมขนาดเล็กแห่งหนึ่ง. ปัญหาพิเศษอุตสาหกรรมศาสตรมหาบัณฑิต. มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าพระนครเหนือ.
- Fanny Indradjaja. (2562). Quantitative Analysis. สืบค้น 27 กันยายน 2564, จาก <https://slideplayer.in.th/slide/16724291/>

ภาษาต่างประเทศ

- Akaike, H. (1973). Information Theory as an Extension of the Maximum Likelihood Principle. In: Petrov, B.N. and Csaki, F., Eds., Second International Symposium on Information Theory, Akademiai Kiado, Budapest, pp.267-281.
- A Salam, A El Hibaoui. (2018). Comparison of Machine Learning Algorithms for the Power Consumption Prediction : - Case Study of Tetouan city-. 2018 6th International Renewable and Sustainable Energy Conference (IRSEC). pp 1210
- Breiman, Leo; Friedman, J. H.; Olshen, R. A.; Stone, C. J. (1984). Classification and regression trees. Monterey, CA: Wadsworth & Brooks/Cole Advanced Books & Software. ISBN 978-0-412-04841-8
- B. Singh P. Kumar N. Sharma and K P Sharma. (2020). Sales Forecast for Amazon Sales with Time Series Modeling. 2020 First International Conference on Power, Control and Computing Technologies (ICPC2T). pp 38-43.
- Canova, F. and Hansen, Bruce E. (1995). Are seasonal patterns constant over time? A test for seasonal stability. Journal of Business & Economic Statistics, 13(3), pp. 237-252.

- Georgios D. (2019). Decision Tree Regressor explained in depth. Retrieved September 19, 2021, from <https://gdcoder.com/decision-tree-regressor-explained-in-depth/>
- Kullback, S., & Leibler, R. A. (1951). On information and sufficiency. *Annals of Mathematical Statistics*, 22, 79–86.
- Taylor G Smith. (2017). pmdarima: ARIMA estimators for Python. Retrieved August 3, 2021, from https://alkaline-ml.com/pmdarima/tips_and_tricks.html
- T. Tanizaki, T. Hoshino, T. Shimmura and T. Takenaka. (2019). Demand forecasting in restaurants using machine learning and statistical analysis. *Procedia CIRP*. vol. 79. pp. 679-683.
- Weather Underground. (2020). Bangkok, Thailand Weather History. Retrieved June 14, 2021 from <https://www.wunderground.com/history/monthly/th/bangkok>

ภาคผนวก

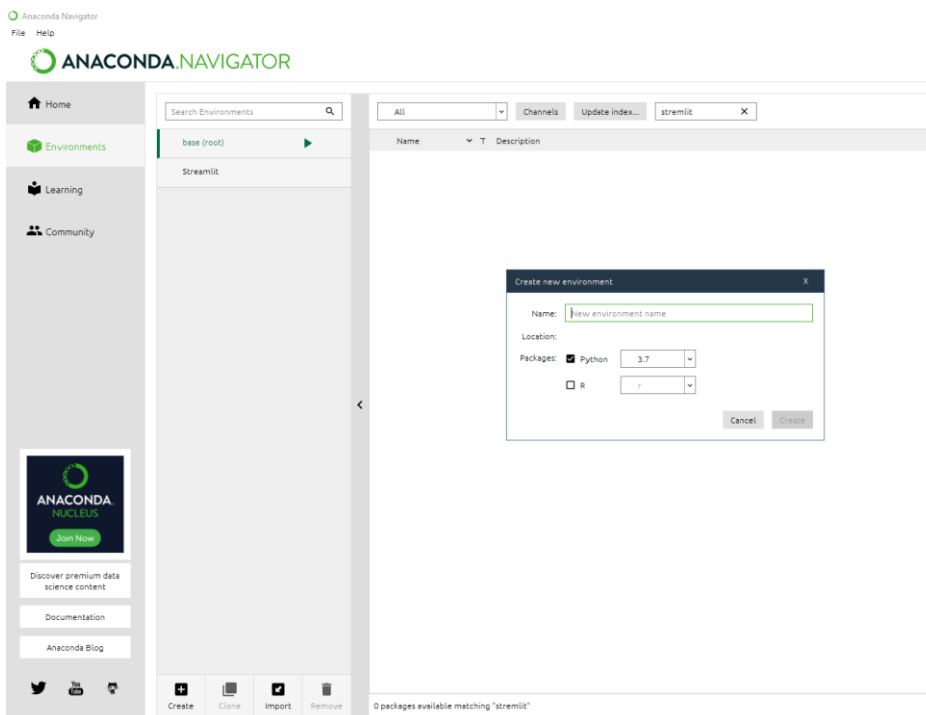
ภาคผนวก ก
ขั้นตอนการพัฒนาระบบ

ขั้นตอนการพัฒนาาระบบ

1. ดาวน์โหลด Anaconda <https://www.anaconda.com/products/individual> และติดตั้ง ให้เลือก bit ตาม Windows ที่เราลง



2. กด Create Environment เลือก package type “Python”



3. เปิด environment ที่สร้างไว้ แล้วทำการลง library ต่าง ๆ ที่เราใช้งาน

ตัวอย่าง การ Install Streamlit

Streamlit เป็น Python library ที่ถูกนำมาใช้พัฒนา Web Application สามารถทำงานกับ Library อื่นๆที่เป็นภาษา Python ด้วยกันได้ เช่น Pandas, Numpy, Matplotlib หรือ Scikit-learn เป็นต้น

<https://docs.streamlit.io/en/stable/troubleshooting/clean-install.html#install-streamlit-on-windows>

หรือสำหรับ OS อื่นๆ <https://docs.continuum.io/anaconda/install/>

Install Streamlit on Windows

Streamlit's officially-supported environment manager on Windows is [Anaconda Navigator](#).

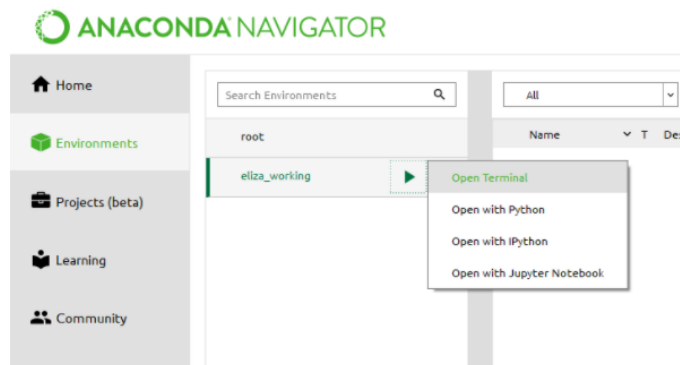
Install Anaconda

If you don't have Anaconda install yet, follow the steps provided on the [Anaconda installation page](#).

Create a new environment with Streamlit

Next you'll need to set up your environment.

1. Follow the steps provided by Anaconda to [set up and manage your environment](#) using the Anaconda Navigator.
2. Select the "▶" icon next to your new environment. Then select "Open terminal":



3. In the terminal that appears, type:

```
pip install streamlit
```

สำหรับ Library อื่น ๆ สามารถดูขั้นตอนการ install ได้จาก documentation website ของ Library นั้นๆ เช่น Scikit learn, pmdarima, pandas, numpy, altair, matplotlib, statsmodels เป็นต้น

4. พัฒนาระบบด้วยการเขียน code python ตัวอย่าง code

4.1 import library

```

1 import streamlit as st
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import altair as alt
6 import matplotlib
7 import matplotlib.pyplot as plt
8 #from sklearn import datasets
9 from sklearn.model_selection import train_test_split
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.naive_bayes import GaussianNB
12 from sklearn.svm import SVC
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.metrics import accuracy_score
16 from sklearn import model_selection
17 #from sklearn.preprocessing import LabelEncoder
18 # evaluate an ARIMA model using a walk-forward validation
19 from pandas import read_csv
20 from datetime import datetime
21 from matplotlib import pyplot
22 from statsmodels.tsa.arima.model import ARIMA
23 from statsmodels.tsa.arima_model import ARIMAResults
24 from math import sqrt
25 from pandas import DataFrame
26 # Compute Seasonal Index
27 from statsmodels.tsa.statespace.sarimax import SARIMAX
28 from statsmodels.tsa.stattools import adfuller
29 import matplotlib.pyplot as plt
30 from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
31 from sklearn.metrics import median_absolute_error, mean_squared_log_error
32 from sklearn.metrics import mean_absolute_percentage_error
33 import csv
34 #%matplotlib inline
35
36 import pandas_profiling
37
38 from datetime import datetime
39 import statsmodels.api as sm
40 from sklearn import metrics
41 #from sklearn.cross_validation import train_test_split
42 from sklearn.linear_model import LinearRegression
43 from sklearn.ensemble import RandomForestRegressor
44 from sklearn.tree import DecisionTreeRegressor
45 from sklearn.ensemble import AdaBoostRegressor
46 from sklearn.ensemble import GradientBoostingRegressor
47 from sklearn.svm import SVC, LinearSVC

```


4.2 เขียน code ด้วยภาษา Python ตัวอย่าง code ด้านล่างคือการสร้าง model ของเทคนิค ARIMA

```
def parser(x):=
def arima(data):
    if data is not None:
        #st.success("Data successfully Loaded")
        data_check = pd.read_csv(data.name)
        data_check.to_csv('train_timeseries.csv', index = False)
        #st.write(data_check.columns[0],data_check.columns[1])

        if 'Date' == data_check.columns[0] and 'Total' == data_check.columns[1]:

            d = pd.read_csv(data.name)

            # Create Training and Test
            size = int(len(d.Total) * 0.80)
            #train, test = X[0:size], X[size:len(d.Total)]
            train = d.Total[0:size]
            test = d.Total[size:len(d.Total)]
            test_date = d.Date[size:len(d.Date)].reset_index()

            testt = test.reset_index()
            kpss_diffs = ndiffs(train, alpha=0.05, test='kpss', max_d=6)
            adf_diffs = ndiffs(train, alpha=0.05, test='adf', max_d=6)
            n_diffs = max(adf_diffs, kpss_diffs)

            model = pm.auto_arima(train, start_p=0, start_q=0,
                test='adf', # use adftest to find optimal 'd'
                max_p=3, max_q=3, # maximum p and q
                d=n_diffs, # Let model determine 'd'
                seasonal=False, # No Seasonality
                start_P=0,
                trace=True,
                error_action='ignore',
                suppress_warnings=True,
                stepwise=True)

            tup = model.order

            # ***** save the model to disk
            filename = "ARIMA_Model.pkl"
            saveModel(model,filename)

            def forecast_one_step():
                fc, conf_int = model.predict(n_periods=1, return_conf_int=True)
                return (
                    fc.tolist()[0],
                    np.asarray(conf_int).tolist()[0])
```

5. Run code script ที่เราพัฒนาด้วยคำสั่ง streamlit run

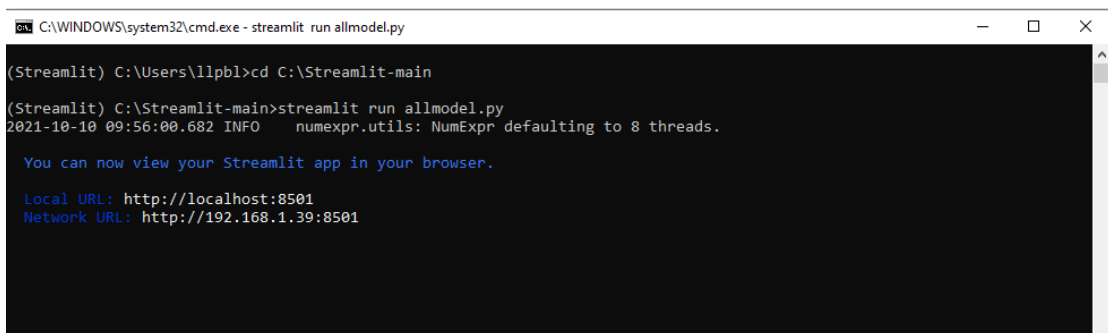
https://docs.streamlit.io/en/stable/main_concepts.html

Create an app

Working with Streamlit is simple. First you sprinkle a few Streamlit commands into a normal Python script, then you run it with `streamlit run`:

```
streamlit run your_script.py [-- script args]
```

As soon as you run the script as shown above, a local Streamlit server will spin up and your app will open in a new tab your default web browser. The app is your canvas, where you'll draw charts, text, widgets, tables, and more.



```
C:\WINDOWS\system32\cmd.exe - streamlit run allmodel.py
(Streamlit) C:\Users\llpbl>cd C:\Streamlit-main
(Streamlit) C:\Streamlit-main>streamlit run allmodel.py
2021-10-10 09:56:00.682 INFO numexpr.utils: NumExpr defaulting to 8 threads.

You can now view your Streamlit app in your browser.

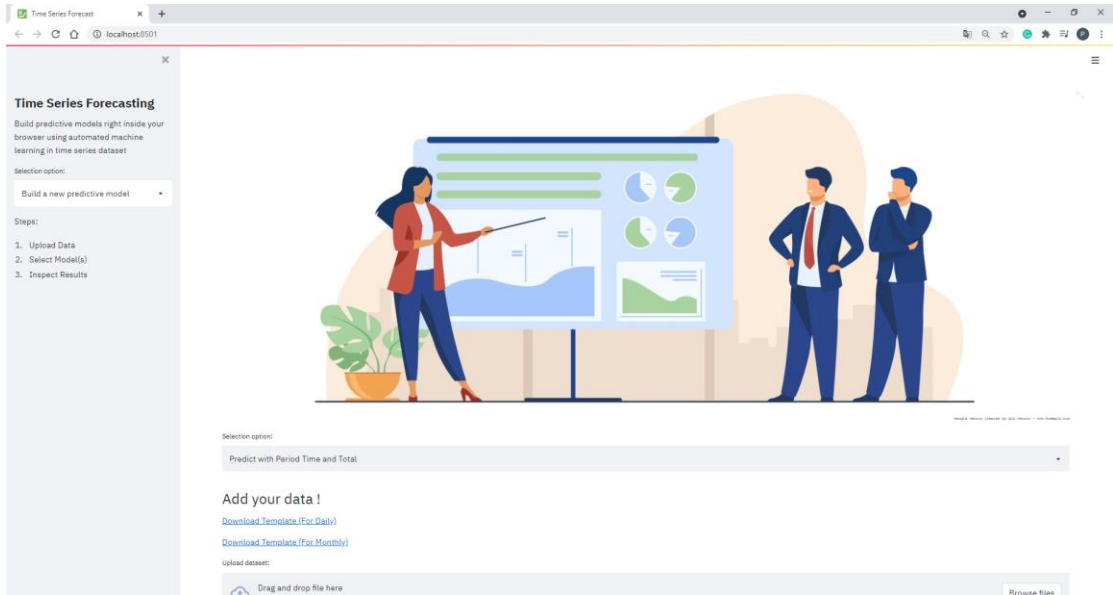
Local URL: http://localhost:8501
Network URL: http://192.168.1.39:8501
```

ภาคผนวก ข

**ขั้นตอนการใช้งานของระบบวิเคราะห์ข้อมูลอนุกรมเวลา
ด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง**

ขั้นตอนการใช้งานระบบวิเคราะห์ข้อมูลอนุกรมเวลาด้วยเทคนิคทางสถิติและการเรียนรู้ของเครื่อง

1. เข้าสู่หน้าจอของระบบจะมีตัวเลือกด้านซ้าย 2 ส่วน



Time Series Forecasting

Build predictive models right inside your browser using automated machine learning in time series dataset

Selection option:

Build a new predictive model |

- Build a new predictive model
- Apply models to new datasets

2. Select Model(s)
3. Inspect Results

1.1 สร้าง model

1.2 ใช้ model ที่สนใจพยากรณ์ข้อมูลจริง

1.1 สร้าง model

เลือกรูปแบบของข้อมูลนำเข้าที่ต้องการพยากรณ์

- แบบที่ 1 คือข้อมูล Date, Total
- แบบที่ 2 คือข้อมูล Date, Total, ID

Selection option:

Predict with Period Time and Total |

- Predict with Period Time and Total
- Predict with Period Time, ID and Total

แบบที่ 1 : รูปแบบนี้จะใช้เทคนิค ARIMA, SARIMA สร้าง model

The screenshot displays a software interface for time series forecasting. On the left, a sidebar titled 'Time Series Forecasting' provides instructions and steps: 'Build a new predictive model', 'Upload Data', 'Select Model(s)', and 'Inspect Results'. The main workspace features an illustration of a woman presenting to two men. Below this, there are sections for 'Add your data!', 'Upload Data', and 'Data successfully loaded' with a list of data points. To the right, two line plots are shown: 'ARIMA Forecast Line Plot' and 'SARIMA Forecast Line Plot', both plotting 'Total VS Total Forecast' against 'Date'. Below the plots are data tables with columns for 'Date', 'Actual', and 'Forecast'. A 'Comparison table' at the bottom right compares the performance of the models.

- สามารถ Download Template สำหรับข้อมูลรายวันหรือรายเดือนได้

Add your data !

[Download Template \(For Daily\)](#)

[Download Template \(For Monthly\)](#)

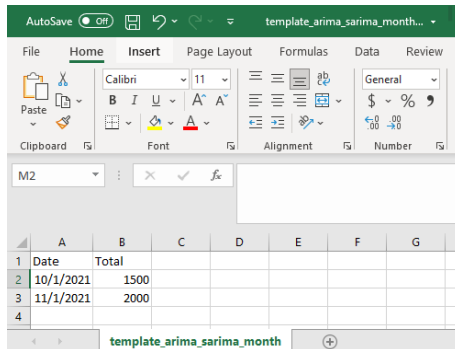
template_arima_sa....csv

ตัวอย่างไฟล์ template ข้อมูลนำเข้าสร้าง model รายวัน (.csv)

The screenshot shows an Excel spreadsheet with the following data:

| Date | Total |
|-----------|-------|
| 10/1/2021 | 1500 |
| 10/2/2021 | 2000 |

ตัวอย่างไฟล์ template ข้อมูลนำเข้าสร้าง model รายเดือน (.csv)



- Upload ข้อมูล .csv เพื่อสร้าง model
- เลือกความถี่ของข้อมูล (รายวัน, รายเดือน)
- เลือก model ที่ต้องการ หรือเลือกทั้งหมด

Selection option:
Predict with Period Time and Total

Add your data !
[Download Template \(For Daily\)](#)
[Download Template \(For Monthly\)](#)

Upload dataset:
 Drag and drop file here
 Limit 200MB per file • CSV Browse files

PPT_train.csv 0.0KB x

Data successfully loaded

| | Date | Total |
|----|-----------|--------|
| 0 | 8/1/2009 | 624472 |
| 1 | 9/1/2009 | 629663 |
| 2 | 10/1/2009 | 650772 |
| 3 | 11/1/2009 | 647287 |
| 4 | 12/1/2009 | 602661 |
| 5 | 1/1/2010 | 634652 |
| 6 | 2/1/2010 | 611039 |
| 7 | 3/1/2010 | 719225 |
| 8 | 4/1/2010 | 691343 |
| 9 | 5/1/2010 | 748456 |
| 10 | 6/1/2010 | 817301 |

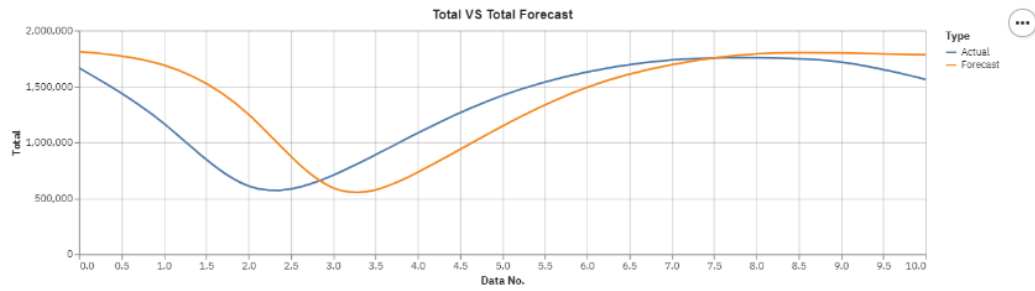
Selection frequency:
Month

Select one or more options:
ARIMA SARIMA o -

Select all models

- ระบบประมวลผลจากข้อมูลที่ upload เข้าไป 80% train, 20% validate test แสดงผลลัพธ์เปรียบเทียบค่าพยากรณ์และ ค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคาดเคลื่อน สามารถ export ค่าข้อมูลเปรียบเทียบในรูปแบบ .xlsx และ Save model ที่สนใจได้

ARIMA Forecast Line Plot

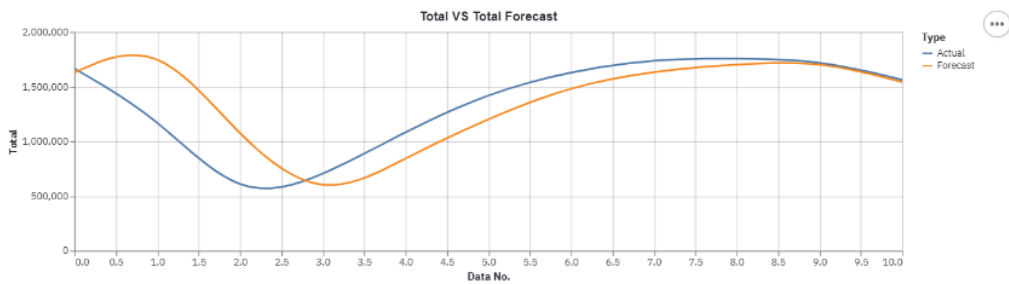


| | Date | Actual | Forecast |
|----|-----------|---------|----------------|
| 0 | 2/1/2563 | 1664132 | 1,810,297.2144 |
| 1 | 3/1/2563 | 1227632 | 1,746,234.9941 |
| 2 | 4/1/2563 | 439198 | 1,346,136.1556 |
| 3 | 5/1/2563 | 674483 | 373,399.7429 |
| 4 | 6/1/2563 | 1094334 | 715,648.3362 |
| 5 | 7/1/2563 | 1442882 | 1,156,167.4705 |
| 6 | 8/1/2563 | 1643622 | 1,520,785.6153 |
| 7 | 9/1/2563 | 1754362 | 1,709,073.1440 |
| 8 | 10/1/2563 | 1759392 | 1,808,604.7293 |
| 9 | 11/1/2563 | 1745557 | 1,799,514.9579 |
| 10 | 12/1/2563 | 1563252 | 1,783,208.5668 |

[Download Forecast Result excel file](#)

[Download Trained Model in ARIMA Model.pkl File](#)

SARIMA Forecast Line Plot



| | Date | Actual | Forecast |
|----|-----------|---------|----------------|
| 0 | 2/1/2563 | 1664132 | 1,629,884.4980 |
| 1 | 3/1/2563 | 1227632 | 1,958,101.9484 |
| 2 | 4/1/2563 | 439198 | 1,014,195.6672 |
| 3 | 5/1/2563 | 674483 | 438,552.8031 |
| 4 | 6/1/2563 | 1094334 | 851,034.3448 |
| 5 | 7/1/2563 | 1442882 | 1,211,913.1160 |
| 6 | 8/1/2563 | 1643622 | 1,506,791.4291 |
| 7 | 9/1/2563 | 1754362 | 1,645,677.0150 |
| 8 | 10/1/2563 | 1759392 | 1,705,748.0657 |
| 9 | 11/1/2563 | 1745557 | 1,740,035.5340 |
| 10 | 12/1/2563 | 1563252 | 1,542,221.9004 |

[Download Forecast Result excel file](#)

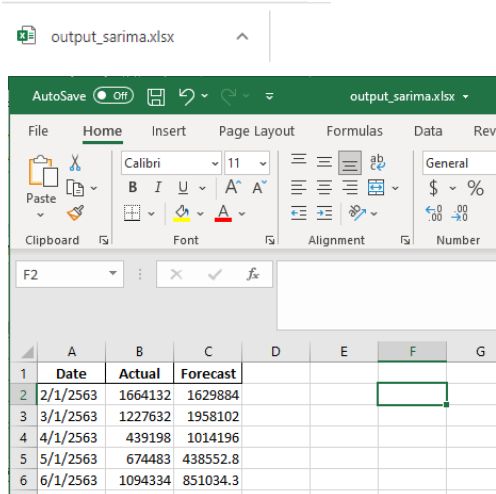
[Download Trained Model in SARIMA Model.pkl File](#)

Comparison table

Below, you can inspect the performance of your model

| | mean_absolute_percentage_error(MAPE) |
|--------|--------------------------------------|
| ARIMA | 35.1504 |
| SARIMA | 25.9028 |

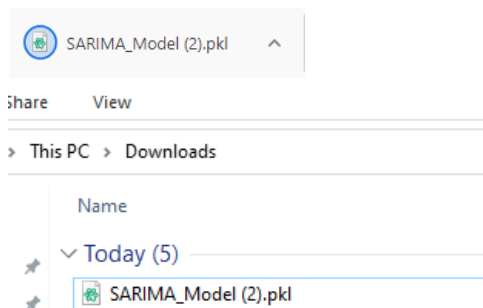
ตัวอย่างไฟล์ข้อมูลเปรียบเทียบในรูปแบบ .xlsx



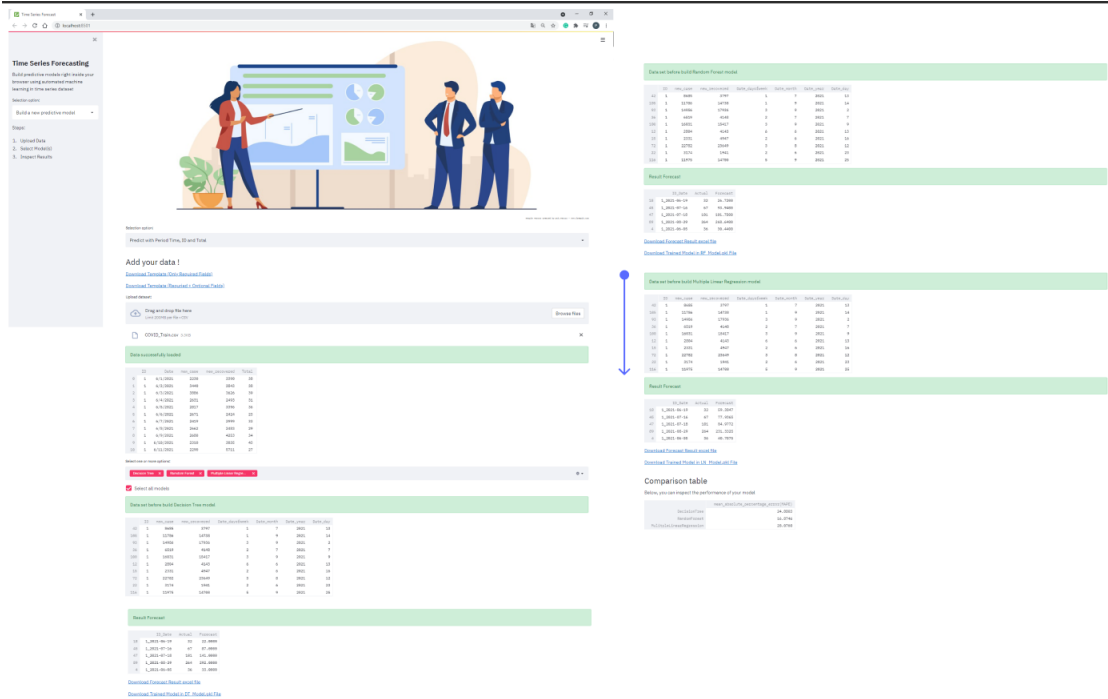
The screenshot shows an Excel spreadsheet with the following data:

| | A | B | C | D | E | F | G |
|---|-------------|---------------|-----------------|---|---|---|---|
| 1 | Date | Actual | Forecast | | | | |
| 2 | 2/1/2563 | 1664132 | 1629884 | | | | |
| 3 | 3/1/2563 | 1227632 | 1958102 | | | | |
| 4 | 4/1/2563 | 439198 | 1014196 | | | | |
| 5 | 5/1/2563 | 674483 | 438552.8 | | | | |
| 6 | 6/1/2563 | 1094334 | 851034.3 | | | | |

ตัวอย่างไฟล์ model ที่ทำการ download ในรูปแบบ .pkl



แบบที่ 2 : รูปแบบนี้จะใช้เทคนิค Decision Tree, Random Forest, Multiple Linear Regression
สร้าง model

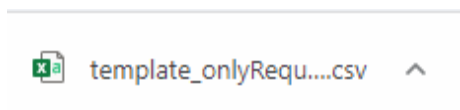


- สามารถ Download Template สำหรับข้อมูลตัวอย่างได้

Add your data !

[Download Template \(Only Required Fields\)](#)

[Download Template \(Required + Optional Fields\)](#)



ตัวอย่างไฟล์ template ข้อมูลนำเข้าสร้าง model เฉพาะข้อมูลที่จำเป็นต้องระบุ (.csv)

ประกอบไปด้วย ID, Date, Total

| | A | B | C | D | E | F | G |
|---|----|----------|-------|---|---|---|---|
| 1 | ID | Date | Total | | | | |
| 2 | 1 | 6/1/2021 | 38 | | | | |
| 3 | 2 | 6/2/2021 | 38 | | | | |

ตัวอย่างไฟล์ template ข้อมูลนำเข้าสร้าง model สามารถเพิ่มเติมข้อมูลประเภทตัวเลขเข้าไปเพิ่มเติมได้จากข้อมูลจำเป็นที่ต้องระบุ (ID, Date, Total) (.csv)

| | A | B | C | D | E | F | G | H |
|---|----|----------|-------|----------|---------------|---|---|---|
| 1 | ID | Date | Total | new_case | new_recovered | | | |
| 2 | 1 | 6/1/2021 | 38 | 2230 | 3390 | | | |
| 3 | 2 | 6/2/2021 | 38 | 3440 | 2843 | | | |

- Upload ข้อมูล .csv เพื่อสร้าง model
- เลือก model ที่ต้องการ หรือเลือกทั้งหมด

Selection option:
Predict with Period Time, ID and Total

Add your data !
[Download Template \(Only Required Fields\)](#)
[Download Template \(Required + Optional Fields\)](#)

Uploaded dataset:
 Drag and drop file here
 Limit 200MB per file - CSV Browse files

COVID_Train.csv 3.5 KB x

Data successfully loaded

| ID | Date | new_case | new_recovered | Total | |
|----|------|-----------|---------------|-------|----|
| 0 | 1 | 6/1/2021 | 2230 | 3390 | 38 |
| 1 | 1 | 6/2/2021 | 3440 | 2843 | 38 |
| 2 | 1 | 6/3/2021 | 3856 | 3106 | 39 |
| 3 | 1 | 6/4/2021 | 2631 | 2103 | 31 |
| 4 | 1 | 6/5/2021 | 2817 | 3396 | 36 |
| 5 | 1 | 6/6/2021 | 2671 | 2424 | 23 |
| 6 | 1 | 6/7/2021 | 2419 | 2099 | 33 |
| 7 | 1 | 6/8/2021 | 2662 | 2183 | 29 |
| 8 | 1 | 6/9/2021 | 2639 | 4233 | 34 |
| 9 | 1 | 6/10/2021 | 2310 | 3035 | 43 |
| 10 | 1 | 6/11/2021 | 2299 | 3711 | 27 |

Select one or more options:
Decision Tree Random Forest Multiple Linear Regr...

Select all models

- ระบบแสดงข้อมูลจริงหลังการแปลงข้อมูลผ่านระบบอัตโนมัติเพื่อสร้าง model

Data set before build Decision Tree model

| | ID | new_case | new_recovered | Date_dayofweek | Date_month | Date_year | Date_day | |
|--|-----|----------|---------------|----------------|------------|-----------|----------|----|
| | 42 | 1 | 8685 | 3797 | 1 | 7 | 2021 | 13 |
| | 105 | 1 | 11786 | 14738 | 1 | 9 | 2021 | 14 |
| | 93 | 1 | 14956 | 17936 | 3 | 9 | 2021 | 2 |
| | 36 | 1 | 6519 | 4148 | 2 | 7 | 2021 | 7 |
| | 100 | 1 | 16031 | 15417 | 3 | 9 | 2021 | 9 |
| | 12 | 1 | 2804 | 4143 | 6 | 6 | 2021 | 13 |
| | 15 | 1 | 2331 | 4947 | 2 | 6 | 2021 | 16 |
| | 72 | 1 | 22782 | 23649 | 3 | 8 | 2021 | 12 |
| | 22 | 1 | 3174 | 1941 | 2 | 6 | 2021 | 23 |
| | 116 | 1 | 11975 | 14700 | 5 | 9 | 2021 | 25 |

- ระบบประมวลผลจากข้อมูลที่ upload เข้าไป 80% train, 20% validate test แสดงผลลัพธ์เปรียบเทียบค่าพยากรณ์และ ค่าเฉลี่ยของค่าสัมบูรณ์เปอร์เซ็นต์ความคลาดเคลื่อน สามารถ export ค่าข้อมูลเปรียบเทียบในรูปแบบ .xlsx และ Save model ที่สนใจได้

Result Forecast

| | ID | Date | Actual | Forecast |
|--|----|--------------|--------|----------|
| | 18 | 1_2021-06-19 | 32 | 59.3847 |
| | 45 | 1_2021-07-16 | 67 | 77.9365 |
| | 47 | 1_2021-07-18 | 101 | 84.9772 |
| | 89 | 1_2021-08-29 | 264 | 231.3325 |
| | 4 | 1_2021-06-05 | 36 | 40.7875 |

[Download Forecast Result excel file](#)

[Download Trained Model in LN Model.pkl File](#)

Comparison table

Below, you can inspect the performance of your model

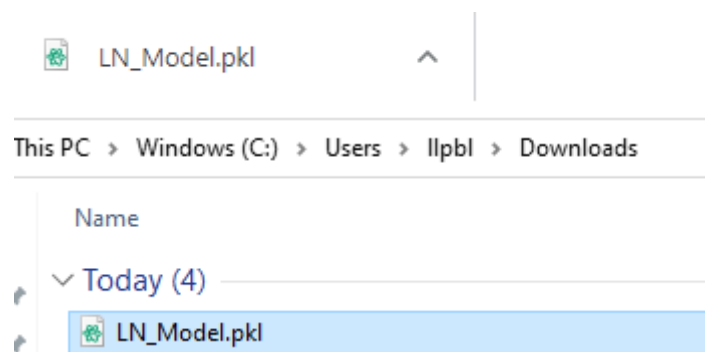
| | mean_absolute_percentage_error(MAPE) |
|--------------------------|--------------------------------------|
| DecisionTree | 24.8803 |
| RandomForest | 16.8746 |
| MuiltpleLinearRegression | 28.0760 |

ตัวอย่างไฟล์ข้อมูลเปรียบเทียบในรูปแบบ .xlsx

The screenshot shows an Excel spreadsheet with the following data:

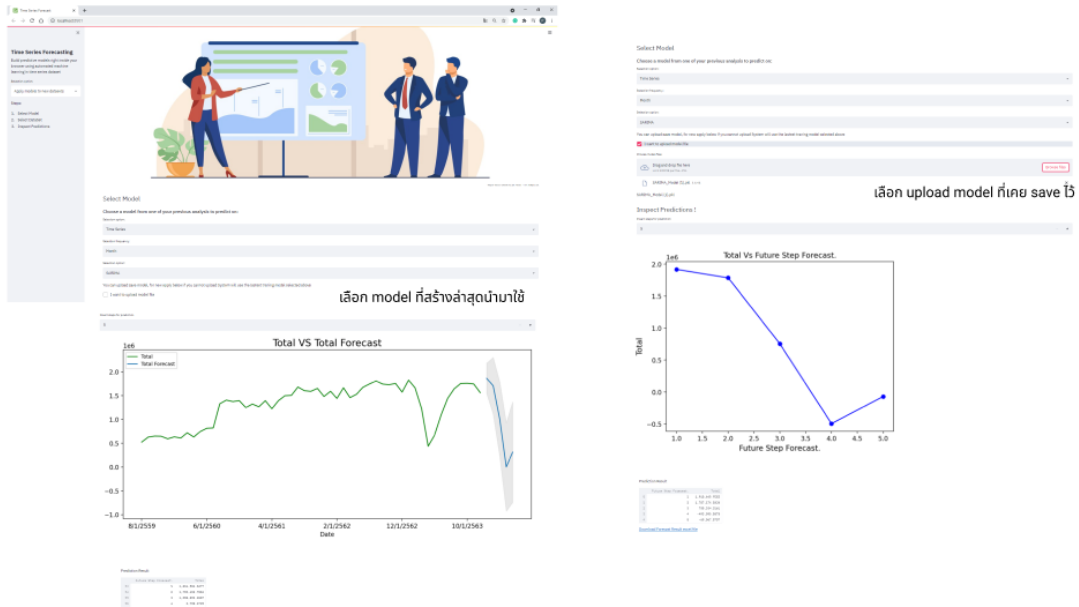
| ID | Date | Actual | Forecast |
|----|------------|--------|----------|
| 1 | 2021-06-19 | 32 | 59.3847 |
| 1 | 2021-07-16 | 67 | 77.93646 |
| 1 | 2021-07-18 | 101 | 84.9772 |
| 1 | 2021-08-29 | 264 | 231.3325 |
| 1 | 2021-06-05 | 36 | 40.7875 |
| 1 | 2021-07-11 | 86 | 65.69357 |
| 1 | 2021-08-02 | 178 | 154.7613 |

ตัวอย่างไฟล์ model ที่ทำการ download ในรูปแบบ .pkl



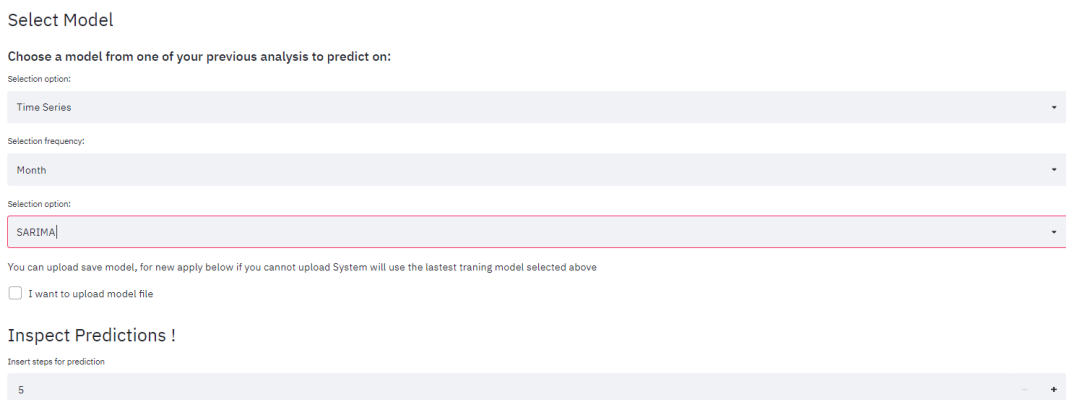
1.2 ใช้ model ที่สนใจพยากรณ์ข้อมูลจริง

แบบที่ 1 : รูปแบบนี้จะใช้เทคนิค ARIMA, SARIMA พยากรณ์



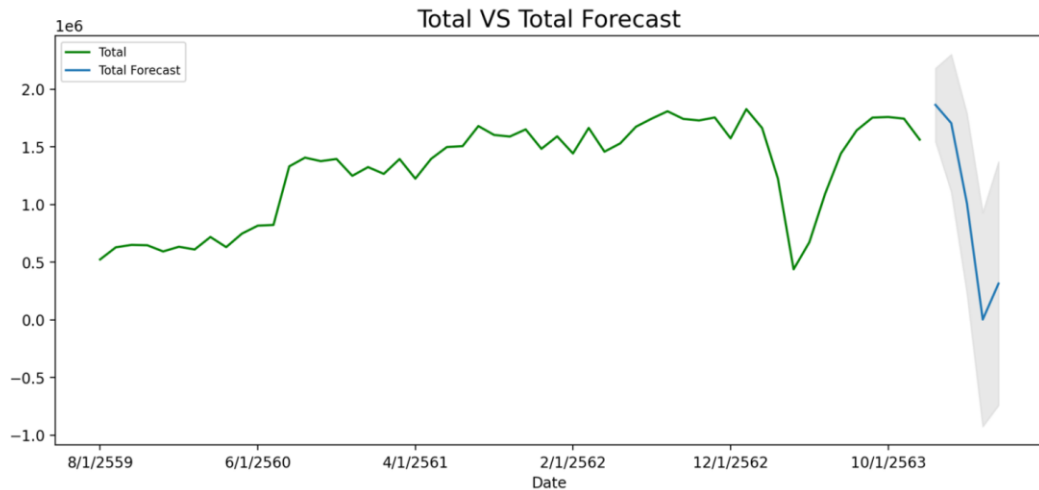
- เลือกประเภทการพยากรณ์ (Time Series)
- เลือกความถี่ของข้อมูล (รายวัน, รายเดือน)
- เลือก model ที่ต้องการพยากรณ์

กรณีที่ 1 ระบบจะทำการเก็บ model ล่าสุดที่เพิ่งสร้างไว้หากต้องการเลือกใช้ model ที่เพิ่งสร้างสามารถเลือก model ที่ต้องการพยากรณ์ได้ดังรูปด้านล่าง
ตัวอย่าง กรณีนี้เลือก SARIMA model ช่วงความถี่รายเดือน และต้องการพยากรณ์ล่วงหน้า 5 เดือน



- ระบบประมวลผลและแสดงผลลัพธ์ค่าพยากรณ์

กรณีนี้กราฟจะแสดงข้อมูล train ที่ระบบให้ทำการจัดเก็บไว้ล่าสุดสำหรับสร้าง model (สีเขียว) แสดงเปรียบเทียบกับค่าพยากรณ์ใหม่ด้วย (สีฟ้า)



Prediction Result

| | Future Step Forecast | Total |
|----|----------------------|--------------|
| 53 | 1 | 1,864,581.63 |
| 54 | 2 | 1,706,402.79 |
| 55 | 3 | 1,008,853.91 |
| 56 | 4 | 3,758.38 |
| 57 | 5 | 315,255.02 |

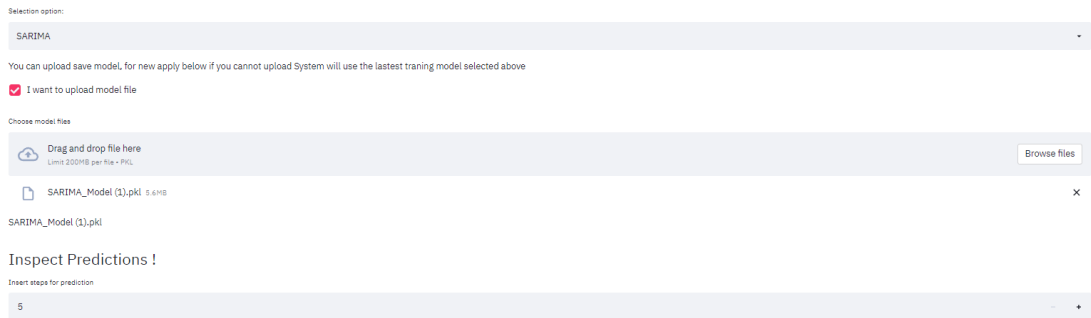
[Download Forecast Result excel file](#)

- Export ข้อมูลออกมาในรูปแบบ .xlsx

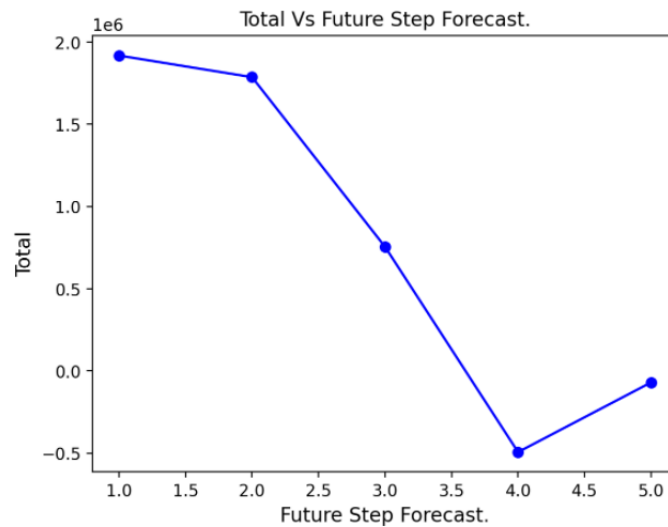
| | Future Step Forecast | Total |
|---|----------------------|--------------|
| 1 | | |
| 2 | 53 | 1,864,581.63 |
| 3 | 54 | 1,706,402.79 |
| 4 | 55 | 1,008,853.91 |
| 5 | 56 | 3,758.38 |
| 6 | 57 | 315,255.02 |
| 7 | | |
| 8 | | |

กรณีที่ 2 สามารถเลือก upload model ที่จัดเก็บไว้ (.pkl) สำหรับการพยากรณ์ได้ดังรูป
ด้านล่าง

ตัวอย่าง กรณีนี้เลือก upload SARIMA model ช่วงความถี่รายเดือน และต้องการพยากรณ์
ล่วงหน้า 5 เดือน



- ระบบประมวลผลและแสดงผลลัพท์ค่าพยากรณ์

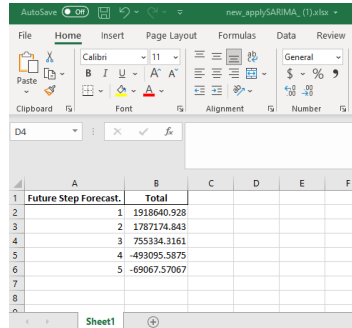


Prediction Result

| | Future Step Forecast. | Total |
|---|-----------------------|----------------|
| 0 | 1 | 1,918,649.9283 |
| 1 | 2 | 1,787,174.8426 |
| 2 | 3 | 755,334.3161 |
| 3 | 4 | -493,095.5875 |
| 4 | 5 | -69,067.5707 |

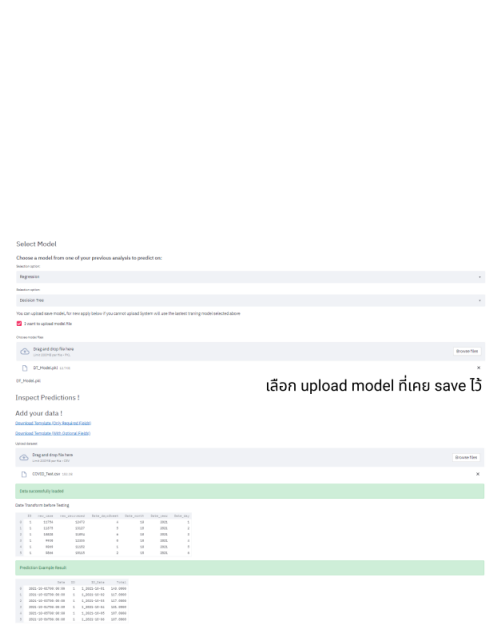
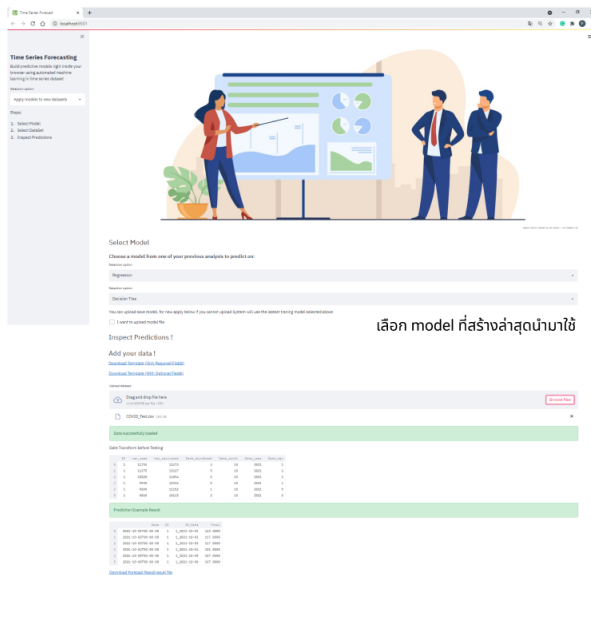
[Download Forecast Result excel file](#)

- Export ข้อมูลออกมาในรูปแบบ .xlsx



แบบที่ 2 : รูปแบบนี้จะใช้เทคนิค Decision Tree, Random Forest, Multiple Linear Regression

พยากรณ์



- เลือกประเภทการพยากรณ์ (Regression)
- เลือก model ที่ต้องการพยากรณ์

กรณีที่ 1 ระบบจะทำการเก็บ model ล่าสุดที่เพิ่งสร้างไว้หากต้องการเลือกใช้ model ที่เพิ่งสร้างสามารถเลือก model ที่ต้องการพยากรณ์ได้ดังรูปด้านล่าง

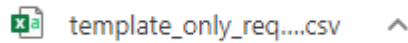
ตัวอย่าง กรณีนี้เลือก Decision tree model

- สามารถ Download Template สำหรับข้อมูลตัวอย่างได้

Add your data !

[Download Template \(Only Required Fields\)](#)

[Download Template \(With Optional Fields\)](#)



ตัวอย่างไฟล์ template ข้อมูลสำหรับพยากรณ์เฉพาะข้อมูลที่เป็นต้องระบุ (.csv)
ประกอบไปด้วย ID, Date

| ID | Date |
|----|----------|
| 1 | 6/1/2021 |
| 2 | 6/2/2021 |

ตัวอย่างไฟล์ template ข้อมูลสำหรับพยากรณ์สามารถเพิ่มเติมข้อมูลประเภทตัวเลขเข้าไป
เพิ่มเติมได้จากข้อมูลจำเป็นที่ต้องระบุ (ID, Date) (.csv)

| ID | Date | new_case | new_recovered |
|----|----------|----------|---------------|
| 1 | 6/1/2021 | 2230 | 3390 |
| 2 | 6/2/2021 | 3440 | 2843 |

- Upload ข้อมูล .csv ที่ต้องการให้พยากรณ์ในอนาคต

ตัวอย่าง : กรณีต้องการพยากรณ์จำนวนผู้เสียชีวิตรายใหม่จากสถานการณ์ COVID-19
ระลอก 3 วันที่ 1-6 ตุลาคม 2564

| ID | Date | new_case | new_recovered |
|----|-----------|----------|---------------|
| 1 | 10/1/2021 | 11754 | 12473 |
| 1 | 10/2/2021 | 11375 | 13127 |
| 1 | 10/3/2021 | 10828 | 11894 |
| 1 | 10/4/2021 | 9930 | 12336 |
| 1 | 10/5/2021 | 9869 | 11152 |
| 1 | 10/6/2021 | 9866 | 10115 |

- ระบบแสดงข้อมูลจริงหลังการแปลงข้อมูลผ่านระบบอัตโนมัติเพื่อนำเข้าพยากรณ์

Date Transform before Testing

| | ID | new_case | new_recovered | Date_dayofweek | Date_month | Date_year | Date_day |
|---|----|----------|---------------|----------------|------------|-----------|----------|
| 0 | 1 | 11754 | 12473 | 4 | 10 | 2021 | 1 |
| 1 | 1 | 11375 | 13127 | 5 | 10 | 2021 | 2 |
| 2 | 1 | 10828 | 11894 | 6 | 10 | 2021 | 3 |
| 3 | 1 | 9930 | 12336 | 0 | 10 | 2021 | 4 |
| 4 | 1 | 9869 | 11152 | 1 | 10 | 2021 | 5 |
| 5 | 1 | 9866 | 10115 | 2 | 10 | 2021 | 6 |

- ระบบประมวลผลและแสดงผลลัพธ์ค่าพยากรณ์

Prediction Example Result

| | Date | ID | ID_Date | Total |
|---|---------------------|----|--------------|----------|
| 0 | 2021-10-01T00:00:00 | 1 | 1_2021-10-01 | 143.0000 |
| 1 | 2021-10-02T00:00:00 | 1 | 1_2021-10-02 | 117.0000 |
| 2 | 2021-10-03T00:00:00 | 1 | 1_2021-10-03 | 117.0000 |
| 3 | 2021-10-04T00:00:00 | 1 | 1_2021-10-04 | 101.0000 |
| 4 | 2021-10-05T00:00:00 | 1 | 1_2021-10-05 | 107.0000 |
| 5 | 2021-10-06T00:00:00 | 1 | 1_2021-10-06 | 107.0000 |

[Download Forecast Result excel file](#)

- Export ข้อมูลออกมาในรูปแบบ .xlsx

| | A | B | C | D | E |
|---|---------------------|----|--------------|-------|---|
| 1 | Date | ID | ID_Date | Total | |
| 2 | 2021-10-01 00:00:00 | 1 | 1_2021-10-01 | 143 | |
| 3 | 2021-10-02 00:00:00 | 1 | 1_2021-10-02 | 117 | |
| 4 | 2021-10-03 00:00:00 | 1 | 1_2021-10-03 | 117 | |
| 5 | 2021-10-04 00:00:00 | 1 | 1_2021-10-04 | 101 | |
| 6 | 2021-10-05 00:00:00 | 1 | 1_2021-10-05 | 107 | |
| 7 | 2021-10-06 00:00:00 | 1 | 1_2021-10-06 | 107 | |

กรณีที่ 2 สามารถเลือก upload model ที่จัดเก็บไว้ (.pkl) สำหรับการพยากรณ์ได้ดังรูป
ด้านล่างตัวอย่าง

กรณีนี้เลือก Upload Decision tree model และ upload ข้อมูลสำหรับพยากรณ์จำนวน
ผู้เสียชีวิตรายใหม่จากสถานการณ์ COVID-19 ระลอก 3 วันที่ 1-6 ตุลาคม 2564

Selection option:
Decision Tree

You can upload save model, for new apply below if you cannot upload System will use the latest training model selected above

I want to upload model file

Choose model file

Drag and drop file here
Limit 200MB per file - PKL

DT_Model.pkl 12.7KB

DT_Model.pkl

Inspect Predictions !

Add your data !

[Download Template \(Only Required Fields\)](#)
[Download Template \(With Optional Fields\)](#)

Upload dataset:

Drag and drop file here
Limit 200MB per file - CSV

COVID_Test.csv 182.0B

Data successfully loaded

- ระบบแสดงข้อมูลจริงหลังการแปลงข้อมูลผ่านระบบอัตโนมัติเพื่อนำเข้าพยากรณ์

Date Transform before Testing

| | ID | new_case | new_recovered | Date_dayofweek | Date_month | Date_year | Date_day |
|---|----|----------|---------------|----------------|------------|-----------|----------|
| 0 | 1 | 11754 | 12473 | 4 | 10 | 2021 | 1 |
| 1 | 1 | 11375 | 13127 | 5 | 10 | 2021 | 2 |
| 2 | 1 | 10828 | 11894 | 6 | 10 | 2021 | 3 |
| 3 | 1 | 9930 | 12336 | 0 | 10 | 2021 | 4 |
| 4 | 1 | 9869 | 11152 | 1 | 10 | 2021 | 5 |
| 5 | 1 | 9866 | 10115 | 2 | 10 | 2021 | 6 |

- ระบบประมวลผลและแสดงผลลัพธ์ค่าพยากรณ์

Prediction Example Result

| | Date | ID | ID_Date | Total |
|---|---------------------|----|--------------|----------|
| 0 | 2021-10-01T00:00:00 | 1 | 1_2021-10-01 | 143.0000 |
| 1 | 2021-10-02T00:00:00 | 1 | 1_2021-10-02 | 117.0000 |
| 2 | 2021-10-03T00:00:00 | 1 | 1_2021-10-03 | 117.0000 |
| 3 | 2021-10-04T00:00:00 | 1 | 1_2021-10-04 | 101.0000 |
| 4 | 2021-10-05T00:00:00 | 1 | 1_2021-10-05 | 107.0000 |
| 5 | 2021-10-06T00:00:00 | 1 | 1_2021-10-06 | 107.0000 |

[Download Forecast Result excel file](#)

Export ข้อมูลออกมาในรูปแบบ .xlsx

The screenshot shows the Microsoft Excel interface with the following data in the worksheet:

| | A | B | C | D | E |
|---|---------------------|----|--------------|-------|---|
| 1 | Date | ID | ID_Date | Total | |
| 2 | 2021-10-01 00:00:00 | 1 | 1_2021-10-01 | 143 | |
| 3 | 2021-10-02 00:00:00 | 1 | 1_2021-10-02 | 117 | |
| 4 | 2021-10-03 00:00:00 | 1 | 1_2021-10-03 | 117 | |
| 5 | 2021-10-04 00:00:00 | 1 | 1_2021-10-04 | 101 | |
| 6 | 2021-10-05 00:00:00 | 1 | 1_2021-10-05 | 107 | |
| 7 | 2021-10-06 00:00:00 | 1 | 1_2021-10-06 | 107 | |

ประวัติผู้เขียน

ชื่อ-นามสกุล

นางสาว พรทิวา วิศิษฏ์สรอรรถ

ประวัติการศึกษา

วิศวกรรมศาสตรบัณฑิต

ภาควิชาวิศวกรรมคอมพิวเตอร์

มหาวิทยาลัยเทคโนโลยีพระจอมเกล้าธนบุรี

ปีการศึกษา 2556

ตำแหน่งและสถานที่ทำงานปัจจุบัน

ผู้จัดการ

ฝ่ายพัฒนาผลิตภัณฑ์และการใช้ดิจิทัลเพื่อลูกค้าธุรกิจ

SME สายงานพัฒนาธุรกิจดิจิทัล โซลูชันเพื่อลูกค้า

ธุรกิจ กลุ่มงานนวัตกรรมดิจิทัลและข้อมูล

ธนาคารกรุงศรีอยุธยา จำกัด (มหาชน)