



การศึกษาการทำนายค่าความเสถียรของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด19  
ด้วยวิธีการเรียนรู้ของเครื่องจักร

PREDICTING STABLE COVID-19 mRNA VACCINE BY MACHINE LEARNING

ณัฐนรี พอสูงเนิน

บัณฑิตวิทยาลัย มหาวิทยาลัยศรีนครินทรวิโรฒ

2566

การศึกษาการทำนายค่าความเสียหายของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด19  
ด้วยวิธีการเรียนรู้ของเครื่องจักร



สารนิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตร  
วิทยาศาสตร์มหาบัณฑิต สาขาวิชาวิทยาการข้อมูล  
คณะวิทยาศาสตร์ มหาวิทยาลัยศรีนครินทรวิโรฒ  
ปีการศึกษา 2566  
ลิขสิทธิ์ของมหาวิทยาลัยศรีนครินทรวิโรฒ

PREDICTING STABLE COVID-19 mRNA VACCINE BY MACHINE LEARNING



NUTNAREE POOSUNGNOEN

A Master's Project Submitted in Partial Fulfillment of the Requirements

for the Degree of MASTER OF SCIENCE

(Data Science)

Faculty of Science, Srinakharinwirot University

2023

Copyright of Srinakharinwirot University

สารนิพนธ์

เรื่อง

การศึกษาการทำนายค่าความเสถียรของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด19

ด้วยวิธีการเรียนรู้ของเครื่องจักร

ของ

ณัฐนรี พอสุงเนิน

ได้รับอนุมัติจากบัณฑิตวิทยาลัยให้นับเป็นส่วนหนึ่งของการศึกษาตามหลักสูตร

ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาวิทยาการข้อมูล

ของมหาวิทยาลัยศรีนครินทรวิโรฒ

(รองศาสตราจารย์ นายแพทย์ฉัตรชัย เอกปัญญาสกุล)

คณบดีบัณฑิตวิทยาลัย

คณะกรรมการสอบปากเปล่าสารนิพนธ์

..... ที่ปรึกษาหลัก  
(อาจารย์ ดร.ศุภร คนธมักค์)

..... ประธาน  
(ผู้ช่วยศาสตราจารย์ ดร.รัตน์ชัยนันท์ ธรรมสุขจิต)

..... กรรมการ  
(ผู้ช่วยศาสตราจารย์ ดร.ศิริสรรพ เหล่าหะเกียรติ)

ชื่อเรื่อง	การศึกษาการทำนายค่าความเสียหายของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด19
	ด้วยวิธีการเรียนรู้ของเครื่องจักร
ผู้วิจัย	ณัฐนรี พอสุงเนิน
ปริญญา	วิทยาศาสตรมหาบัณฑิต
ปีการศึกษา	2566
อาจารย์ที่ปรึกษา	อาจารย์ ดร. ศุภร คนธภักดิ์

จากการระบาดของโรคโควิด19 มีการนำวัคซีนเอ็มอาร์เอ็นเอที่มีประสิทธิภาพมาใช้ในการป้องกันโรค งานวิจัยนี้นำเสนอวิธีการเรียนรู้ของเครื่องจักรโดยใช้เทคนิค 4 แบบ คือ XGboost, Random Forest, Catboost และ LightGBM เพื่อทำนายความเสียหายของการย่อยสลายของลำดับอาร์เอ็นเอ โดยเทคนิค LightGBM เป็นตัวทำนายที่ดีที่สุดโดยได้ MCRMSE เท่ากับ 0.31265 โดยผลลัพธ์ที่ได้ชี้ให้เห็นความเป็นไปได้ในการใช้วิธีการดังกล่าวในการพัฒนาวัคซีนเอ็มอาร์เอ็นเอในอนาคต

คำสำคัญ : โรคโควิด-19, วัคซีนเอ็มอาร์เอ็นเอ, การเรียนรู้ของเครื่องจักร, การย่อยสลาย

Title PREDICTING STABLE COVID-19 mRNA VACCINE BY  
MACHINE LEARNING

Author NUTNAREE POOSUNGNOEN

Degree MASTER OF SCIENCE

Academic Year 2023

Thesis Advisor Dr. Subhorn Khonthapagdee

As a result of the COVID-19 outbreak, an effective mRNA vaccine has been used to prevent disease. This research proposes a machine learning method using four techniques: XGboost, Random Forest, Catboost, and LightGBM to predict the stability of RNA sequence degradation. It was found that LightGBM was the best predictor. The results found that the MCRMSE is 0.31265, with the results suggesting the possibility of using this method in the future development of the MRNA vaccine.

Keyword : COVID-19, mRNA vaccine, Machine learning, Degradation

## กิตติกรรมประกาศ

สารนิพนธ์นี้สำเร็จลุล่วงได้ด้วยความอนุเคราะห์จาก ดร.ศุภร คนธภักดี และ ผศ.ดร.จันทร์ ผลประเสริฐ อาจารย์ที่ปรึกษาทั้ง 2 ท่านที่ให้คำปรึกษา คำแนะนำในการทำสารนิพนธ์ ตลอดจนสนับสนุนข้อมูลทางวิชาการ

ขอกราบขอบพระคุณคณะกรรมการสอบสารนิพนธ์ ที่ได้ให้คำแนะนำและข้อเสนอแนะ สำหรับการปรับปรุงสารนิพนธ์ และการใช้เทคนิคการเรียนรู้ของเครื่องเป็นเครื่องมือในการวิเคราะห์ และแก้ไขปัญหาทางข้อมูลต่อไป

ขอกราบขอบพระคุณ บัณฑิตวิทยาลัย และ คณะวิทยาศาสตร์ มหาวิทยาลัยศรีนครินทรวิโรฒ สำหรับการช่วยเหลือด้านเอกสารและการจัดทำเล่มสารนิพนธ์

ขอกราบขอบพระคุณ บิดา มารดา และ ขอขอบคุณ พี่เอกปรีญา ไบสนิ สำหรับกำลังใจในการทำสารนิพนธ์นี้

ณัฐนรี พอสุงเนิน

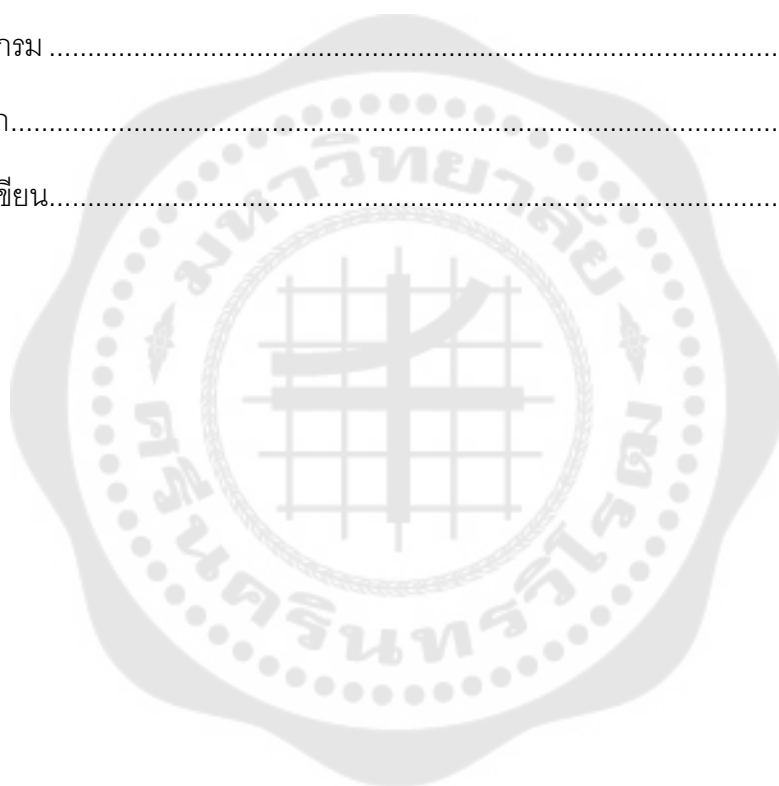
## สารบัญ

	หน้า
บทคัดย่อภาษาไทย .....	ง
บทคัดย่อภาษาอังกฤษ .....	จ
กิตติกรรมประกาศ.....	ฉ
สารบัญ .....	ช
สารบัญตาราง.....	ญ
สารบัญรูปภาพ .....	ฎ
บทที่ 1 บทนำ.....	1
ความสำคัญและความเป็นมาของงานวิจัย.....	1
วัตถุประสงค์.....	3
ขอบเขตงานวิจัย .....	3
วิธีการดำเนินการวิจัย.....	3
ประโยชน์ที่คาดว่าจะได้รับจากงานวิจัย .....	3
บทที่ 2 วรรณกรรมและงานวิจัยที่เกี่ยวข้อง.....	4
โรคโควิด 19 (Covid-19) .....	4
ไวรัสโคโรนา 2019 .....	4
อาร์เอ็นเอ (RNA) .....	5
เอ็มอาร์เอ็นเอ (mRNA).....	6
วัคซีน (vaccine).....	6
การพัฒนาวัคซีน COVID-19 .....	6
วัคซีนเอ็มอาร์เอ็นเอ .....	9
การทำงานของวัคซีนเอ็มอาร์เอ็นเอ.....	10



การออกแบบวัคซีนเอ็มอาร์เอ็นเอที่มีความเสถียร .....	10
ปัญหาเกี่ยวกับวัคซีนเอ็มอาร์เอ็นเอ .....	11
การเรียนรู้ของเครื่องจักร (Machine Learning) .....	13
ความแตกต่างระหว่าง Machine Learning และการเขียนโปรแกรมดั้งเดิม .....	13
ประเภทของการเรียนรู้โดยเครื่องจักร .....	14
ทฤษฎีของแบบจำลองที่นำมาใช้ในงานวิจัย .....	16
XGBoost (eXtreme Gradient Boosting) .....	16
Random Forest .....	16
Light GBM (Light Gradient Boosting) .....	16
CatBoost (Gradient boosting with categorical features support) .....	16
Cross validation .....	16
งานวิจัยที่เกี่ยวข้อง .....	17
สรุปงานวิจัยที่ศึกษาทั้งหมด .....	20
บทที่ 3 การดำเนินการวิจัย .....	21
แผนขั้นตอนดำเนินการวิจัย .....	21
ชุดข้อมูล (data set) เพื่อใช้ในงานวิจัย .....	22
กระบวนการสำรวจข้อมูลเบื้องต้น (Exploratory Data Analysis) .....	24
ศึกษาวิธีการสร้างและสร้างโมเดลในการทำนายข้อมูล .....	29
ทดสอบการทำแบบจำลองด้วยเทคนิค XGboost (baseline) .....	29
ทดสอบการทำแบบจำลองด้วยเทคนิค Random Forest .....	29
ทดสอบการทำแบบจำลองด้วยเทคนิค Catboost .....	30
ทดสอบการทำแบบจำลองด้วยเทคนิค Light GBM .....	30
ปรับปรุงโมเดลและวัดประสิทธิภาพของแบบจำลอง .....	31

MCRMSE (mean columnwise root mean squared error).....	32
บทที่ 4 ผลการดำเนินงานวิจัย.....	33
การหา Feature importance.....	34
บทที่ 5 สรุปผลการวิจัย อภิปรายผลและข้อเสนอแนะ .....	40
5.1 สรุปผลการวิจัย และ อภิปรายผลการวิจัย.....	40
5.2 ข้อเสนอแนะ.....	41
บรรณานุกรม .....	43
ภาคผนวก.....	46
ประวัติผู้เขียน.....	61



## สารบัญตาราง

	หน้า
ตาราง 1 รายละเอียดคอลลัมน์ของชุดฝึกฝน .....	23
ตาราง 2 พารามิเตอร์ที่เหมาะสมของแต่ละโมเดลที่ใช้ในการเรียนรู้.....	31
ตาราง 3 Feature Importance ที่สำคัญที่สุดของแต่ละเป้าหมาย .....	41

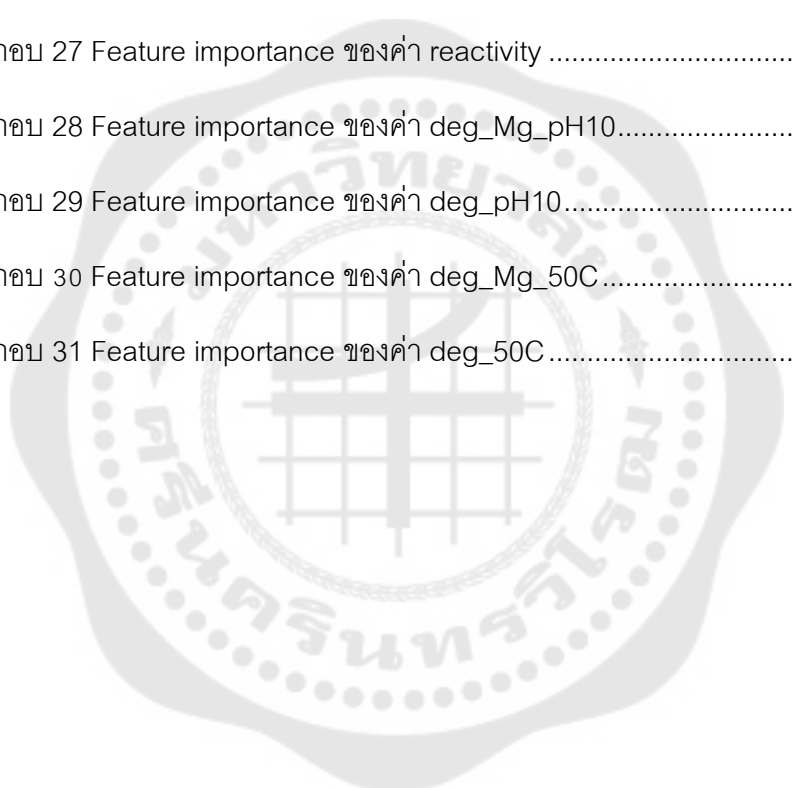


## สารบัญรูปภาพ

หน้า

ภาพประกอบ 1 แผนผังการแพร่ระบาดครั้งใหญ่ของเชื้อไวรัสไข้หวัดใหญ่ และเชื้อ SARS-CoV-2	2
ภาพประกอบ 2 โครงสร้างไวรัสโคโรนา.....	4
ภาพประกอบ 3 อาร์เอ็นเอหรือกรดไรโบนิวคลีอิก.....	5
ภาพประกอบ 4 แสดงโครงสร้างโมเลกุล mRNA ที่สมบูรณ์.....	6
ภาพประกอบ 5 ชนิดของวัคซีน COVID-19 ที่มีการพัฒนา.....	7
ภาพประกอบ 6 การทำงานของวัคซีนเอ็มอาร์เอ็นเอต่อไวรัส.....	10
ภาพประกอบ 7 mRNA of a mini protein, with bases colored by the probability they are unpaired. ....	11
ภาพประกอบ 8 ประเภทของการเรียนรู้โดยเครื่องจักร.....	13
ภาพประกอบ 9 ความแตกต่างระหว่าง Machine Learning และการเขียนโปรแกรมดั้งเดิม.....	14
ภาพประกอบ 10 แสดงตารางสรุปงานวิจัยที่ศึกษา.....	20
ภาพประกอบ 11 แผนผังการทำงานของระบบ.....	21
ภาพประกอบ 12 ดูชนิดข้อมูล.....	24
ภาพประกอบ 13 ดู Summary ของข้อมูล.....	24
ภาพประกอบ 14 ตัวอย่างของข้อมูลตามคอลัมน์.....	25
ภาพประกอบ 15 การใช้คำสั่ง count และ len.....	25
ภาพประกอบ 16 จำนวน SN_filter และ signal to noise.....	26
ภาพประกอบ 17 จำนวนความยาว sequence ข้อมูลทดสอบ.....	27
ภาพประกอบ 18 คำนวณค่าเฉลี่ยสำหรับคอลัมน์เป้าหมายแบบไม่ใช้แมกนีเซียม.....	28
ภาพประกอบ 19 คำนวณค่าเฉลี่ยสำหรับคอลัมน์เป้าหมายแบบใช้แมกนีเซียม.....	28
ภาพประกอบ 20 การกำหนดค่าพารามิเตอร์เทคนิค XGboost.....	29

ภาพประกอบ 21 พารามิเตอร์ที่ได้จาก GridSearchCV ในเทคนิค XGboost.....	29
ภาพประกอบ 22 การพารามิเตอร์ของเทคนิค Random Forest.....	30
ภาพประกอบ 23 การกำหนดพารามิเตอร์ของเทคนิค Catboost.....	30
ภาพประกอบ 24 การกำหนดพารามิเตอร์ของเทคนิค Light GBM .....	30
ภาพประกอบ 25 ค่า MCRMSE ของแต่ละแบบจำลอง .....	33
ภาพประกอบ 26 parameter ของ feature importance .....	34
ภาพประกอบ 27 Feature importance ของค่า reactivity .....	35
ภาพประกอบ 28 Feature importance ของค่า deg_Mg_pH10.....	36
ภาพประกอบ 29 Feature importance ของค่า deg_pH10.....	37
ภาพประกอบ 30 Feature importance ของค่า deg_Mg_50C.....	38
ภาพประกอบ 31 Feature importance ของค่า deg_50C.....	39



## บทที่ 1

### บทนำ

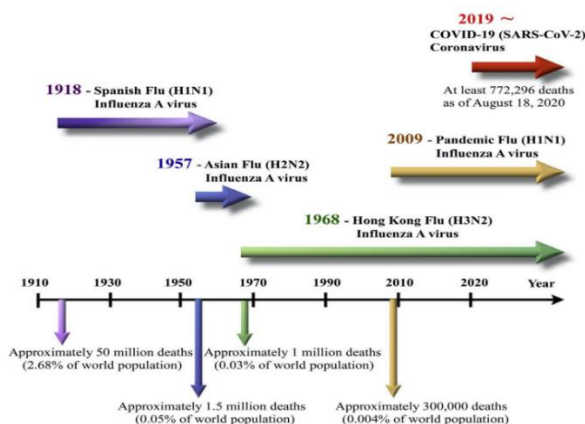
#### ความสำคัญและความเป็นมาของงานวิจัย

โรคโควิด-19 (covid-19) คือโรคติดต่อเกิดจากไวรัสโคโรนา 2019-nCoV พบครั้งแรกที่เมืองอู่ฮั่น ประเทศจีน ไวรัสนี้มีการแพร่เชื้อระหว่างคนเหมือนกับไข้หวัดใหญ่ ซึ่งไวรัส 2019-nCoV ระบาดทั่วประเทศจีนภายในเวลาเพียงไม่กี่สัปดาห์ และยังมีการระบาดไปยังประเทศต่าง ๆ ในทุกทวีปทั่วโลกภายในเวลาไม่นาน การแพร่ระบาดของเชื้อไวรัส 2019-nCoV นี้เป็นการระบาดใหญ่ (pandemic) ที่สร้างผลกระทบในทุกประเทศทั่วโลก ดังภาพประกอบที่ 1 ซึ่งครั้งนี้จัดเป็นการระบาดใหญ่ครั้งที่ 5 หลังจากการระบาดใหญ่ครั้งที่ 4 เมื่อปี ค.ศ. 2009 ซึ่งเกิดจากเชื้อไวรัสไข้หวัดใหญ่สายพันธุ์ H1N1

โดยสถานการณ์ในประเทศไทย ในช่วงเดือนมิถุนายน 2023 พบจำนวนผู้ป่วยรายใหม่ที่รักษาในโรงพยาบาล จำนวน 3,085 ราย และพบเสียชีวิตจำนวน 68 ราย ซึ่งผู้เสียชีวิตทั้งหมดเป็นผู้ที่ไม่ได้รับวัคซีนครบ 2 เข็ม หรือไม่ได้รับวัคซีนเข็มกระตุ้น ดังนั้นการฉีดวัคซีนเข็มกระตุ้นหรือวัคซีนประจำปี จึงเป็นสิ่งสำคัญในการลดอาการป่วยหนักและเสียชีวิตจากโควิด 19 ได้

วัคซีนโควิด 19 มีประสิทธิภาพในการป้องกันการติดเชื้อไม่สมบูรณ์ ช่วง 4 เดือนแรก และสำหรับปัสายพันธุ์กลายพันธุ์ จะต้องทำการฉีดอย่างน้อย 3 เข็ม ถึงจะสามารถป้องกันความรุนแรงได้ระยะเวลานาน ป้องกันการเสียชีวิต อีกทั้งวัคซีนยังช่วยลดการเกิดภาวะ Long COVID ได้ (เจาะลึกระบบสุขภาพ, 2023)

ดังนั้นการแพร่ระบาดของโควิดจะต้องมีวัคซีนที่มีประสิทธิภาพซึ่งสามารถกระจายได้อย่างเท่าเทียมกัน วัคซีนเอ็มอาร์เอ็นเอเป็นหนึ่งในตัวเลือกของการทำวัคซีนที่เร็วที่สุดสำหรับโรคโควิด ข้อจำกัดที่เกิดขึ้นอย่างหนึ่งในตอนนี้คือการออกแบบโมเลกุลเอ็มอาร์เอ็นเอ ที่มีความเสถียรสูง วัคซีนทั่วไป (เช่นไข้หวัดใหญ่ตามฤดูกาล) วัคซีนทั่วไปสามารถบรรจุในกระบอกฉีดยาแบบใช้แล้วทิ้งและจัดส่งผ่านตู้เย็นทั่วโลก แต่ปัจจุบันวัคซีนเอ็มอาร์เอ็นเอยังไม่สามารถทำได้เนื่องจากไม่สามารถรักษาความเสถียรจากคุณสมบัติการย่อยสลายตามธรรมชาติของตัววัคซีนเอง



ภาพประกอบ 1 แผนผังการแพร่ระบาดครั้งใหญ่ของเชื้อไวรัสไข้หวัดใหญ่ และเชื้อ SARS-CoV-2

ที่มา : (รัตนโรจน์พงศ์, 2021)

ในการระบาดของโรคโควิด-19 วัคซีนเป็นส่วนสำคัญในการช่วยยับยั้งและป้องกันการระบาดได้มากขึ้น โดยดูจาก Neutralizing Antibody (ความสามารถในการสร้างภูมิ) และดูจากความสามารถในการป้องกันโรคที่เกิดขึ้น โดยพบว่าวัคซีน mRNA ได้แก่ Pfizer และ Moderna มีประสิทธิภาพเกิน 90% , วัคซีน AstraZeneca ประสิทธิภาพ ประมาณ 80% , วัคซีน Johnson&Johnson ประสิทธิภาพ 65-70% และ วัคซีน Sinovac ประสิทธิภาพประมาณ 50-70% (ข้างแก้ว)

สำหรับวัคซีน mRNA ผลิตด้วยเทคโนโลยีใหม่ล่าสุด ทำให้สามารถผลิตได้ในปริมาณมาก ประสิทธิภาพสูงและมีราคาค่อนข้างถูกแต่มีข้อเสียที่ต้องเก็บในอุณหภูมิต่ำกว่าวัคซีนทั่วไป (-20°C ขึ้นไป)

ชุมชน Eterna เป็นแพลตฟอร์มในการออกแบบข้อมูลของอาร์เอ็นเอโดยโซลูชันดังกล่าวได้รับการสังเคราะห์และทดสอบโดยนักวิจัยที่ Stanford เพื่อให้ได้ข้อมูลเชิงลึกใหม่ ๆ เกี่ยวกับโมเลกุลอาร์เอ็นเอการปรับปรุงความเสถียรของวัคซีนเอ็มอาร์เอ็นเอ โดยส่งมอบข้อมูลในทาง Kaggle ซึ่งเป็นแพลตฟอร์มสำหรับการแข่งขันการเรียนรู้โดยเครื่องจักรเพื่อให้ทุกคนได้ร่วมแข่งขันในการออกแบบแบบจำลองในการหาอัตราการย่อยสลายของวัคซีน การแก้ปัญหาที่ความท้าทายทางวิทยาศาสตร์ ภายในไม่กี่เดือนเพื่อเร่งการวิจัยวัคซีน เอ็มอาร์เอ็นเอ และส่งมอบวัคซีนป้องกัน

จากที่กล่าวมาในข้างต้นทางผู้วิจัยสนใจที่จะศึกษาการทำนายอัตราการย่อยสลายที่เป็นไปได้ที่แต่ละฐานของโมเลกุล RNA ซึ่งได้รับการฝึกฝนจากชุดข้อมูลย่อยของชุดข้อมูล Eterna

ซึ่งประกอบด้วยโมเลกุล RNA มากกว่า 3,000 โมเลกุลของโรค covid-19โดยใช้เทคนิคการเรียนรู้ของเครื่องจักรเพื่อเพิ่มโอกาสในการนำไปทดสอบเพื่อป้องกันโรสดังกล่าวต่อไป

### วัตถุประสงค์

1. เพื่อศึกษาการสร้างแบบจำลองการทำนายค่าความเสถียรของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด 19 ด้วยวิธีการเรียนรู้ของเครื่องจักร
2. เพื่อเปรียบเทียบการทำนายค่าความเสถียรของโมเดลการเรียนรู้เครื่องจักร
3. หาปัจจัยที่ส่งผลต่อความเสถียรของวัคซีน

### ขอบเขตงานวิจัย

1. ศึกษาโดยใช้ข้อมูลของโครงสร้างอาร์เอ็นเอโดยโมเดลจะบอกว่ารูปร่างของอาร์เอ็นเอแบบใดที่จะมีความเสถียรซึ่งคืออัตราการย่อยสลายจากสภาวะภายนอก
2. ทำนายค่าความเสถียรของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด 19 มีทั้งหมด 5 เอาท์พุท คือ ค่า reactivity, ค่า deg\_pH10 ค่า deg\_mg\_pH10,ค่า deg\_50C และ ค่า deg\_mg\_50C ซึ่งเป็นอัตราการย่อยสลายในสภาวะต่างๆ โดยใช้หลักการของ mean columnwise root mean squared error (MCRMSE)

### วิธีการดำเนินการวิจัย

1. ทำ Literature Review ที่เกี่ยวข้อง
2. หาข้อมูลเพื่อใช้ในงานวิจัย
3. เตรียมข้อมูล
4. ศึกษาวิธีการสร้างโมเดลในการทำนายข้อมูล
5. นำค่าวัดประสิทธิภาพของโมเดล
6. ประเมินและสรุปงานวิจัย

### ประโยชน์ที่คาดว่าจะได้รับจากงานวิจัย

1. วิธีการเรียนรู้ของเครื่องจักรสามารถนำมาใช้ในการพัฒนาทางการแพทย์ได้
2. สามารถศึกษาการทำนายค่าความเสถียรของวัคซีนเอ็มอาร์เอ็นเอสำหรับไวรัสโควิด
3. สามารถนำข้อมูลที่ได้จากการศึกษาไปต่อยอดในงานวิจัยอื่นๆต่อไป



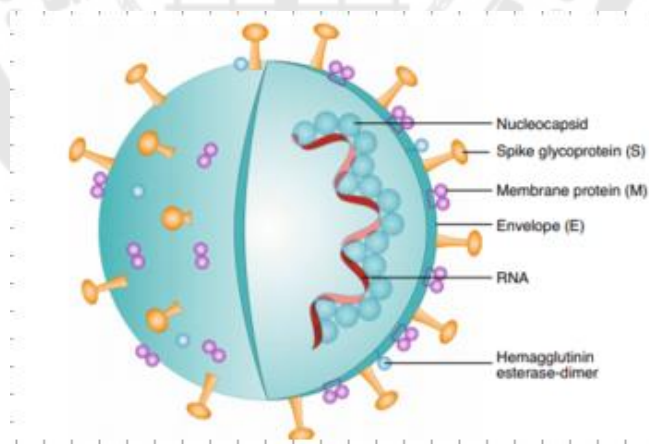
## บทที่ 2 วรรณกรรมและงานวิจัยที่เกี่ยวข้อง

### โรคโควิด 19 (Covid-19)

โรคโควิด 19 คือโรคเกิดจากไวรัสโคโรนาที่มีการค้นพบในเดือนธันวาคมปี 2019 โดยผู้ป่วยที่พบป่วยด้วยโรคปอดบวม ที่เมืองอู่ฮั่น ประเทศจีน ต่อมานักวิจัยชาวจีนพบว่าโรคติดเชื้อเชื้อไวรัสโคโรนา ชนิดใหม่ (novel coronavirus) และตั้งชื่อว่า “2019 novel Coronavirus” หรือเรียก อย่างย่อว่า “2019-nCoV “ โดยจะสามารถแพร่จากคนไปสู่คนได้โดยตรง (human-to-human transmission) และแพร่ระบาดไปทั่วโลกได้ โดยแพร่จากคนสู่คนผ่านทางฝอยละอองที่ออกจากจมูกหรือปากด้วยการ ไอหรือจาม โดยผู้ป่วยจะเริ่มมีอาการหลังจากติดเชื้อไม่เกิน 14 วัน

### ไวรัสโคโรนา 2019

ไวรัสโคโรนาเป็นไวรัสในวงศ์ใหญ่มีหลากหลายสายพันธุ์ทำให้เกิดโรกระบบทางเดินหายใจ เช่น โรคทางเดินหายใจตะวันออกกลาง (MERS) และโรกระบบทางเดินหายใจเฉียบพลันร้ายแรง (SARS) และโรคติดเชื้อไวรัสโคโรนา 2019 หรือโควิด 19 ลักษณะโครงสร้างไวรัสโดยรวมจะเป็นทรงกลมมีหนามแหลม

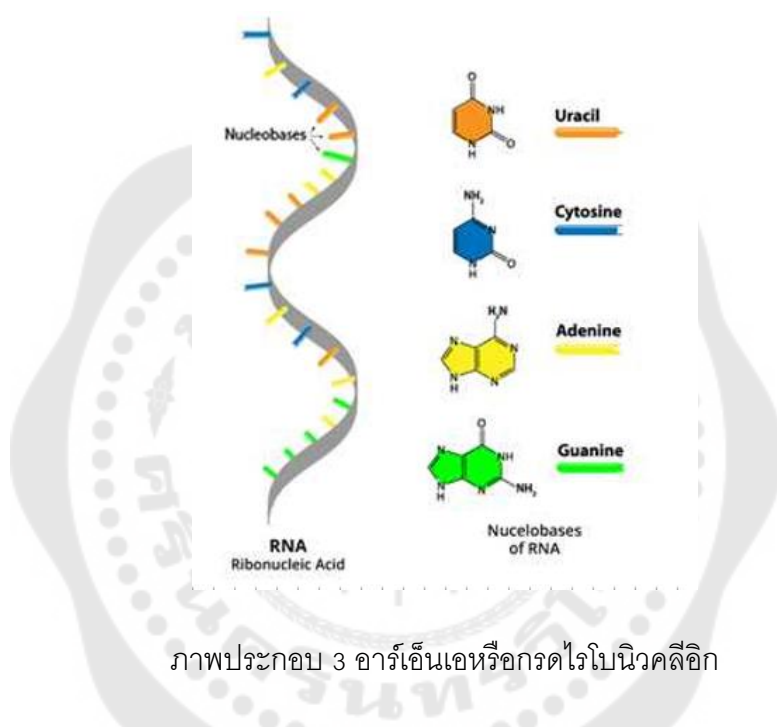


ภาพประกอบ 2 โครงสร้างไวรัสโคโรนา

ที่มา: (Florindo และคนอื่น ๆ, 2020)

## อาร์เอ็นเอ (RNA)

อาร์เอ็นเอ (RNA) หรือ กรดไรโบนิวคลีอิก (Ribonucleic acid – RNA) เป็นสายพอลิเมอร์ของนิวคลีโอไทด์ (Nucleotide) ไม่มีการแตกกิ่งก้านและมีความยาวสั้นกว่าโมเลกุลของ DNA มีส่วนประกอบด้วย น้ำตาลไรโบส (Ribose) และเบส 4 ชนิด คืออะดีนีน (Adenine, A) , ยูราซิล (Uracil, U) , ไซโตซีน (Cytosine, C) และกัวนีน (Guanine, G) และฟอสเฟต



ภาพประกอบ 3 อาร์เอ็นเอหรือกรดไรโบนิวคลีอิก

ที่มา : (Mackenzie, 2020)

โดยอาร์เอ็นเอ (RNA) จะเป็นโพลีนิวคลีโอไทด์สายเดี่ยว (single strand) ซึ่งนิวคลีโอไทด์ (Nucleotide) เชื่อมต่อกันด้วยพันธะฟอสโฟไดเอสเทอร์ (phosphodiester bond) อาร์เอ็นเอ (RNA) จะมีการถอดรหัส (transcription) จากดีเอ็นเอ (DNA) ด้วย RNA Polymerase ซึ่งตัวอาร์เอ็นเอ (RNA) เหมือนเป็นแม่แบบ (Template) ในการแปลข้อมูลจากยีนไปเป็นโปรตีน แล้วขนย้ายกรดอะมิโนเข้าไปใน ribosome ของเซลล์ เพื่อผลิตโปรตีน และแปลรหัส (translation) เป็นข้อมูลในโปรตีน

อาร์เอ็นเอ (RNA) มีทั้งหมด 3 ชนิด คือ เอ็มอาร์เอ็นเอ หรือ เมสเซนเจอร์ อาร์เอ็นเอ (messenger RNA, mRNA) , อาร์อาร์เอ็นเอ หรือ ไรโบโซมอล อาร์เอ็นเอ (ribosomal RNA, rRNA) และทีอาร์เอ็นเอ หรือ ทรานสเฟอร์ อาร์เอ็นเอ (transfer RNA, tRNA)

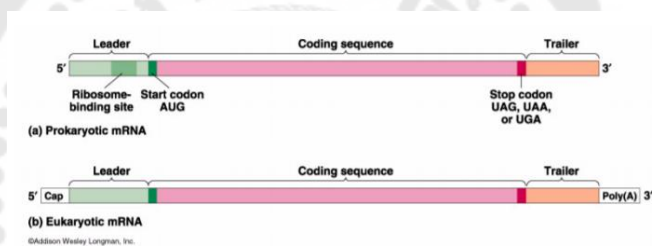
## เอ็มอาร์เอ็นเอ (mRNA)

เอ็มอาร์เอ็นเอ (mRNA) หรือ เมสเซนเจอร์ อาร์เอ็นเอ (messenger RNA) เป็น อาร์เอ็นเอ (RNA) ชนิดหนึ่ง โดยมีหน้าที่การเรียงลำดับของนิวคลีโอไทด์ที่ Complement กับข้อมูลบริเวณจำเพาะของ DNA ซึ่ง mRNA ทำหน้าที่เป็น template ในการสร้างโปรตีน

โครงสร้างของเอ็มอาร์เอ็นเอ ประกอบด้วย 3 ส่วน

1. Leader sequence หรือ 5 untranslated leader sequence (5/ULR) เป็นลำดับเบสจำเพาะที่อยู่บริเวณปลาย 5' P ของ mRNA เกาะกับสาย mRNA เพื่อเริ่มกระบวนการแปลรหัส
2. Coding sequence บริเวณกำหนดชนิดของกรดอะมิโน ในกระบวนการแปลรหัส
3. trailer sequence หรือ 3 untranslated trailer sequence (3/ UTR) คือ ตำแหน่งปลาย

3



ภาพประกอบ 4 แสดงโครงสร้างโมเลกุล mRNA ที่สมบูรณ์

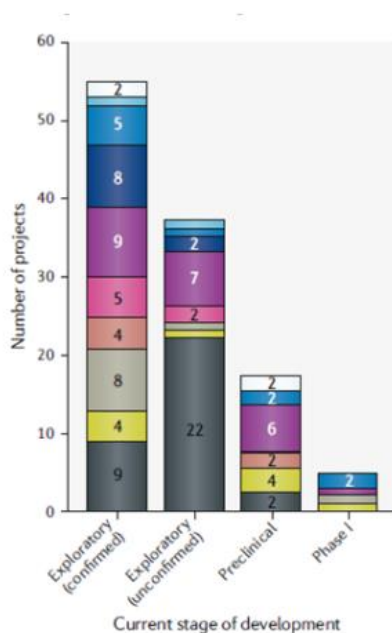
ที่มา : ("โมเลกุลอาร์เอ็นเอ (RNA Molecules) ")

## วัคซีน (vaccine)

วัคซีน คือสารชนิดหนึ่งที่ถูกฉีดเข้าไปร่างกาย เพื่อสร้างภูมิคุ้มกันโรคต่างๆ ทำมาจากเชื้อโรค แบ่งเป็น 2 ประเภท คือ ทำจากเชื้อโรคที่ตาย และ ทำจากเชื้อโรค เมื่อฉีดเข้าไปในร่างกายก็จะสร้างภูมิคุ้มกัน กับโรคนั้นๆ(มหาวิทยาลัยมหิดล, 2020)

## การพัฒนาวัคซีน COVID-19

ข้อมูลด้านพันธุกรรมของไวรัส SARS-CoV-2 ได้รับการเผยแพร่ในวันที่ 11 มกราคม 2020 ซึ่งส่งผลให้เกิดการเริ่มต้นพัฒนาวัคซีนในหลายๆที่ทั่วโลก โดยมีรายงานในเดือนเมษายน 2020 ว่า ที่อยู่ในกระบวนการพัฒนาโดยที่วัคซีน 78 ชนิดมีการระบุอย่างชัดเจนในการพัฒนา (รัตนโรจน์พงศ์, 2021)โดยมีรายละเอียดชนิดของวัคซีนที่พัฒนาดังแสดงในรูปที่ 5



ภาพประกอบ 5 ชนิดของวัคซีน COVID-19 ที่มีการพัฒนา

ที่มา: (รัตนโรจน์พงศ์, 2021)

สำหรับการพัฒนาวัคซีนในประเทศไทยนั้นมีทั้งหมด 7 เทคโนโลยีการผลิต ได้แก่

1. mRNA Vaccine ได้มาจากการสังเคราะห์เอ็มอาร์เอ็นเอ ที่เฉพาะเจาะจงกับเชื้อไวรัส วัคซีนจะทำหน้าที่พาเอ็มอาร์เอ็นเอเข้าเซลล์ และทำให้เซลล์ผลิตสารโปรตีนสไปค์ของเชื้อไวรัสโดยโปรตีนนี้จะกระตุ้นระบบภูมิคุ้มกันของร่างกาย

2. DNA Vaccine การนำสารพันธุกรรมดีเอ็นเอฉีดเข้ากล้ามเนื้อ สามารถที่จะทำให้เกิดการแสดงออกของยีนได้ ดีเอ็นเอวัคซีนประกอบด้วยยีน ที่ต้องการซึ่งสามารถสร้างโปรตีนหรือแอนติเจนที่ได้รับการ ดัดแปลงอย่างเหมาะสมภายใต้การควบคุมของส่วนที่เป็น regulatory sequences โปรตีนที่สร้างจากดีเอ็นเอวัคซีนจะ อยู่ในลักษณะ native conformation

3. Viral-like particle (VLP) วัคซีนอนุภาค ไวรัสเสมือนเป็นเทคโนโลยีการสร้างโครงสร้างเลียนแบบอนุภาคไวรัส แต่ไม่มีสารพันธุกรรมของไวรัส

4. Protein Sub Unit Vaccine วัคซีนที่ผลิตโดยการสร้างโปรตีนของเชื้อไวรัสด้วยระบบ cell culture, yeast เป็นต้น ใช้ในการผลิตวัคซีนหลายชนิด เช่น วัคซีนป้องกันไข้หวัดใหญ่ วัคซีนป้องกันไวรัสตับอักเสบบี เป็นต้น (ศักดิ์ทองจีน, 2021)

5. Viral vector Vaccine สร้างจากไวรัสที่ถูกทำให้อ่อนลงหรือไม่สามารถแบ่งตัวได้อีก ตัดแต่งพันธุกรรมเพื่อใช้เป็นพาหะ สามารถกระตุ้นภูมิคุ้มกันได้ดี เนื่องจากเลียนแบบการติดเชื้อที่ใกล้เคียงกับการติดเชื้อตามธรรมชาติ (โรงพยาบาลศิรินครินทร์)

6. Inactivated Vaccine วัคซีนกลุ่มนี้ผลิตโดยนำไวรัสมาเลี้ยงขยายจำนวนมาก และทำให้เชื้อตาย

7. Live attenuated Vaccine วัคซีนที่ผลิตขึ้นโดยใช้เชื้อโรคมาทำให้อ่อนฤทธิ์ลงจนไม่สามารถทำให้เกิดโรคแต่เพียงพอที่จะกระตุ้นภูมิคุ้มกันของร่างกายได้

โดย หน่วยงานที่เกี่ยวข้องทั้งหมด 9 หน่วยงาน ทั้งภาครัฐ รัฐวิสาหกิจ เอกชน และ มหาวิทยาลัย โดยการวิจัยที่มี ความก้าวหน้ามากที่สุด และอยู่ระหว่างเตรียมเข้าสู่การทดสอบในมนุษย์ระยะที่ 1 นั้นมี 3 เทคโนโลยีการผลิต ได้แก่

1. วัคซีนชนิดเอ็มอาร์เอ็นเอซึ่งพัฒนาโดยศูนย์เชี่ยวชาญเฉพาะทางการวิจัยและพัฒนาวัคซีน คณะแพทยศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย
2. วัคซีนชนิด DNA โดยบริษัท ไบโอเนท-เอเชีย จำกัด
3. วัคซีนที่ใช้เทคโนโลยีการ สกัดโปรตีนจากพืช (Plant based) โดยคณะเภสัชศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย และบริษัท ไบยาไฟโต ฟาร์ม จำกัด

การวิจัยการพัฒนาวัคซีน COVID-19 ในอดีตการพัฒนาวัคซีนป้องกันโรคติดเชื้อในกลุ่ม coronavirus เช่น SARS-CoV ไม่ค่อยสัมฤทธิ์ ผลเนื่องจากวัคซีนยังไม่มีประสิทธิภาพที่ดีพอในการป้องกันโรค ถึงแม้ว่าวัคซีนที่พัฒนาขึ้นมา มีความสามารถในการกระตุ้นการตอบสนองของระบบภูมิคุ้มกันในสัตว์ทดลองและสิ่งที่กังวลคือวัคซีนอาจไม่นำไปสู่การป้องกันโรคได้ในระยะยาว และส่งผลก่อให้เกิดการติดเชื้อไวรัสซ้ำได้

ดังนั้นการพัฒนาวัคซีน COVID-19 จึงต้องคำนึงถึงความปลอดภัยของวัคซีนในการที่จะไม่ก่อให้เกิดพยาธิ สภาพจากการตอบสนองของระบบภูมิคุ้มกัน การป้องกันโรคได้อย่างยาวนาน และไม่ก่อให้เกิดการติดเชื้อซ้ำใน ผู้ที่ได้รับวัคซีน โดยศึกษาการพัฒนาวัคซีนในรูปแบบต่างๆที่สามารถนำไปสู่การผลิตวัคซีนที่มีประสิทธิภาพเพื่อ ลดข้อจำกัดของวัคซีนที่เคยมีงานวิจัยมาก่อนหน้านี้ในการพัฒนาวัคซีนต่อ SARS-CoV และ MERS-CoV ความปลอดภัยของผู้ที่ได้รับวัคซีนเป็นสิ่งที่สำคัญที่สุดในการพัฒนาวัคซีนเพื่อป้องกันโรคต่างๆ รวมทั้งวัคซีนป้องกันโรค COVID-19 เนื่องจากมีรายงาน การพัฒนาวัคซีนป้องกันไวรัส SARS-CoV และ MERS-CoV ที่พบว่า วัคซีนก่อให้เกิดพยาธิสภาพที่ไม่พึงประสงค์โดยไปก่อให้เกิดพยาธิในปอด

อย่างไรก็ตามยังไม่มีรายงานว่า วัคซีนต่อ COVID-19 ที่พัฒนาและกำลังดำเนินการทดสอบทางคลินิกอยู่นั้นบางชนิดได้ผ่านการรับรองเพื่อนำมาใช้ในการป้องกันการระบาดของโรค สามารถส่งผลทำให้เกิดการตอบสนองของระบบภูมิคุ้มกันที่มี ผลกระทบต่อผู้ได้รับวัคซีนหรือไม่ อย่างไรก็ตามประชากรทั้งหมดที่ได้รับวัคซีน โดยเฉพาะผู้สูงอายุและผู้ที่มีโรคแทรกซ้อนเช่นโรคเบาหวาน ความดันสูงหรือโรคหัวใจที่มีความเสี่ยงที่จะมีอาการรุนแรงถ้าติดเชื้อไวรัส SARS-CoV-2 ซึ่งการทดสอบทาง preclinic ในกลุ่มผู้สูงอายุหรือการใช้สัตว์ทดลองที่มีอายุจึงไม่ควรมองข้ามเพื่อยืนยันว่า วัคซีนที่พัฒนา มีความปลอดภัยในกลุ่มเสี่ยงนี้ รวมทั้งองค์ความรู้ในการตอบสนอง ของระบบภูมิคุ้มกันในกลุ่มนี้ต่อวัคซีนแต่ละประเภทที่พัฒนาขึ้นเพื่อให้สามารถพัฒนาวัคซีนให้มีความปลอดภัย และสามารถกระตุ้นระบบภูมิคุ้มกันที่นำไปสู่การป้องกันโรคของคนกลุ่มนี้ได้ อย่างมีประสิทธิภาพ ในปัจจุบันได้มีการนำเอาเทคโนโลยีใหม่ๆมาใช้ในการพัฒนาวัคซีนในรูปแบบต่างๆ

จากการร่วมมือกันหน่วยงานต่างๆระหว่างประเทศเพื่อให้ได้วัคซีนที่มีประสิทธิภาพดีที่สุด และผลิตได้เร็วที่สุดเพื่อนำมาใช้ในการป้องกันการแพร่ระบาดของโรค COVID-19 ที่ยังเกิดขึ้นอย่างต่อเนื่องและไม่มีท่าทีว่าการระบาดจะเริ่มชะลอตัวหรือลดน้อยลงโดยวัคซีนที่พัฒนาขึ้นจากเทคโนโลยีใหม่ๆนี้ เช่น nucleic acid vaccine (อาจเป็นชนิด DNA หรือ RNA), viral vector vaccine (ใช้เชื้อไวรัสชนิดอื่นๆให้มีการแสดงออกของโปรตีนจากไวรัส SARS-CoV-2), Inactivated virus vaccine, attenuated virus vaccine, virus like particle vaccine (VLP) และ protein subunit/recombinant protein subunit vaccine (ใช้เทคนิคทางด้านพันธุวิศวกรรมในการผลิตโปรตีนของไวรัสที่เป็นเป้าหมายในการกระตุ้นการทำงานของระบบภูมิคุ้มกัน)

### วัคซีนเอ็มอาร์เอ็นเอ

วัคซีนเอ็มอาร์เอ็นเอเป็นพื้นฐานทางเทคโนโลยีของการบำบัดรักษาและวัคซีนนั้นมีความยืดหยุ่นสูงในด้านการผลิตและการใช้งานโปรตีนใด ๆ สามารถเข้ารหัสและแสดงโดยเอ็มอาร์เอ็นเอเป็นทางเลือกที่ดีสำหรับวิธีการฉีดวัคซีนแบบเดิมเนื่องจากมีศักยภาพสูงความสามารถในการพัฒนาอย่างรวดเร็วและศักยภาพในการผลิตต้นทุนต่ำ (Schlake, Thess, Fotin-Mleczek, และ Kallen, 2012) วัคซีนเอ็มอาร์เอ็นเอได้กระตุ้นภูมิคุ้มกันที่มีศักยภาพต่อเป้าหมายของโรคติดเชื้อ เช่น ไวรัสไข้หวัดใหญ่ หรือมะเร็ง โดยเฉพาะอย่างยิ่งในช่วงไม่กี่ปีที่ผ่านมาโดยใช้เอ็มอาร์เอ็นเอ ที่ได้รับการปรับแต่งตามลำดับด้วยไขมันหรือรูปแบบเปลือยเปล่า

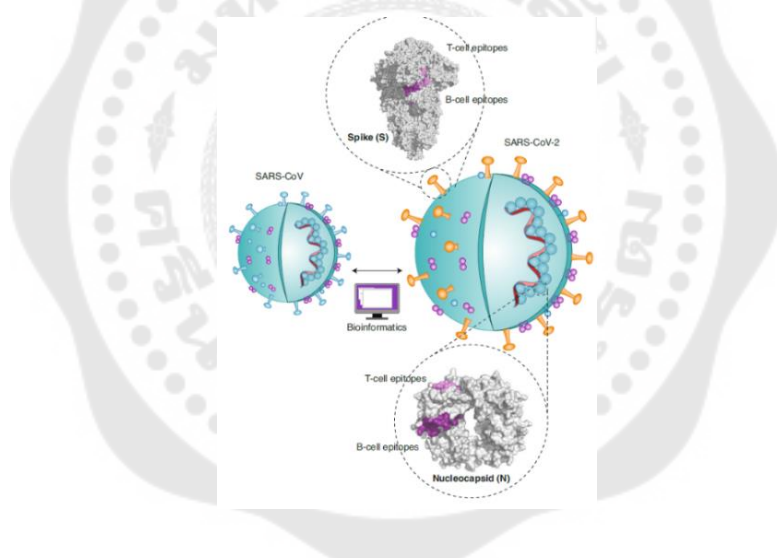
แนวทางที่หลากหลายในการฉีดวัคซีนเอ็มอาร์เอ็นเอเพื่อยับยั้งมะเร็งรวมถึงวัคซีนเซลล์เดนไดรติกและ วัคซีนเอ็มอาร์เอ็นเอชนิดฉีดโดยตรงหลายชนิดได้ถูกนำไปใช้ในการทดลองทางคลินิก

มะเร็งจำนวนมากโดยมีผลการทดลองบางอย่างที่แสดงให้เห็นถึงการตอบสนองของ T ของแอนติเจนที่จำเพาะต่อแอนติเจนและการอยู่รอดโดยปราศจากโรคเป็นเวลานานในบางกรณี

โดยในขณะนี้การผลิตวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด 19 มีบริษัทผู้ผลิตยาโมเดอร์นาของสหรัฐอเมริกาเปิดเผยผลเบื้องต้นจากการทดลองทางคลินิกกับผู้ทดลองมากกว่า 30,000 คน ได้ผลในการป้องกันไวรัสถึง 94.5 % และคงสภาพได้เมื่อเก็บในอุณหภูมิระหว่าง 2-8 องศาเซลเซียส (ไทยโพสต์, 2020)

### การทำงานของวัคซีนเอ็มอาร์เอ็นเอ

การทำงานของวัคซีนเอ็มอาร์เอ็นเอไปกระตุ้นให้ร่างกายสร้างโปรตีนที่มีรูปทรงแหลมคล้าย "หนามแหลม" (Spike) บนไวรัส SARS-CoV-2 เมื่อเราฉีดวัคซีนเข้าสู่ร่างกายจะทำให้ระบบภูมิคุ้มกันในร่างกายของเราเรียนรู้รูปร่างของไวรัสชนิดนี้ไว้ล่วงหน้าก่อนที่จะโดนไวรัสตัวจริงเข้าสู่ร่างกาย (Florindo และคนอื่น ๆ, 2020)



ภาพประกอบ 6 การทำงานของวัคซีนเอ็มอาร์เอ็นเอต่อไวรัส

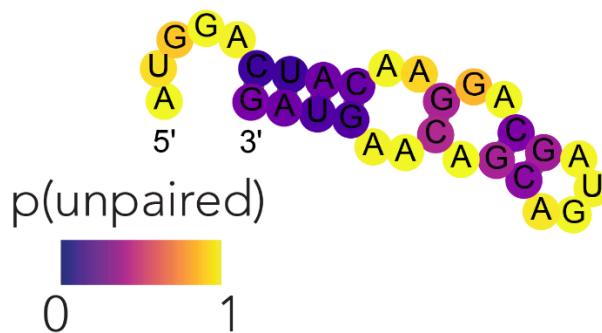
ที่มา: (Florindo และคนอื่น ๆ, 2020)

### การออกแบบวัคซีนเอ็มอาร์เอ็นเอที่มีความเสถียร

จากบทความของชุมชน Eterna ที่อธิบายเกี่ยวกับวัคซีนอาร์เอ็นเอ โดย RNA จะสร้างคู่เบสเช่นเดียวกับ DNA แทนที่จะเป็นคู่เบส G-C และคู่เบส A-T ใน DNA RNA จะเป็นคู่เบส G-C และ A-U

RNA แบบจับคู่ฐานมีแนวโน้มที่จะย่อยสลายน้อยกว่า RNA แบบเกลียวเดี่ยว โดยแสดงตัวอย่างตาม รูปที่ 6 เอ็มอาร์เอ็นเอ ของโปรตีนขนาดเล็กโดยแต่ละฐานจะมีสีตามความน่าจะเป็น

ที่ไม่ได้จับคู่ เราจะเห็นว่าโมเลกุลมี "จุดร้อน" เป็นสีเหลืองซึ่งมีแนวโน้มที่จะย่อยสลายได้มากกว่า แต่ยังคงมีความคลาดเคลื่อนจากการทดลองจริงอยู่อีกพอสมควร (Eternagame, 2020)



ภาพประกอบ 7 mRNA of a mini protein, with bases colored by the probability they are unpaired.

ที่มา : (Eternagame, 2020)

จากการอ้างอิงงานวิจัยของ Wadhwa และคณะ (Wadhwa A, 2020) อ้างว่าภายใต้เงื่อนไขการขนส่งแบบ cool chain ในตู้เย็น ครึ่งชีวิตของวัคซีน mRNA อาจมีครึ่งชีวิต 900 วัน โดยมีอัตราการย่อยสลายอย่างน้อย 2% ทุกๆ 30 วัน ยิ่งกว่านั้นที่อุณหภูมิประมาณ 37° เชื่อกันว่าจะลดครึ่งชีวิตของ วัคซีนถึง 5 วันและ 10 วันตามลำดับ นอกเหนือไปจากนี้, เป็นที่น่าสังเกตว่าเมื่อมี Mg<sup>2+</sup> ที่อุณหภูมิ 37°C ซึ่งเป็นสภาวะตามปกติสำหรับการทดสอบในหลอดทดลอง ครึ่งชีวิตของ วัคซีนจะลดลงเหลือไม่เกิน 2 ชั่วโมงและผลลัพธ์ยังคงเหมือนเดิมแม้หลังจากลดความเข้มข้นของ Mg<sup>2+</sup>, ค่า pH หรืออุณหภูมิเพื่อลดการไฮโดรไลซิสซึ่ง วัคซีนเอ็มอาร์เอ็นเอจะไม่เสถียรเนื่องจากการย่อยสลายยังคงเกิดขึ้นในระหว่างกระบวนการถอดรหัส (transcription) (S. H. Ing, 2021)

### ปัญหาเกี่ยวกับวัคซีนเอ็มอาร์เอ็นเอ

จากการออกแบบโมเลกุลอาร์เอ็นเอของสารที่มีความเสถียรสูง โมเลกุลของอาร์เอ็นเอมีแนวโน้มที่จะย่อยสลายเองตามธรรมชาติเป็นข้อจำกัด วัคซีนเอ็มอาร์เอ็นเอในปัจจุบันสำหรับป้องกันโรคโควิด 19 จะต้องได้รับการเตรียมและจัดส่งภายใต้การแช่เย็นที่เข้มข้น วัคซีนเอ็มอาร์เอ็นเอเผชิญกับความท้าทายเนื่องจากความไม่เสถียรทางความร้อน ซึ่งทำให้ไวต่อการย่อยสลายทางเคมี ด้วยเหตุนี้ วัคซีน เอ็มอาร์เอ็นเอ จึงต้องมีเงื่อนไขที่เข้มงวดในการผลิต การจัดเก็บ และการ



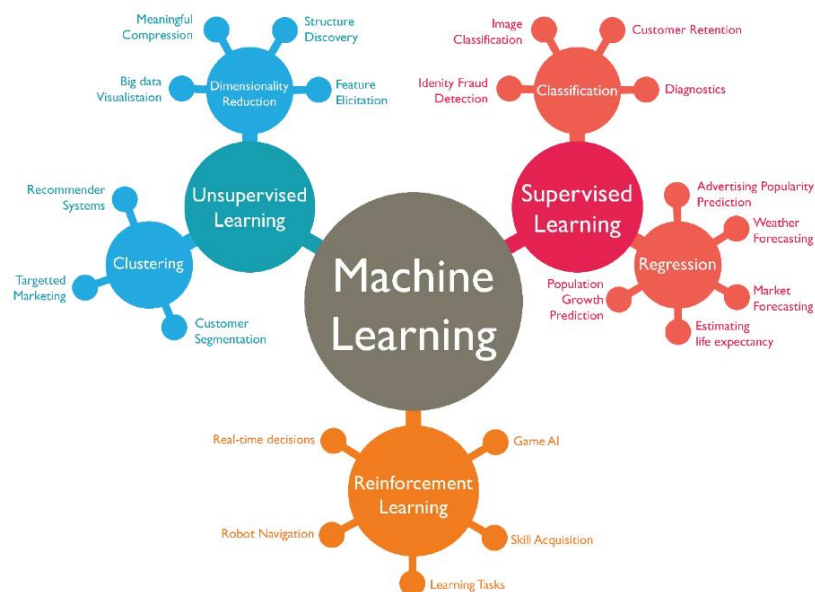
จัดส่งทั่วโลก ในการทำให้วัคซีน เอ็มอาร์เอ็นเอ สามารถเข้าถึงได้ในวงกว้างมากขึ้น จำเป็นอย่างยิ่งที่จะต้องทำความเข้าใจและปรับปรุงความเสถียรของวัคซีน

นอกจากเหตุผลทางด้านอุณหภูมิเอ็มอาร์เอ็นเอมีความไวสูงต่อเอนไซม์ RNase ซึ่งย่อยสลายได้ง่าย ไวต่อปฏิกิริยาออกซิเดชันมากกว่า ไวต่อการไฮโดรไลซิสที่ค่า pH สูงกว่า 6 สำหรับวัคซีน RNA เป็นสิ่งสำคัญในการป้องกันความสมบูรณ์ของโมเลกุลที่สมบูรณ์เนื่องจากความผิดปกติเพียงครั้งเดียวสามารถสลายได้ (Uddin MN, 2021)



## การเรียนรู้ของเครื่องจักร (Machine Learning)

การเรียนรู้ของเครื่องจักร (Machine Learning) คือระบบที่สามารถเรียนรู้ได้จากตัวอย่างด้วยตนเองโดยไม่จำเป็นต้องมีคำสั่งของโปรแกรมเมอร์ สามารถเรียนรู้เพียงแค่จากข้อมูลอย่างเดียว แบ่งออกเป็น 3 ประเภท คือ Supervised Learning, Unsupervised Learning, Reinforcement Learning

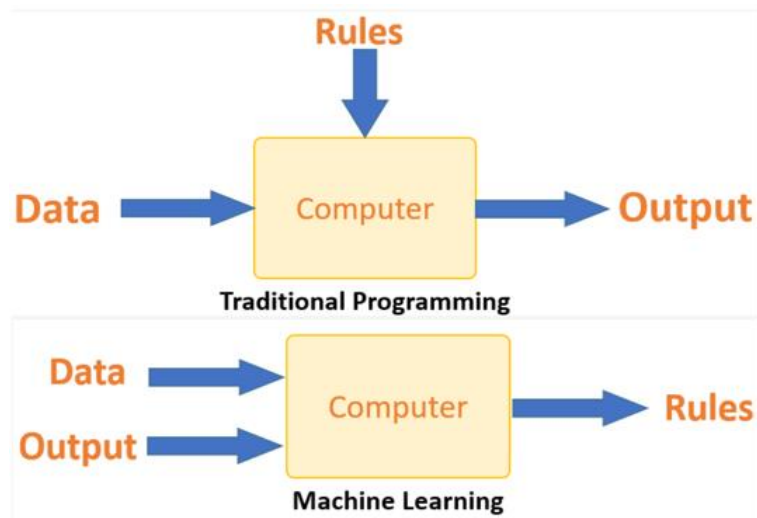


ภาพประกอบ 8 ประเภทของการเรียนรู้โดยเครื่องจักร

ที่มา : <https://medium.com/investic/machine-learning>

### ความแตกต่างระหว่าง Machine Learning และการเขียนโปรแกรมดั้งเดิม

การเขียนโปรแกรมในสมัยก่อนนั้นได้ทั้งหมดจะต้องถูกกำหนดแนวทางไว้ชัดเจน ซึ่งแต่ ละกฎจะขึ้นอยู่กับพื้นฐานความเข้าใจด้านตรรกศาสตร์ (Logic Foundation) ทำให้ผลลัพธ์ ออกมาตามคำสั่ง เมื่อระบบเริ่มซับซ้อนมากขึ้นกฎจะมากขึ้นตาม ส่วนเครื่อง (machine) จะทำ การเรียนรู้ว่าข้อมูลขาเข้าและข้อมูลขาออกเกี่ยวข้องกับอย่างไร และจะเขียนกฎขึ้นมา โดย โปรแกรมเมอร์ไม่จำเป็นต้องเขียนกฎใหม่ทุกครั้งที่ อัลกอริทึมจะปรับเข้ากับข้อมูลใหม่เอง



ภาพประกอบ 9 ความแตกต่างระหว่าง Machine Learning และการเขียนโปรแกรมดั้งเดิม

### ประเภทของการเรียนรู้โดยเครื่องจักร

#### 1. การเรียนรู้แบบมีผู้สอน (Supervised Learning)

จะอยู่ในลักษณะการทำนายผลลัพธ์ จะมีชุดข้อมูลฝึกฝน (Training Data) เป็นตัวฝึก โดยมนุษย์มาคอยแยกประเภทหรือบอกผลลัพธ์ (Label) ไว้ โดยนำข้อมูลที่ใช้ฝึกไปผ่านอัลกอริทึม สำหรับสร้างแบบจำลอง จากนั้นนำข้อมูลที่เครื่องไม่เคยเห็นให้เครื่องทำนาย (predict)

ประเภทของ supervised learning อยู่ 2 ประเภท

- การแบ่งแยกประเภท (Classification)
- การถดถอย (Regression)

### การแบ่งแยกประเภท (Classification)

เป็นวิธีที่จะต้องมี Target ไว้สำหรับให้ตัว Model เรียนรู้จาก Input Data เพื่อหาคำตอบออกมาตาม Target ที่เราได้วางเอาไว้ ซึ่งผลลัพธ์จากการวิเคราะห์ข้อมูลด้วย Classification Model จะเป็นในรูปแบบของการจำแนกข้อมูลเพื่อให้ได้คำตอบที่เป็นตัวเลือก หรือกลุ่มข้อมูล ตัวอย่างของ Target ของโมเดล

### การถดถอย (Regression)

เป็นการศึกษาความสัมพันธ์ระหว่างลักษณะตัวแปรตั้งแต่สองตัวแปรขึ้นไปโดยหาสมการความสัมพันธ์ระหว่างตัวแปร 2 ตัวเพื่อที่จะนำไปสู่การคาดการณ์หรือประมาณค่าการวิเคราะห์พยากรณ์เหตุการณ์ในอนาคต เช่น การทำนายมูลค่าของหุ้น เป็นต้น

## 2. การเรียนรู้แบบไม่มีผู้สอน (Unsupervised Learning)

เป็นการให้เครื่องตามหาโครงสร้างข้อมูลที่เราไม่รู้จักร (Unknown Structure) โดยอัลกอริทึมจะตรวจสอบเฉพาะข้อมูลที่ป้อนเข้ามาเท่านั้นโดยปราศจากการให้ผลลัพธ์ที่จะเกิดขึ้น (เช่น การสำรวจข้อมูลประชากรเพื่อหาแบบแผน(pattern)ของข้อมูลนั้น)

### 3. การเรียนรู้แบบเสริมแรง (Reinforcement Learning)

การทำ Action ตามลำดับต่างๆ ออกมาแล้วคะแนนดี เราก็จะได้รางวัล (Reward) เครื่องจะเริ่มจดจำลำดับการทำ Action แล้วได้คะแนนดี และพยายามจะทำ Action นั้นเรื่อยๆ

#### ทฤษฎีของแบบจำลองที่นำมาใช้ในงานวิจัย

##### XGBoost (eXtreme Gradient Boosting)

XGBoost เป็นแบบจำลองที่นำเอาต้นไม้ตัดสินใจมาฝึกสอนต่อกันหลาย ๆ ต้น โดยที่ต้นไม้ตัดสินใจแต่ละต้นจะเรียนรู้จากค่าความผิดพลาดของต้นก่อนหน้าต่อเนื่องกันจนมีความลึกมากพอ แบบจำลองจะหยุดเรียนรู้เมื่อไม่เหลือค่าความผิดพลาดจากต้นไม้ตัดสินใจต้นก่อนหน้าให้เรียนรู้แล้ว

##### Random Forest

เป็นแบบจำลองที่นำ Decision Tree หลายๆ tree มา Train ร่วมกัน (ตั้งแต่ 10 ต้น ถึงมากกว่า 1000 ต้น) โดยที่แต่ละ tree จะได้รับ feature และ data เป็น subset ของ feature และ data ทั้งหมด แบบสุ่มตอนคาดการณ์ แต่ละ Decision Tree ทำการคาดการณ์ของตัวเอง และเลือกผล final prediction

##### Light GBM (Light Gradient Boosting)

Light GBM เป็น Gradient Boosting ใช้อัลกอริทึม Leaf-wise ซึ่งสามารถลดการสูญเสียได้ดีกว่าอัลกอริทึมแบบ Level-wise ให้ผลลัพธ์ที่แม่นยำและมีความเร็วมากกว่า

##### CatBoost (Gradient boosting with categorical features support)

เป็นอัลกอริทึมการเรียนรู้ของเครื่องจากที่เปิดให้ใช้งานแบบโอเพนซอร์ส CatBoost สามารถใช้งานร่วมกับเฟรมเวิร์กการเรียนรู้เชิงลึก เช่น Google's TensorFlow ได้ ซึ่ง โดยสามารถจัดการกับตัวแปรโดยอัตโนมัติ และไม่ต้องการการแปลงชุดข้อมูลให้เป็นรูปแบบเฉพาะ นอกจากนี้ CatBoost ยังสามารถจัดการกับตัวแปรและค่าข้อมูลบางส่วนที่หายไป (Missing Values) ได้อย่างมีประสิทธิภาพอีกด้วย

##### Cross validation

เป็นเทคนิคในการแบ่งข้อมูลเป็น k ส่วนเท่าๆกัน เพื่อสร้างและทดสอบแบบจำลอง (train + validate) คำนวณค่าเฉลี่ย accuracy หรือ error (i.e. model performance) ใช้เปรียบเทียบว่าข้อมูลชุดไหนดีที่สุดก่อนที่จะนำไปใช้ทำนายข้อมูลทดสอบ (test set) หรือ การเรียนรู้ (train) และทดสอบ (test) แบบจำลองในแต่ละรอบ เราสามารถทดสอบ hyper parameters หลายๆค่าพร้อม

กันได้ เพื่อหาค่าที่ดีที่สุดสำหรับแบบจำลองของเรา ให้เปรียบเทียบระหว่างแบบจำลองได้ว่าแบบจำลองไหนดีกว่ากัน

### งานวิจัยที่เกี่ยวข้อง

1. Elisabeth J. Wurtmann, Sandra L. Wolin (2009) ได้ศึกษาความเสียหายต่ออาร์เอ็นเอเกิดจากแสงอัลตราไวโอเล็ต, การออกซิเดชัน, คลอรีน, ไนเตรตและอะซิเลชันอาจรวมถึงการตัดแปลงทางเคมีของนิวคลีโอเบสเช่นเดียวกับการเชื่อมขวางอาร์เอ็นเอ-อาร์เอ็นเอและ อาร์เอ็นเอ-โปรตีน การศึกษาในหลอดทดลองได้อธิบายถึงปัจจัยที่สร้างความเสียหายได้หลายประเภทซึ่งบางส่วนได้รับการสนับสนุนว่าเกี่ยวข้องกับสรีรวิทยาโดยการสังเกตในร่างกายในการเจริญเติบโตปกติสภาวะความเครียดหรือสถานะของโรค ความเสียหายต่อทั้งเอ็มอาร์เอ็นเอและ อาร์เอ็นเอที่ไม่มีมีการเข้ารหัสอาจมีผลในการทำงานและงานได้เริ่มอธิบายถึงบทบาทของเส้นทางการหมุนเวียนของอาร์เอ็นเอและเส้นทางการรับรู้ความเสียหายเฉพาะในการล้างเซลล์ของอาร์เอ็นเอเหล่านี้ โดยงานวิจัยนี้ช่วยบอกปัจจัยที่มีผลต่อการสลายตัวของอาร์เอ็นเอ (Wurtmann และ Wolin, 2009)
2. Norbert Pardi และคณะ (2018) ได้ศึกษาว่าวัคซีนเอ็มอาร์เอ็นเอเป็นทางเลือกที่ดีที่สุดสำหรับวิธีการฉีดวัคซีนทั่วไป เนื่องจากมีศักยภาพสูงความสามารถในการพัฒนาอย่างรวดเร็วและศักยภาพในการใช้ต้นทุนต่ำ การผลิตและการบริหารที่ปลอดภัย แต่ถูกจำกัด โดยความไม่เสถียรและการส่งมอบเอ็มอาร์เอ็นเอที่ไม่มีประสิทธิภาพวัคซีนเอ็มอาร์เอ็นเอหลายตัวต่อต้านโรคติดเชื้อและมะเร็งหลายชนิดบทวิจารณ์นี้ให้ภาพรวมโดยละเอียดของวัคซีนเอ็มอาร์เอ็นเอและพิจารณาทิศทางและความท้าทายในอนาคตในการพัฒนาแพลตฟอร์มวัคซีนใช้ในการรักษาอย่างแพร่หลาย (Pardi, Hogan, Porter, และ Weissman, 2018)
3. Edison Ong และคณะ (2020) พัฒนาวัคซีนที่มีประสิทธิภาพและปลอดภัย 15 วัคซีนสำหรับโรคติดต่อที่เกิดจากโคโรนาไวรัส ได้ใช้เครื่องมือ Vaxign reverse Vaccinology 19 และเครื่องมือการเรียนรู้ของเครื่อง Vaxign-ML ที่พัฒนาขึ้นมาใหม่เพื่อทำนายผู้ป่วยที่ได้รับวัคซีน COVID-19 อัลกอริทึมการจำแนก ML ได้แก่ การถดถอยโลจิสติก, ซัพพอร์ตเวกเตอร์แมชชีน, วิธีการเพื่อนบ้านใกล้ที่สุด, แรนด้อมฟอเรส และ xgboost ประสิทธิภาพที่ดีที่สุดคือ xgboost (Ong, Wong, Huffman, และ He, 2020)
4. Brett N. Bowman และคณะ (2011) ได้ศึกษาหาการออกแบบวัคซีนหน่วยย่อยโดยการคาดการณ์อย่างรวดเร็วว่าโปรตีนในโปรตีนโอมของแบคทีเรียชนิดใดเป็นแอนติเจนที่ป้องกันได้โดยใช้วิธีเครื่องจักรเวกเตอร์สนับสนุน (Bowman และคนอื่น ๆ, 2011)

5. Mohammad N. Uddin และคณะพยายามหากลยุทธ์เพื่อเพิ่มความคงตัวของวัคซีนที่ใช้ mRNA ในอุณหภูมิที่ค่อนข้างสูงขึ้น เนื่องจากการใช้ทั่วโลกถูกจำกัดโดย ultracold ข้อกำหนดในการจัดเก็บ ดังนั้นประเทศที่ยากจนด้านทรัพยากรส่วนใหญ่ไม่มีที่เก็บห่วงโซ่ความเย็นเพื่อดำเนินการเก็บวัคซีนจำนวนมากดังนั้น การหากลยุทธ์เพื่อเพิ่มความคงตัวของวัคซีนที่ใช้เอ็มอาร์เอ็นเอ ในอุณหภูมิที่ค่อนข้างสูงขึ้นสามารถช่วยในการจัดการป้องกันโรคระบาดทั่วโลกในปัจจุบัน (Uddin MN, 2021)
6. Hannah K. Wayment-Steele และ คณะวัดความละเอียดของนิวคลีโอไทด์เดี่ยวบนโครงสร้าง RNA ที่หลากหลายของนิวคลีโอไทด์ 102-130 นิวคลีโอไทด์ โดยใช้ชุดข้อมูลของ Eterna การทดลองทั้งหมดเสร็จสิ้นภายในเวลาไม่ถึง 6 เดือน และ 41% ของการคาดการณ์ระดับนิวคลีโอไทด์จากแบบจำลองที่ชนะนั้นเกิดจากข้อผิดพลาดในการทดลองของการวัดค่าความจริงพื้นฐาน ยิ่งไปกว่านั้น แบบจำลองเหล่านี้ทำให้คาดการณ์ข้อมูลการสลายตัวแบบมุมฉากแบบสุ่มสี่สุ่มห้าบนโมเลกุล mRNA ที่ยาวกว่ามาก (นิวคลีโอไทด์ 504-1588) ด้วยความแม่นยำที่ดีขึ้นเมื่อเทียบกับแบบจำลองที่เคยแพร่ก่อนหน้านี้ (Wayment-Steele, 2022)
7. Amgad Muneer และ คณะได้ศึกษาคุณสมบัติการย่อยสลายของเบส RNA แต่ละตัวในโมเลกุล เพื่อตรวจสอบว่าการเรียนรู้เชิงลึกแบบผสมผสานสามารถทำนายการย่อยสลาย RNA จากลำดับ RNA ได้หรือไม่ โดยเสนอแบบจำลองโครงข่ายประสาท ได้แก่ GCN\_GRU และ GCN\_CNN ผลการทดลองพบว่า GCN\_GRU มีประสิทธิภาพดีกว่าโมเดล GCN\_CNN แบบจำลองได้รับคะแนนสูงสุดจากคะแนนการทดสอบ MCRMSE 0.22614 และ 0.34152 ตามลำดับ ในที่สุด (Amgad Muneer, 2022)
8. Feiran Cheng ศึกษาปัจจัยที่มีอิทธิพลต่อความคงตัว ตั้งแต่โครงสร้าง mRNA สารเพิ่มปริมาณระบบการนำส่งอนุภาคนาโนไขมัน (LNP) และกระบวนการผลิตเป็นปัจจัยหลักส่งผลกระทบต่อความคงตัวของวัคซีนเอ็มอาร์เอ็นเอ การปรับโครงสร้างเอ็มอาร์เอ็นเอให้เหมาะสมและคัดกรองสารเพิ่มปริมาณได้อย่างมีประสิทธิภาพปรับปรุงความเสถียรของวัคซีนเอ็มอาร์เอ็นเอ นอกจากนี้ การปรับปรุงกระบวนการผลิตยังสามารถเตรียมความพร้อมวัคซีนเอ็มอาร์เอ็นเอที่มีความเสถียรทางความร้อนพร้อมความปลอดภัยและสรุปปัจจัยสำคัญที่มีผลต่อความคงตัวของวัคซีนเอ็มอาร์เอ็นเอและเสนอแนวทางการวิจัยที่เป็นไปได้เพื่อปรับปรุงความเสถียรของวัคซีนเอ็มอาร์เอ็นเอ(Cheng F, 2023)
9. จากชุดข้อมูลของ OpenVaccine ที่มีฐานข้อมูล RNA พร้อมการวัดอัตราการย่อยสลายโดยนักวิจัยของ Stanford งานวิจัยนี้ Sheikh Asif Imran และคณะได้พัฒนาแบบจำลองโครงข่าย

ประสาทเพื่อช่วยให้นักวิจัยด้านชีวสารสนเทศสามารถระบุได้ว่าเอ็มอาร์เอ็นเอที่จะไม่เสถียรมีแนวโน้มที่จะเกิดขึ้นที่ตำแหน่งด้วยแบบจำลอง LSTM พบว่าค่าการทดสอบ MCRMSE อยู่ที่ 0.38796 (S. Asif Imran, 2020)

10. ปัญญาประดิษฐ์และการเรียนรู้โดยเครื่องจักร โดยเฉพาะอย่างยิ่งการเรียนรู้เชิงลึกนำไปสู่การปรับปรุงครั้งใหญ่ในสาขาวิทยาศาสตร์และวิศวกรรมศาสตร์หลายสาขา เนื่องจากความสามารถในการเรียนรู้คุณลักษณะต่างๆ อย่างลึกซึ้ง โดยเฉพาะการค้นพบวัคซีนสำหรับโรค (Arash Keshavarzi Arshadi, 2020)

จากงานวิจัยข้างต้นที่กล่าวมาจะเห็นได้ว่าการศึกษ้อัตราการสลายตัวของวัคซีนในหลายๆ งานเพื่อความเข้าใจในหลักการทำงานของวัคซีนเอ็มอาร์เอ็นเอ และพยายามที่ทำให้วัคซีนคงตัวนานที่สุด และในการศึกษานี้ ผู้วิจัยได้ทดลองแบบเปรียบเทียบแบบจำลองที่ไม่ซับซ้อนมากนักว่าจะสามารถทำงานได้ดีและให้ผลลัพธ์ใกล้เคียงกับงานวิจัยอื่นๆ ที่ใช้แบบจำลองซับซ้อน โดยเน้นที่การดูประสิทธิภาพของ Xg boost ,Random Forest,Cat boost และLightGBMโดยจะกล่าวถึงในส่วนถัดไป



## สรุปรงานวิจัยที่ศึกษาทั้งหมด

ลำดับที่	เรื่อง	ปีที่ตีพิมพ์	แบบจำลอง					อื่นๆ
			LR	RF	XgB	SVM	KNN	
1	RNA under attack: cellular handling of RNA damage	2009						ปัจจัยที่มีผลต่อการสลายตัวของอาร์เอ็นเอ
2	mRNA vaccines - a new era in vaccinology	2018						การพัฒนาวัคซีนเอ็มอาร์เอ็นเอ
3	COVID-19 coronavirus vaccine design using reverse vaccinology and machine learning	2020	/	/	/	/	/	
4	Improving reverse vaccinology with a machine learning approach	2011				/		
5	Challenges of Storage and Stability of mRNA-Based COVID-19 Vaccines	2021						หากกลยุทธ์เพื่อเพิ่มความคงตัวของวัคซีนที่ใช้เอ็มอาร์เอ็นเอ
6	Deep learning models for predicting RNA degradation via dual crowdsourcing.	2022						GNN
7	Vaccine-Deep: Prediction of COVID-19 mRNA vaccine degradation using deep learning	2022						GCN_GRU และ GCN_CNN
8	Research Advances on the Stability of mRNA Vaccines	2023						ศึกษาปัจจัยที่มีอิทธิพลต่อความคงตัวของโครงสร้าง mRNA
9	COVID-19 mRNA Vaccine Degradation Prediction using Regularized LSTM Model,	2022						LSTM
10	Artificial Intelligence for COVID-19 Drug Discovery and Vaccine Development	2020						วิเคราะห์การใช้ปัญญาประดิษฐ์ในการค้นพบวัคซีน

ภาพประกอบ 10 แสดงตารางสรุปรงานวิจัยที่ศึกษา

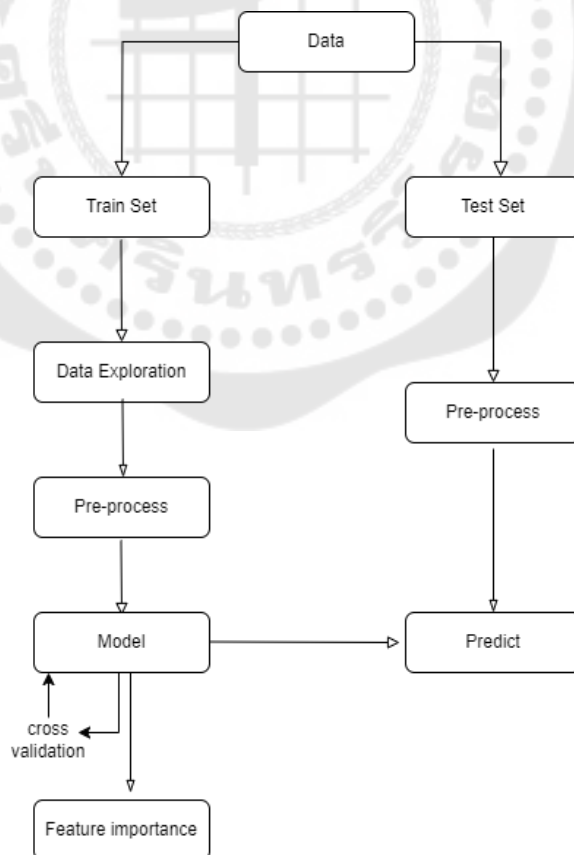
### บทที่ 3 การดำเนินการวิจัย

#### แผนขั้นตอนดำเนินการวิจัย

ในการวิจัยครั้งนี้ ผู้วิจัยดำเนินงานตามหัวข้อดังนี้

1. ทำ Literature Review ที่เกี่ยวข้อง
2. หาชุดข้อมูล(data set)เพื่อใช้ในการวิจัย
3. เตรียมข้อมูล
4. ศึกษาวิธีการสร้างและสร้างโมเดลในการทำนายข้อมูล
5. ปรับจูนโมเดลและวัดประสิทธิภาพของโมเดล
6. ประเมินและสรุปงานวิจัย

#### แผนผังดำเนินงานระบบ



ภาพประกอบ 11 แผนผังการทำงานของระบบ

จากภาพประกอบ 11 ได้อธิบายถึงกระบวนการสร้างแบบจำลองและการวัดผลลัพธ์จากการวิเคราะห์ข้อมูล โดยเริ่มจากการนำข้อมูล (Data) จาก Kaggle เข้ามาทำการแบ่งชุดข้อมูลออกเป็นชุดเรียนรู้ (train) และชุดทดสอบ (test) นำข้อมูลชุดเรียนรู้มาทำการสำรวจเบื้องต้น ก่อนนำเข้าสู่การเตรียมข้อมูลให้พร้อม (Pre-processing) เพื่อทำไปใช้ในการสร้างแบบจำลอง (Model) ทำ Cross validation เพื่อหาค่าที่ทำให้แบบจำลองดีที่สุดแล้วทำการรันแบบจำลองที่หากผลลัพธ์ของแบบจำลองยังไม่เป็นที่น่าพอใจ จะกลับไปปรับการปรับพารามิเตอร์ เพื่อให้ได้มาซึ่งแบบจำลองที่มีประสิทธิภาพดีที่สุด และนำแบบจำลองไปใช้กับชุดข้อมูลทดสอบที่ทำการเตรียมข้อมูลแล้ว

### ชุดข้อมูล(data set)เพื่อใช้ในงานวิจัย

ในการวิจัยครั้งนี้ผู้วิจัยได้ใช้ข้อมูลโครงการ OpenVaccine: COVID-19 mRNA Vaccine Degradation Prediction จาก Kaggle โดยความร่วมมือของมหาวิทยาลัยแสตนฟอร์ด และชุมชนอีเทอน่า (Stanford, 2020) โดยจากข้อมูลจะแบ่งเป็น

train.json – ข้อมูลชุดเรียนรู้มีค่าของการเก็บข้อมูลภาคสนาม (Ground Truth) ประกอบด้วย 5 เมตริกของการเกิดปฏิกิริยาย่อยสลาย ดังที่ระบุไว้ในตารางที่ 1 อย่างไรก็ตาม การศึกษาครั้งนี้มุ่งเน้นไปที่การประเมิน 3 เมตริกแรก ได้แก่ ปฏิกิริยา deg\_Mg\_ph10 และ deg\_Mg\_50C

test.json - ชุดทดสอบโดยไม่มีคอลลัมน์ใด ๆ ที่เกี่ยวข้องกับค่าของการเก็บข้อมูลภาคสนาม

ตาราง 1 รายละเอียดคอลลัมน์ของชุดฝึกฝน

คอลลัมน์	รายละเอียดคอลลัมน์
Reactivity	ค่าความสามารถในการเกิดปฏิกิริยาทั่วไป
deg_pH10	ความเป็นไปได้ของการสลายตัวหลังการบ่มที่ pH สูง (pH10) โดยไม่มีแมกนีเซียม
deg_Mg_pH10	ความเป็นไปได้ของการสลายตัวหลังจากบ่มด้วยแมกนีเซียมใน pH สูง (pH 10)
deg_50C	ความเป็นไปได้ของการสลายตัวหลังจากการบ่มโดยไม่มีแมกนีเซียมที่อุณหภูมิสูง (50 องศาเซลเซียส)
deg_Mg_50C	ความเป็นไปได้ของการสลายตัวหลังจากบ่มด้วยแมกนีเซียมที่อุณหภูมิสูง (50 องศาเซลเซียส)
Id	รหัสสำหรับแต่ละตัวอย่าง
seq_scored	ค่าจำนวนเต็มแสดงถึงจำนวนตำแหน่งที่ใช้ในการให้คะแนนด้วยค่าที่ทำนาย
seq_length	ความยาวของลำดับ
Sequence	ลำดับ RNA การรวมกันของ A, G, U และ C สำหรับแต่ละตัวอย่าง

## กระบวนการสำรวจข้อมูลเบื้องต้น (Exploratory Data Analysis)

มีขั้นตอนดังต่อไปนี้

1. ดูค่า Summary แต่ละคอลัมน์ และ ตรวจสอบรายละเอียดของข้อมูล
2. ตรวจสอบ Signal to noise และ SN filter
3. ตรวจสอบค่าเฉลี่ยการย่อยสลายของเป้าหมาย

### ดูชนิดของข้อมูลและ เช็ค Summary ของแต่ละคอลัมน์

ตรวจสอบว่าข้อมูลมีส่วนที่ว่างหรือสูญหายหรือไม่ โดยข้อมูลชุดนี้ไม่มีการสูญหายของข้อมูล และดูชนิดของข้อมูลว่าประกอบด้วยอะไรบ้างในแต่ละคอลัมน์

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2400 entries, 0 to 2399
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  ---                ---
0   index                 2400 non-null   int64
1   id                    2400 non-null   object
2   sequence              2400 non-null   object
3   structure             2400 non-null   object
4   predicted_loop_type   2400 non-null   object
5   signal_to_noise       2400 non-null   float64
6   SN_filter             2400 non-null   int64
7   seq_length            2400 non-null   int64
8   seq_scored            2400 non-null   int64
9   reactivity_error     2400 non-null   object
10  deg_error_Mg_pH10    2400 non-null   object
11  deg_error_pH10       2400 non-null   object
12  deg_error_Mg_50C    2400 non-null   object
13  deg_error_50C        2400 non-null   object
14  reactivity            2400 non-null   object
15  deg_Mg_pH10          2400 non-null   object
16  deg_pH10              2400 non-null   object
17  deg_Mg_50C           2400 non-null   object
18  deg_50C               2400 non-null   object
dtypes: float64(1), int64(4), object(14)
memory usage: 356.4+ KB
```

ภาพประกอบ 12 ดูชนิดข้อมูล

เช็คค่าสถิติพื้นฐานของข้อมูลว่าระบุมาเพียงพอสำหรับการทำแบบจำลองไหม

	index	signal_to_noise	SN_filter	seq_length	seq_scored
count	2400.000000	2400.000000	2400.000000	2400.0	2400.0
mean	1199.500000	4.530456	0.662083	107.0	68.0
std	692.964646	2.835142	0.473099	0.0	0.0
min	0.000000	-0.103000	0.000000	107.0	68.0
25%	599.750000	2.391000	0.000000	107.0	68.0
50%	1199.500000	4.442500	1.000000	107.0	68.0
75%	1799.250000	6.294250	1.000000	107.0	68.0
max	2399.000000	17.194000	1.000000	107.0	68.0

ภาพประกอบ 13 ดู Summary ของข้อมูล

ทำการตรวจสอบรายละเอียดของข้อมูลในตารางเพื่อดูเอาท์พุตที่ต้องใช้งานโดยยกตัวอย่างมา 1 ตัวอย่าง ตามรูปที่ 14

```
One_sequence : GGAAAAGCUCUAUAACAGGAGACUAGGACUACGUUUUCUAGGUAACUGGAUAACCCAUACCAGCAGUUAGAGUUCGCUCUAACAAAAGAAACAACAACAAC
One_structure : .....((((((.....)))))).((.....((.....((((.....)))))).((.....))....((((.....))))).((.....))....((((.....))))).((.....))....
One_predicted_loop_type : EEEEESSSSSH####HSSSSBSSSSIIIISSIISSSSSH####HSSSSSISSIIIISSXXXSSSSSH####HSSSSSEEEEEEEEEEEEEEEEEEE
One_reactivity : [0.3297, 1.5693000000000001, 1.1227, 0.8686, 0.7217, 0.4384, 0.256, 0.3364000000000003, 0.2168000000000002, 0.3583,
```

ภาพประกอบ 14 ตัวอย่างของข้อมูลตามคอลัมน์

ใช้คำสั่ง Count และ len ในการนับจำนวนนิวคลีโอไทด์แรก 68 (ตัวอักษรในลำดับเอ็มอาร์เอ็นเอ 107 ตัวอักษร)

```
len(One_reactivity)
68

Counter(One_sequence)
Counter({'A': 45, 'C': 23, 'G': 19, 'U': 20})

sum_one_sequence = sum(Counter(One_sequence).values())
sum_one_sequence
107

Counter(One_structure)
Counter({'(': 23, ')': 23, '.' : 61})

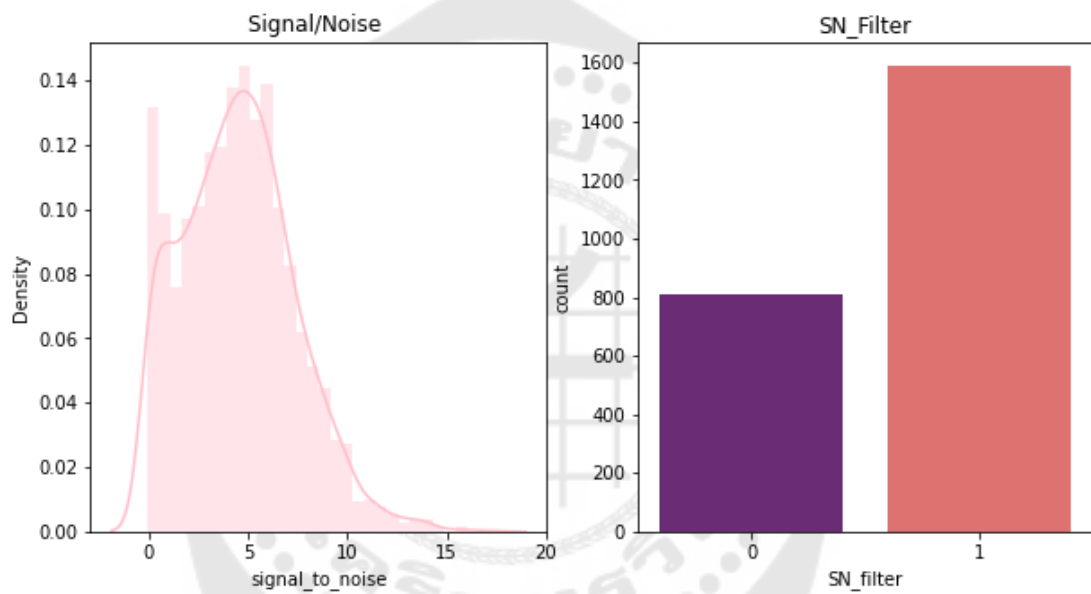
sum_one_structure = sum(Counter(One_structure).values())
sum_one_structure
107
```

ภาพประกอบ 15 การใช้คำสั่ง count และ len

### ตรวจสอบ Signal to noise และ SN filter

Signal to noise มาจากอัตราส่วนสัญญาณต่อสัญญาณรบกวน (หรือ SNR, S/N) และแสดงควมกว้างของสัญญาณตามสัญญาณรบกวนเป็นการวัดมาตรฐานที่ใช้โดยทั่วไปในสัญญาณความถี่สูงเนื่องจากชุดฝึกฝนถูกสร้างขึ้นจากผลการทดลองในห้องปฏิบัติการ เราต้องการอธิบายว่าค่าต่างๆ ได้ผ่านขั้นตอนการควบคุมคุณภาพแล้ว

ในส่วนของ SN\_filter มีไว้เพื่อเป็นการตรวจสอบว่าขั้นตอนการควบคุมคุณภาพทำงานได้ตามปกติ

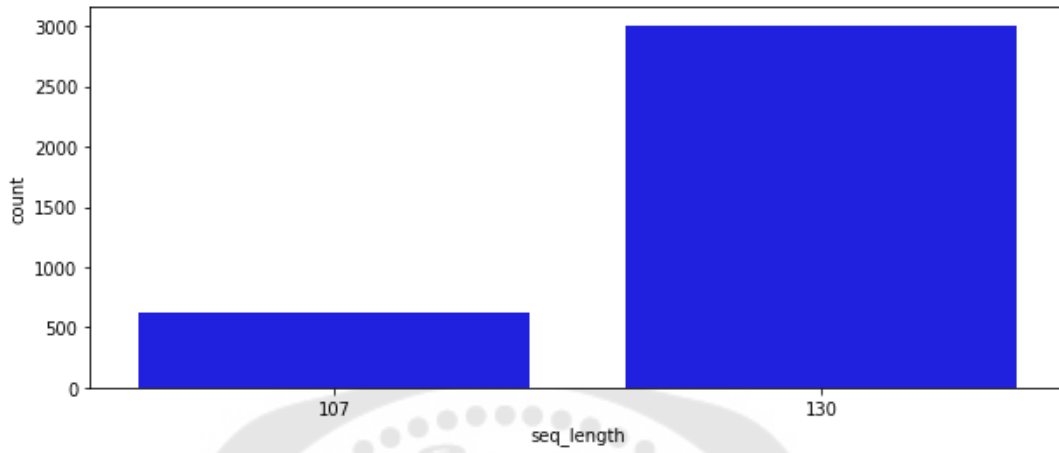


ภาพประกอบ 16 จำนวน SN\_filter และ signal to noise

ในงานวิจัยนี้ ค่าต่ำสุดของเงื่อนไขทั้ง 5 ตามตารางที่ 1 ต้องมากกว่า -0.5 Signal/Noise เฉลี่ยทั้ง 5 เงื่อนไขต้องมากกว่า 1.0 หากเป็นไปตามเงื่อนไข SN\_filter จะเท่ากับ 1

จากรูปที่ 16 ยังพบ signal\_to\_noise มีค่าน้อยกว่า - 0.5 และ SN\_filter เท่ากับ 0 ซึ่งอาจจะเป็น outlier ที่เกิดขึ้น โดย signal\_to\_noise มากกว่า 1 จำนวน 2096 ตัวอย่าง ,SN\_filter เท่ากับ 1 จำนวน 1589 ตัวอย่าง

ข้อมูลฝึกฝนประกอบด้วยความยาว sequence 107 เท่านั้น ข้อมูลการทดสอบส่วนใหญ่  
มีความยาว sequence 130



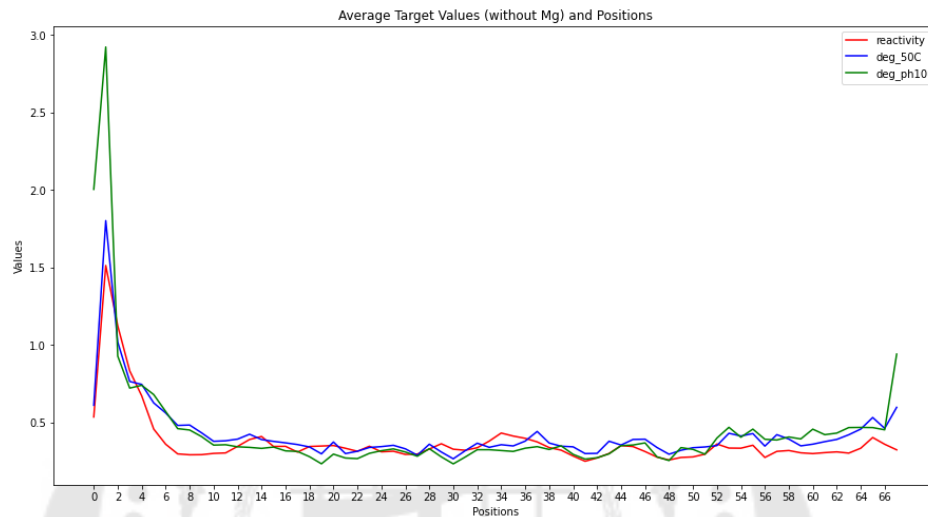
ภาพประกอบ 17 จำนวนความยาว sequence ข้อมูลทดสอบ





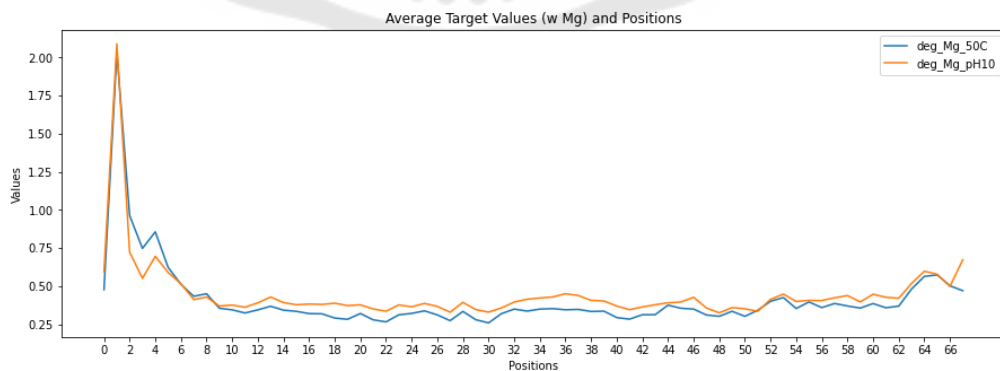
### ตรวจสอบค่าเฉลี่ยการย่อยสลายของเป้าหมาย

ทำการตรวจสอบค่าเฉลี่ยการย่อยสลายของเป้าหมายโดยเราจะเห็นได้ว่าตำแหน่งการย่อยสลายและการเกิดปฏิกิริยาจะเกิดสูงสุดที่ตำแหน่งต้นสาย sequence โดยเฉพาะที่ pH10 ตามรูปที่ 18 จากภาพเราจะเห็นความสัมพันธ์ระหว่างอุณหภูมิสูง 50C และ pH10



ภาพประกอบ 18 คำนวณค่าเฉลี่ยสำหรับคอลัมน์เป้าหมายแบบไม่ใช้แมกนีเซียม

และในสภาวะการใช้แมกนีเซียมร่วมโดยส่วนใหญ่แล้ว deg\_Mg\_pH10 จะมีการลดลยตัวมากกว่าที่ deg\_Mg\_50C แต่ทั้ง 2 เงื่อนไขมีความสัมพันธ์กันอย่างเห็นได้ชัดจากรูปที่ 19



ภาพประกอบ 19 คำนวณค่าเฉลี่ยสำหรับคอลัมน์เป้าหมายแบบใช้แมกนีเซียม

## ศึกษาวิธีการสร้างและสร้างโมเดลในการทำนายข้อมูล

นำข้อมูลที่ได้จากกระบวนการเตรียมข้อมูลมาใช้ทดสอบโมเดลโดยงานวิจัยนี้ใช้โมเดลเปรียบเทียบกันทั้งหมด 4 โมเดล คือ Xgboost ใช้เป็น Baseline ,Random Forest, Catboost และ LightGBM ในการตั้งค่าพารามิเตอร์ใช้เทคนิค 5-Fold cross- validation และทำการปรับโมเดลเพื่อหาค่าพารามิเตอร์ให้เหมาะสมด้วยการใช้ GridSearchCV

### ทดสอบการทำแบบจำลองด้วยเทคนิค XGboost (baseline)

เป็นแบบจำลองที่นำเอา Decision Tree มาต่อกันโดยที่แต่ละ decision tree จะเรียนรู้จาก error ของ tree ก่อนหน้า เมื่อมีการเรียนรู้ของ tree ต่อเนื่องกันจนมีความลึกมากพอ และโมเดลจะหยุดเรียนรู้เมื่อไม่เหลือรูปแบบของ error จาก tree ก่อนหน้าให้เรียนรู้แล้ว

```
xgb = XGBRegressor(
    subsample=0.8,
    colsample_bytree=0.75,
    reg_lambda=2,
    reg_alpha=1,
    random_state=28
)

params = {
    'estimator__n_estimators': [ 800, 900, 1000],
    'estimator__learning_rate': [0.1, 0.25, 0.30],
    'estimator__max_depth': [3, 4, 5]
}
```

ภาพประกอบ 20 การกำหนดค่าพารามิเตอร์เทคนิค XGboost

```
gs.best_params_
```

```
{'estimator__learning_rate': 0.1,
 'estimator__max_depth': 5,
 'estimator__n_estimators': 800}
```

ภาพประกอบ 21 พารามิเตอร์ที่ได้จาก GridSearchCV ในเทคนิค XGboost

### ทดสอบการทำแบบจำลองด้วยเทคนิค Random Forest

เป็นแบบจำลองที่นำ Decision Tree หลายๆ tree มา Train ร่วมกัน (ตั้งแต่ 10 ต้น ถึงมากกว่า 1000 ต้น) โดยที่แต่ละ tree จะได้รับ feature และ data เป็น subset ของ feature และ data ทั้งหมด แบบสุ่มตอนคาดการณ์ แต่ละ Decision Tree ทำการคาดการณ์ของตัวเอง และเลือกผล final prediction

```

rf = RandomForestRegressor(random_state=42,verbose=1)
params = {
    'estimator__n_estimators': [ 100, 200, 300],
    'estimator__max_depth': [8, 9, 10],
}
reg = MultiOutputRegressor(rf)
# Perform 5 fold cross validation on set 1 features
rfr = GridSearchCV(estimator=reg, param_grid = params, cv = 5, n_jobs=-1, scoring=custom_scorer)

```

ภาพประกอบ 22 การพารามิเตอร์ของเทคนิค Random Forest

### ทดสอบการทำแบบจำลองด้วยเทคนิค Catboost

เป็นแบบจำลองการจัดการสำหรับคุณสมบัติตามหมวดหมู่ได้รับความนิยมเมื่อเทียบกับอัลกอริทึม Gradient boosting อื่น ๆ เนื่องจากคุณสมบัติการใช้ Oblivious Trees หรือ Symmetric Trees เพื่อให้โมเดลทำงานได้เร็วขึ้น

```

[ ] cbr = CatBoostRegressor(random_state=28)
    params = {
        'estimator__n_estimators': [ 800, 900, 1000],
        'estimator__learning_rate': [0.1, 0.25, 0.30],
        'estimator__max_depth': [3, 4, 5]
    }
    reg = MultiOutputRegressor(cbr)
    # Perform 5 fold cross validation on set 1 features
    cb = GridSearchCV(estimator=reg, param_grid = params, cv = 5, n_jobs=-1, scoring=custom_scorer)

```

ภาพประกอบ 23 การกำหนดพารามิเตอร์ของเทคนิค Catboost

### ทดสอบการทำแบบจำลองด้วยเทคนิค Light GBM

เป็นเฟรมเวิร์ก gradient boosting framework ที่รวดเร็ว, ประสิทธิภาพสูงโดยใช้อัลกอริทึมแผนผังและงานการเรียนรู้ของเครื่องอื่น ๆ อีกมากมาย ใช้วิธีการแยกความลึกของ Trees

```

FOLD_N = 5
gkf = GroupKFold(n_splits=FOLD_N)

```

```

params = {'objective': 'regression',
          'boosting': 'gbdt',
          'metric': 'rmse',
          'learning_rate': 0.05,
          'seed' : SEEDS}

```

ภาพประกอบ 24 การกำหนดพารามิเตอร์ของเทคนิค Light GBM

## ปรับจูนโมเดลและวัดประสิทธิภาพของแบบจำลอง

ในการศึกษานี้ ผู้วิจัยเปรียบเทียบแบบจำลองทั้ง 4 แบบ แล้วนำโมเดลที่ดีที่สุดในงานวิจัยไปหา feature importance เพื่อหาตำแหน่งที่มีการย่อยสลาย

สำหรับงานวิจัยนี้มีการปรับค่าพารามิเตอร์หลักในการสร้างแบบจำลองทำนายดังนี้

learning_rate	Hyperparameter ควบคุมปรับ Weight ใน Step ของการฝึกฝน
n_estimators	weight function ที่ใช้ในการฝึกฝนและทำนาย
subsample	สัดส่วนของตัวอย่างข้อมูลสำหรับเรียนรู้
colsample_bytree	สัดส่วนในการสุ่มคอลัมน์
reg_lambda	L1 regularization เป็นตัวคูณหน้า weight
reg_alpha	L2 regularization เป็นตัวคูณหน้า weight
max_depth	เมตริกที่ใช้วัดระยะใน tree

และมีการปรับจูนพารามิเตอร์รายละเอียดตามตารางที่ 2

ตาราง 2 พารามิเตอร์ที่เหมาะสมของแต่ละโมเดลในที่ใช้ในการเรียนรู้

แบบจำลอง	พารามิเตอร์
XGBoost	learning_rate: 0.1, n_estimators: 800 subsample=0.8 colsample_bytree=0.75 reg_lambda=2 reg_alpha=1
Random Forest	max_depth: 11 n_estimators: 800
Catboost	learning_rate: 0.1 max_depth: 8 n_estimators: 300
LigthGBM	learning_rate: 0.05

การวิเคราะห์และประเมินผลลัพธ์เพื่อเปรียบเทียบโมเดลจะใช้วิธีการเปรียบเทียบด้วยค่า mean columnwise root mean squared

### MCRMSE (mean columnwise root mean squared error)

โดยปกติเราสามารถคำนวณ root mean squared error (RMSE) เพื่อให้ได้เมตริกการประเมินตัวเลขเดียวสำหรับการคาดคะเนของเรา แต่ถ้าเราคาดการณ์หลายค่าพร้อมกัน ในกรณีของการแข่งขัน OpenVaccine เราจำเป็นต้องทำนายอัตราการย่อยสลายภายใต้เงื่อนไขหลายประการ เราจะได้ค่า RMSE ที่แตกต่างกันหลายค่า หนึ่งค่าสำหรับแต่ละคอลัมน์

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

ดังนั้น MCRMSE เป็นเพียงค่าเฉลี่ยของค่า RMSE ทั้งหมดสำหรับแต่ละคอลัมน์ของเรา ดังนั้นเราจึงยังคงสามารถใช้เมตริกการประเมินด้วยตัวเลขเดียวได้ แม้ในกรณีที่มีหลายผลลัพธ์

$$MCRMSE = \frac{1}{N_t} \sum_{j=1}^{N_t} \sqrt{\frac{1}{n} \sum_{i=1}^n (y_{ij} - \hat{y}_{ij})^2} \quad (2)$$

ที่มา:(Stanford, 2020)

จากสมการตัวแปร  $N_t$  คือจำนวนตัวแปรที่ใช้คาดการณ์

$n$  คือ จำนวนของชุดทดสอบ

ตัวแปร  $y$  และตัวแปร  $\hat{y}$  คือค่าจริงและค่าที่ทำนายตามลำดับค่าเฉลี่ย

## บทที่ 4

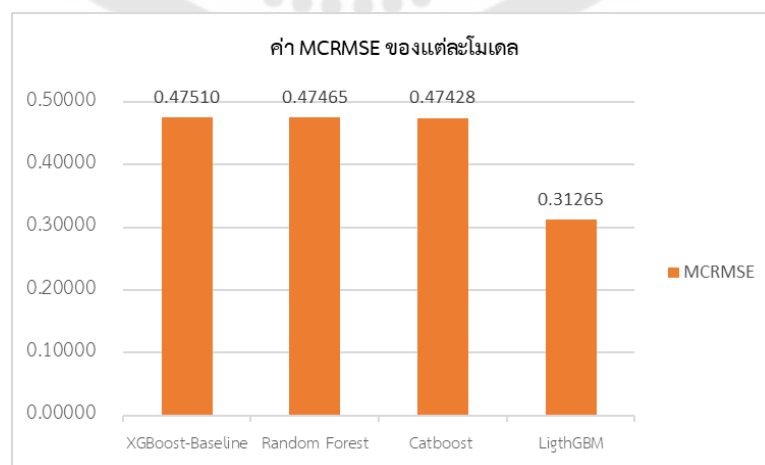
### ผลการดำเนินงานวิจัย

ในการศึกษาการทำนายค่าความเสียหายของวัคซีนเอ็มอาร์เอ็นเอ ซึ่งใช้ข้อมูลเกี่ยวกับการอัตราการย่อยสลายและตำแหน่งการย่อยสลายของ RNA ของข้อมูลที่ได้จากชุมชน Kaggle ผู้วิจัยได้วิจัยโดยการศึกษารุ่นต่อนต่างๆ ตามจุดประสงค์ในการวิจัยคือ

1. ศึกษาการสร้างแบบจำลองการทำนายค่าความเสียหายของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด 19 ด้วยวิธีการเรียนรู้ของเครื่องจักร
2. เปรียบเทียบการทำนายความเสียหายของโมเดลการเรียนรู้เครื่องจักร

ผลการเปรียบเทียบค่า MCRMSE จาก baseline ได้ค่า MCRMSE 0.4751 นั้น เมื่อนำเทคนิค Random Forest และ Catboost มาใช้งานและหาพารามิเตอร์ที่เหมาะสมพบว่าค่า MCRMSE ลดลงเล็กน้อยคือ 0.47465 และ 0.47428 แต่เมื่อใช้ LigthGBM พบว่าค่าลดลงจาก baseline -คือ 0.31265 เมื่อค่า MCRMSE น้อยแสดงให้เห็นว่าเทคนิค LigthGBM มีความคลาดเคลื่อนที่ค่อนข้างดีกว่าเทคนิคอื่นๆ ซึ่งแสดงตามภาพประกอบที่ 25 ดังนั้นจึงให้ผลลัพธ์ที่แม่นยำกว่า

ทางผู้วิจัยจึงเลือกเทคนิคนี้ในการนำไปหา Feature Importance ต่อไป



ภาพประกอบ 25 ค่า MCRMSE ของแต่ละแบบจำลอง

### การหา Feature importance

Feature importance เป็นฟังก์ชันที่ถูกนำมาใช้เพื่อบอกว่า feature ใดที่มีผลต่อแบบจำลองโดยนำมาเทียบกับ feature อื่นๆ ในการทดลองนี้ผู้วิจัยนำแบบจำลอง LightGBM ที่ได้มาหา Feature importance โดยกำหนด parameter ตามภาพประกอบ 26 ดังนี้

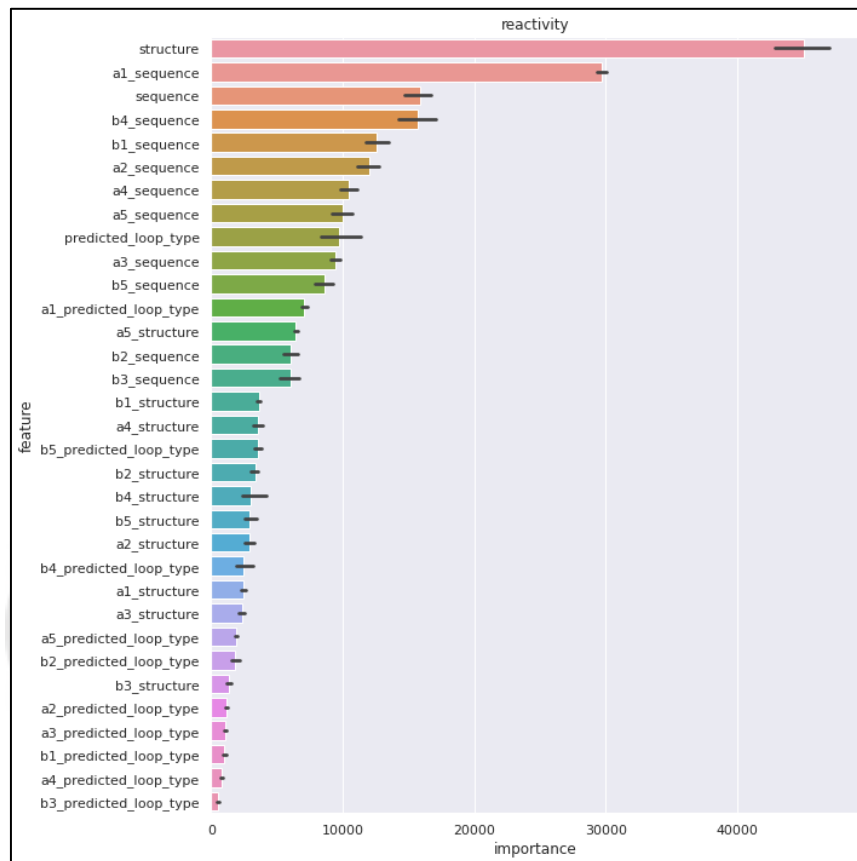
```
def feature_importances_(self):
    if self.model_type == 'lgb':
        return self.model.feature_importance(importance_type='gain')
```

ภาพประกอบ 26 parameter ของ feature importance

ปัจจัย (importance) ที่ส่งผลของประสิทธิภาพในการทำนายอัตราการย่อยสลายของ sequence จะมีทั้งหมด 5 ชนิด

1. การเกิด reactivity
2. การย่อยสลายโดยการใช้แมกนีเซียมที่ pH10
3. การย่อยสลายโดย pH10 อย่างเดียว
4. การย่อยสลายโดยการใช้แมกนีเซียมอุณหภูมิ 50 องศาเซลเซียส
5. การย่อยสลายที่ อุณหภูมิ 50 องศาเซลเซียส อย่างเดียว

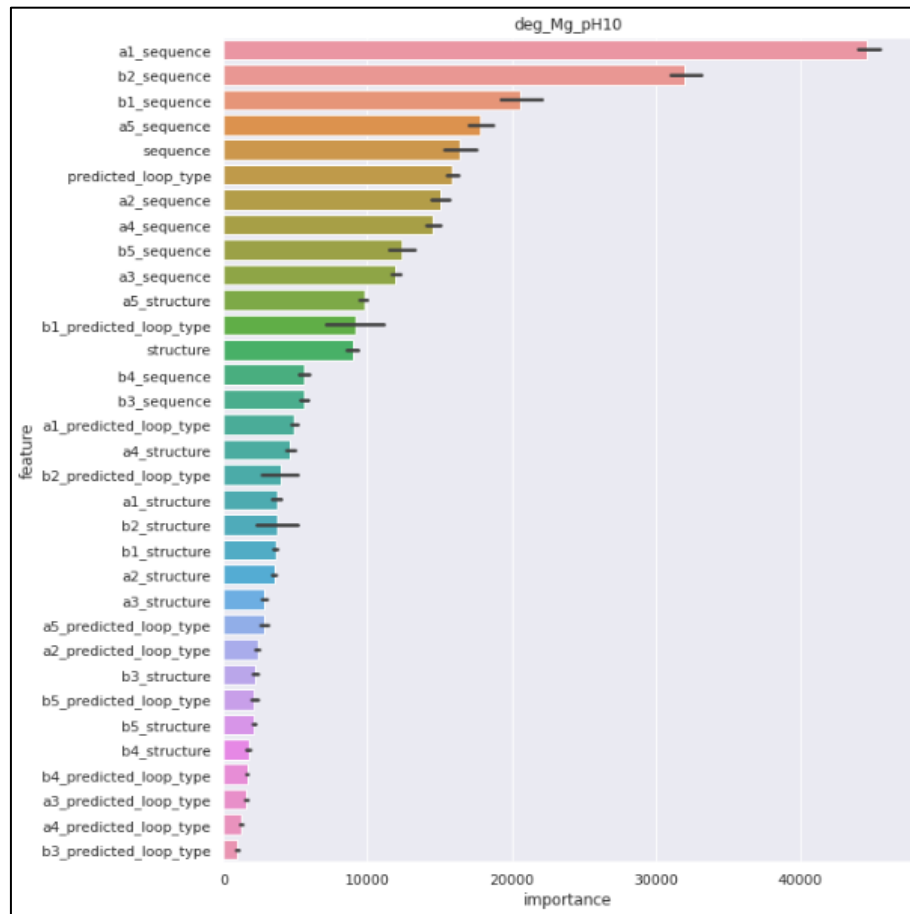
Feature ส่งผลของประสิทธิภาพในการทำนายอัตราการย่อยสลายของ sequence โดยการเกิด reactivity เป็นส่วนที่เรียกว่า structure (ส่วนที่ไว้ใช้ในการจับคู่สายของ เอ็มอาร์เอ็นเอ) มีค่า importance สูงที่สุด รองลงมาคือส่วนของ a1\_sequence จากรูปที่ 27



ภาพประกอบ 27 Feature importance ของค่า reactivity

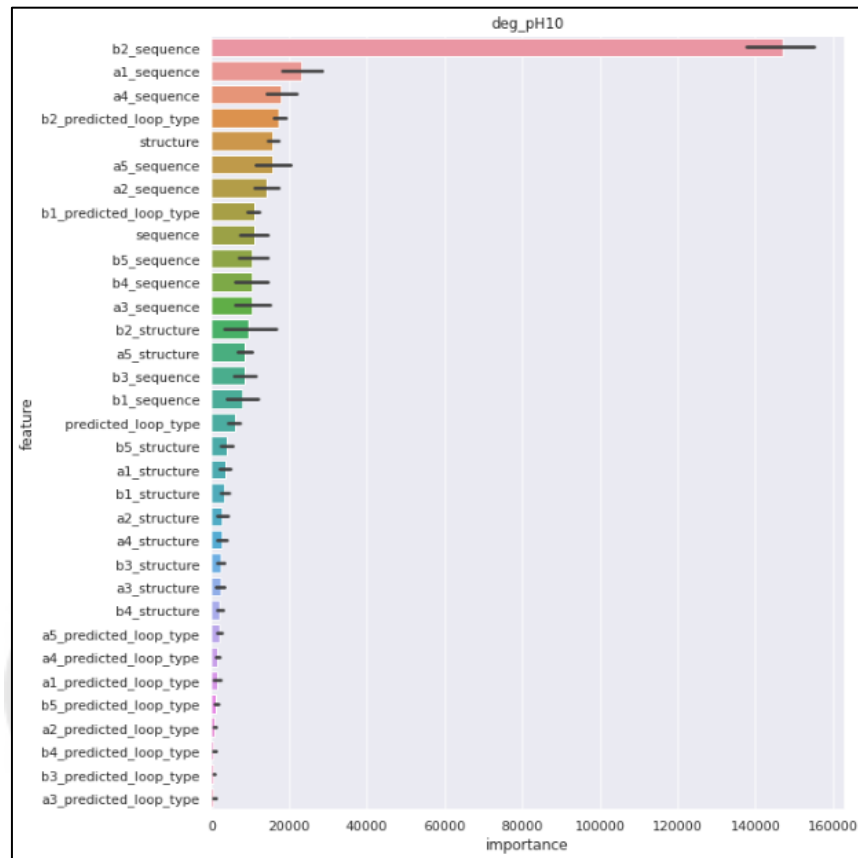


Feature ที่ส่งผลต่อประสิทธิภาพในการทำนายอัตราการย่อยสลายโดยการใส่แมกนีเซียมที่ pH10 เป็นส่วนที่เรียกว่า a1 sequence มีค่า importance สูงที่สุด เช่นกัน ตามรูปที่ 28



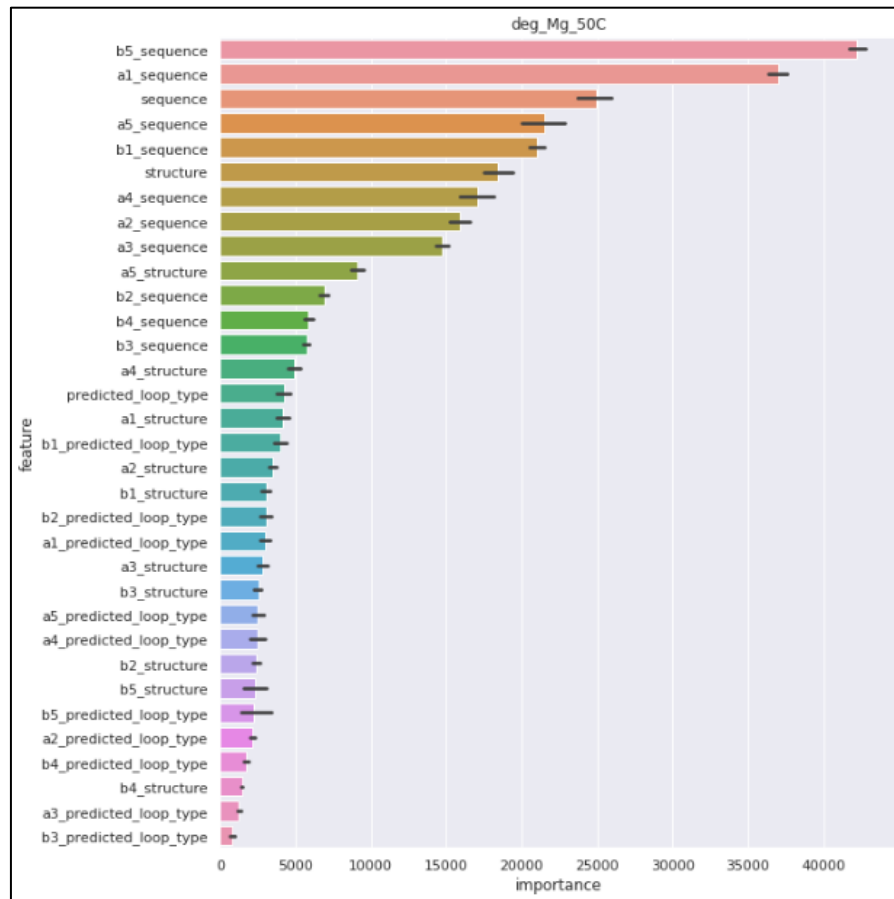
ภาพประกอบ 28 Feature importance ของค่า deg\_Mg\_pH10

ในทางกลับกัน Feature ที่ส่งผลของประสิทธิภาพในการทำนายอัตราการย่อยสลายโดย pH10 เพียงอย่างเดียว เป็นส่วนที่เรียกว่า b2\_sequence ซึ่งมีผลเยอะกว่า a1 sequence ตามรูปที่ 29



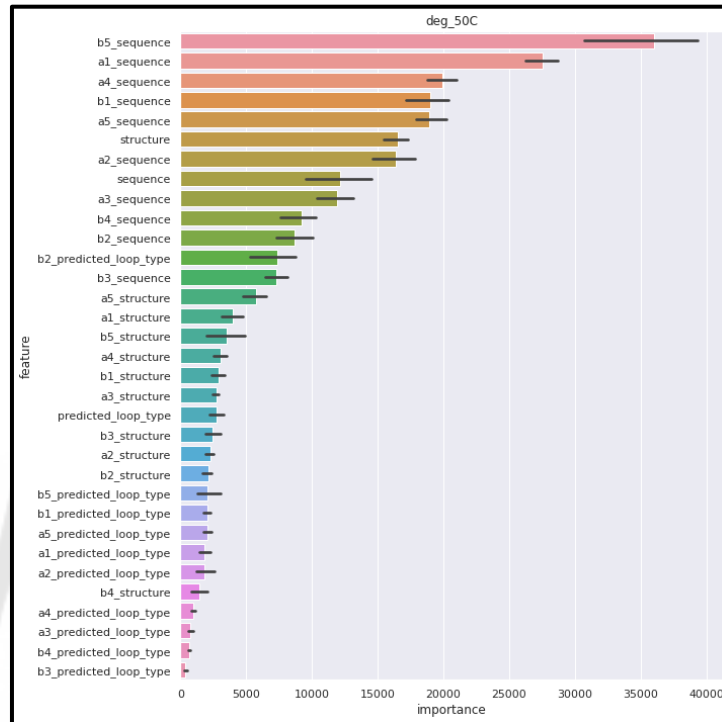
ภาพประกอบ 29 Feature importance ของค่า deg\_pH10

Feature ที่ส่งผลของประสิทธิภาพในการทำนายอัตราการย่อยสลายโดยใช้แมงนีเซียม  
 อุณหภูมิ 50 องศาเซลเซียส คือ b5\_sequence รองลงมาคือส่วนของ a1\_sequence ตามรูปที่ 30



ภาพประกอบ 30 Feature importance ของค่า deg\_Mg\_50C

Feature ที่ส่งผลของประสิทธิภาพในการทำนายอัตราการย่อยสลายที่อุณหภูมิ 50 องศาเซลเซียส  
 อย่างเดียว คือ b5\_sequence รองลงมาคือส่วนของ a1\_sequence อีกเช่นกัน ตามรูปที่ 31



ภาพประกอบ 31 Feature importance ของค่า deg\_50C

## บทที่ 5

### สรุปผลการวิจัย อภิปรายผลและข้อเสนอแนะ

ในการวิจัยเพื่อศึกษา เพื่อศึกษาการสร้างแบบจำลองการทำนายค่าความเสถียรของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด 19 ด้วยวิธีการเรียนรู้ของเครื่องจักรทั้งหมด 4 แบบจำลอง คือ Xgboost ,Random Forest, Catboost และ LightGBM ทำนายค่าความเสถียรของวัคซีนเอ็มอาร์เอ็นเอสำหรับโรคโควิด 19 ซึ่งเป็นอัตราการย่อยสลายในสภาวะต่างๆ โดยใช้หลักการของ MCRMSE สามารถแบ่งหัวข้อในการสรุปผลได้ดังต่อไปนี้

1. สรุปผลการวิจัย และ อภิปรายผลการวิจัย
2. ข้อเสนอแนะ

#### 5.1 สรุปผลการวิจัย และ อภิปรายผลการวิจัย

วัคซีนชนิดเอ็มอาร์เอ็นเอเป็นหนึ่งในเทคโนโลยีที่ถูกนำมาใช้ในการป้องกันการแพร่ระบาดของโรคโควิด 19 โดยวัคซีนชนิดนี้มีประสิทธิภาพในการป้องกันโรคได้สูงถึง 95% งานวิจัยนี้ นำเสนอการเรียนรู้ของเครื่องจักรทั้งหมด 4 แบบ คือ XGboost, Random Forest, Catboost และ Light GBM ในการทำนายความเสถียรของการย่อยสลายของลำดับอาร์เอ็นเอ โดยจากผลการทดลองพบว่า โมเดลการเรียนรู้ของเครื่องจักรชนิด Light GBM เป็นตัวทำนายที่ดีที่สุดโดยมีค่า Mean Column-wise Root Mean Squared (MCRMSE) เท่ากับ 0.31265 เมื่อเทียบกับ best score ของ challenge อยู่ที่ 0.23162 จากการค้นคว้าเพิ่มเติมเกี่ยวกับเทคนิคของ LightGBM พบว่าถูกใช้ในงานเกี่ยวกับการ identify tissue-specific genes โดยอ้างอิงจากงานของ Zijie Wang และคณะ (Wang, 2023) ใช้ parameter 'gain' ในการค้นหาชิ้นเฉพาะของพืช และเป็นเทคนิคที่เหมาะสมกับขนาดของข้อมูลที่มีความเยอะประมาณหนึ่ง อีกทั้งตัวของ LightGBM ยังมีอัลกอริทึม Leaf wise ซึ่งสามารถลดการสูญเสียได้มากกว่าอัลกอริทึม Level-wise ดังนั้นจึงน่าจะเป็นเหตุผลที่ LightGBM ให้ผลลัพธ์ที่ค่อนข้างดีกว่าเทคนิคอื่นๆ และเมื่อนำมาหาส่วนสำคัญที่มีผลกับประสิทธิภาพของโมเดลจะได้ตามตารางที่ 3

ตาราง 3 Feature Importance ที่สำคัญที่สุดของแต่ละเป้าหมาย

Target	Feature importance
Reactivity	Structure
deg_Mg_pH10	a1_sequence
deg_Mg_50C	b5_sequence
deg_pH10	b2_sequence
deg_50C	b5_sequence

การทำ Feature Importance แต่ละปัจจัยหรือเป้าหมายมีตำแหน่งของเอ็มอาร์เอ็นเอ ในการย่อยสลายที่แตกต่างกันไป จากผลที่ได้จะเห็นได้ว่าที่ ปัจจัยอุณหภูมิ 50 องศาเซลเซียสและ ปัจจัยอุณหภูมิ 50 องศาเซลเซียสในแมกนีเซียม sequence ที่มีผลต่อประสิทธิภาพของโมเดลเป็น sequence เดียวกันคือ b5\_sequence ซึ่งเป็นสิ่งที่น่าสนใจแสดงให้เห็นว่าถ้าในห้องปฏิบัติการจริงที่อุณหภูมิสูงอาจจะมีผลต่อการย่อยสลายของ sequence ตำแหน่งดังกล่าวเป็นพิเศษ และอีก ส่วนที่มองข้ามไม่ได้คือ a1\_sequence ซึ่งเป็นส่วนที่มีผลอยู่ในอันดับต้นๆ ของทุกปัจจัย การใช้วิธีการเรียนรู้ของเครื่องจักรในงานวิจัยนี้มีค่า MCRMSE อยู่ประมาณ 0.31-0.48 การใช้เทคนิค ดังกล่าวมีความสามารถในการเพิ่มความเร็วและประสิทธิภาพของการค้นพบวัคซีนเอ็มอาร์เอ็นเอ ที่เสถียรและมีผลต่อการวิจัยในสาขาอื่น ๆ ที่เกี่ยวข้อง เนื่องจากงานวิจัยนี้มุ่งเน้นสำหรับการศึกษา แบบจำลองที่ไม่ซับซ้อนมากนักนำมาเปรียบเทียบกัน และหาส่วนที่มีผลต่อการย่อยสลายของ วัคซีนเอ็มอาร์เอ็นเอ ซึ่งเทคนิคนี้อาจถูกใช้เพื่อทดสอบคัดกรองเพื่อลบลำดับที่ไม่เสถียรออกไป ผู้วิจัยใช้ผลของ Public Leaderboard จากชุมชน Kaggle เมื่อนำไปเปรียบเทียบกับผลการทดลอง อื่นๆซึ่งอาจจะเหมาะเมื่อนำไปต่อยอดด้วยการวิเคราะห์ด้านอื่นๆประกอบกัน

## 5.2 ข้อเสนอแนะ

ข้อจำกัดของงานนี้ที่พบคือความยาวของลำดับของโมเลกุลเอ็มอาร์เอ็นเอที่ใช้ความยาว ที่ใช้ในงานวิจัยนี้อยู่ระหว่าง 107-130 ฐานและเป็นข้อมูลที่ถูกจำลองขึ้นในห้องปฏิบัติการเพื่อใช้ในการทดสอบ ในขณะที่วัคซีน โควิด 19 ที่แท้จริงน่าจะมีขนาดยาวตั้งแต่ 3,000-4,000 ฐาน แต่ผลของงานนี้แสดงให้เห็นว่าอัลกอริทึมการทำนายดังกล่าวมีความเป็นไปได้และมีศักยภาพที่ช่วยย่น เวลา ในระหว่างกระบวนการวิจัยในระหว่างการระบาดของโรค ซึ่งในระยะยาวเทคนิคดังกล่าวอาจ ช่วยในการวิจัยได้ดีขึ้นในการทำความเข้าใจเหตุผลเบื้องหลังความเสถียรของโมเลกุลเอ็มอาร์เอ็นเอ

เอบางชนิดและช่วยในการพัฒนาเทคโนโลยีที่เกี่ยวข้อง โดยหวังเป็นอย่างยิ่งว่างานนี้จะเป็นประโยชน์ต่อนักวิทยาศาสตร์ข้อมูลคนอื่น ๆ ในการสร้างการคาดการณ์ที่ดีขึ้นในอนาคต ในการปรับปรุงเสถียรภาพของวัคซีนเอ็มอาร์เอ็นเอจริงๆเป็นปัญหาที่ได้รับการสำรวจก่อนเกิดโรคระบาด การแก้ปัญหาโดยกระบวนการเรียนรู้ของเครื่องจะถูกใช้ในห้องปฏิบัติการ ก่อนที่จะนำเข้าสู่ภาคอุตสาหกรรม เพื่อเร่งการวิจัยวัคซีน mRNA และส่งมอบวัคซีนที่มีความเสถียรในตู้เย็นเพื่อป้องกันไวรัส SARS-CoV-2 ซึ่งเป็นไวรัสที่เป็นตัวก่อโรคโควิด-19



## บรรณานุกรม

- Amgad Muneer, S. M. F., Nur Arifin Akbar, David Agustriawan, Setyanto Tri Wahyudi,. (2022). iVaccine-Deep: Prediction of COVID-19 mRNA vaccine degradation using deep learning,. *Journal of King Saud University - Computer and Information Sciences*,.
- Arash Keshavarzi Arshadi, J. W., Milad Salem, Emmanuel Cruz,. (2020). Artificial Intelligence for COVID-19 Drug Discovery and Vaccine Development.
- Bowman, B. N., McAdam, P. R., Vivona, S., Zhang, J. X., Luong, T., Belew, R. K., . . . Woelk, C. H. (2011). Improving reverse vaccinology with a machine learning approach. *Vaccine*, 29(45), 8156-8164.
- Cheng F, W. Y., Bai Y, Liang Z, Mao Q, Liu D, Wu X, Xu M. (2023). Research Advances on the Stability of mRNA Vaccines. *Viruses*. .
- Eternagame. (2020). How to build a better vaccine from the comfort of your own web browser. สืบค้นจาก <https://medium.com/eternaproject/how-to-build-a-better-vaccine-from-the-comfort-of-your-own-web-browser-233343e0210d>
- Florindo, H. F., Kleiner, R., Vaskovich-Koubi, D., Acúrcio, R. C., Carreira, B., Yeini, E., . . . Satchi-Fainaro, R. (2020). Immune-mediated approaches against COVID-19. *Nature Nanotechnology*, 15(8), 630-645.
- Mackenzie, R. J. (2020). DNA vs. RNA – 5 Key Differences and Comparison. <https://www.technologynetworks.com/genomics/articles/what-are-the-key-differences-between-dna-and-rna-296719>
- Ong, E., Wong, M. U., Huffman, A., และ He, Y. (2020). COVID-19 coronavirus vaccine design using reverse vaccinology and machine learning. *bioRxiv*.
- Pardi, N., Hogan, M. J., Porter, F. W., และ Weissman, D. (2018). mRNA vaccines - a new era in vaccinology. *Nat Rev Drug Discov*, 17(4), 261-279.
- S. Asif Imran, M. T. I., C. Shahnaz, M. Tafhimul Islam, O. Tawhid Imam and M. Haque, . (2020). COVID-19 mRNA Vaccine Degradation Prediction using Regularized LSTM Model,. *2020 IEEE International Women in Engineering (WIE) Conference on*



*Electrical and Computer Engineering (WIECON-ECE).*

S. H. Ing, A. A. A. a. S. K. (2021). Development of COVID-19 mRNA Vaccine Degradation Prediction System. *IEEE*.

Schlake, T., Thess, A., Fotin-Mieczek, M., และ Kallen, K. J. (2012). Developing mRNA-vaccine technologies. *RNA Biol*, 9(11), 1319-1330.

Stanford. (2020). OpenVaccine: COVID-19 mRNA Vaccine Degradation Prediction. สืบค้นจาก <https://www.kaggle.com/c/stanford-covid-vaccine/data>

Uddin MN, R. M. (2021). Challenges of Storage and Stability of mRNA-Based COVID-19 Vaccines. *Vaccines (Basel)*.

Wadhwa A, A. A., Lokras A, Foged C, Thakur A. . (2020). Opportunities and Challenges in the Delivery of mRNA-based Vaccines. *Pharmaceutics*.

Wang, Z. (2023). Comparative analysis of tissue-specific genes in maize based on machine learning models: CNNperforms technicallybest, LightGBM performs biologically soundest. *Frontiers in Genetics*.

Wayment-Steele, H. K., Kladwang, W., Watkins, A.M. et al. (2022). Deep learning models for predicting RNA degradation via dual crowdsourcing.

Wurtmann, E. J., และ Wolin, S. L. (2009). RNA under attack: cellular handling of RNA damage. *Crit Rev Biochem Mol Biol*, 44(1), 34-49.

เจาะลึกระบบสุขภาพ, ส. H. (2023). กรมควบคุมโรค เผยเคสโควิดยังเพิ่ม ผู้เชี่ยวชาญห่วงโควิดรุนแรงในเด็กอาจส่งผลกระทบต่อพัฒนาการ ย้ำวัคซีนจำเป็นในเด็กเล็ก.

<https://www.hfocus.org/content/2023/06/27769>

ช่างแก้ว, น. พ. (2021). ตอบทุกข้อสงสัย “วัคซีนโควิด-19”. <https://thainakarin.co.th/covid-19-vaccine-knowledge/>

ไทยโพสต์. (2020). ชาวดีมาก วัคซีน'โมเดอร์นา'ป้องกันโควิดได้ผลเกือบ95%. สืบค้นจาก <https://www.thaipost.net/main/detail/84053>

มหาวิทยาลัยมหิดล, ค. (2020). วัคซีน คืออะไร...แล้วทำไมเราถึงต้องฉีดวัคซีน ? สืบค้นจาก <https://www.gj.mahidol.ac.th/main/knowledge-2/what-is-vaccine/>

โมเลกุลอาร์เอ็นเอ (RNA Molecules) ).

รัตนโรจน์พงศ์, โ. (2021). โควิน่าไวรัสสายพันธุ์ใหม่ (โควิด-19): องค์ความรู้ด้านงานวิจัยไวรัสวิทยา

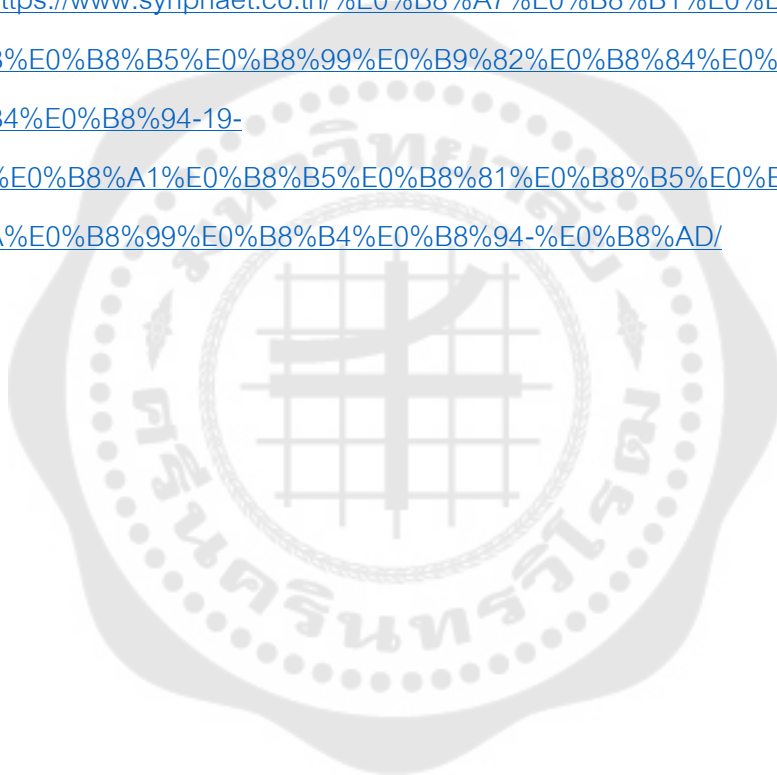
การศึกษาข้อมูลด้านการตรวจวินิจฉัย การรักษา และการพัฒนาวัคซีน.

โรงพยาบาลศิริรินทร์. วัคซีนโควิดมีกี่ชนิด?

<https://www.sikarin.com/health/%E0%B8%A7%E0%B8%B1%E0%B8%84%E0%B8%8B%E0%B8%B5%E0%B8%99%E0%B9%82%E0%B8%84%E0%B8%A7%E0%B8%B4%E0%B8%94%E0%B8%A1%E0%B8%B5%E0%B8%81%E0%B8%B5%E0%B9%88%E0%B8%8A%E0%B8%99%E0%B8%B4%E0%B8%94>

ศักดิ์ทองจีน, พ. ร. (2021). วัคซีนโควิด-19 มีกี่ชนิด อะไรบ้าง.

<https://www.synphaet.co.th/%E0%B8%A7%E0%B8%B1%E0%B8%84%E0%B8%8B%E0%B8%B5%E0%B8%99%E0%B9%82%E0%B8%84%E0%B8%A7%E0%B8%B4%E0%B8%94-19-%E0%B8%A1%E0%B8%B5%E0%B8%81%E0%B8%B5%E0%B9%88%E0%B8%8A%E0%B8%99%E0%B8%B4%E0%B8%94-%E0%B8%AD/>





## รายละเอียดได้คือการตรวจสอบข้อมูลและสำรวจข้อมูลเบื้องต้น (Exploratory Data Analysis)

```

1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 import numpy as np # linear algebra
5 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import warnings
9 import json
10 from collections import Counter
11 warnings.simplefilter(action='ignore')
12 pd.set_option('display.max_rows', None)
13 pd.set_option('display.max_columns', None)
14 plt.rcParams.update({'figure.max_open_warning': 0})
15
16 train = pd.read_json('/content/drive/My Drive/IS/train.json', lines=True)
17 test = pd.read_json('/content/drive/My Drive/IS/test.json', lines=True)
18
19 print('train shapes: ', train.shape)
20 print('test shapes: ', test.shape)
21
22 train.head(5)
23 test.head(5)
24
25 train.head(3)
26 train.dtypes.value_counts().plot.bar()
27
28 train.info()
29
30 train.describe()
31
32 test.info()
33
34 train.isna().sum()
35
36 print('train shape :',train.shape)
37 print('test shape :',test.shape)
38
39 train.head()
40 test.head()
41 train.describe()
42
43 fig, axe = plt.subplots(1, 2, figsize=(10, 5))
44 sns.countplot(train['SN_filter'], palette='magma',ax=axe[1])
45 sns.distplot(train['signal_to_noise'],color='pink',ax=axe[0])
46
47
48 axe[0].set_title('Signal/Noise ')
49 axe[1].set_title('SN_Filter');
50
51
52 Run Cell | Run Below | Debug cell
53 # %%
54 plt.figure(figsize=(10,4))
55 sns.countplot(test['seq_length'],color='blue')
56
57 Run Cell | Run Above | Debug cell
58 # %%
59 plt.figure(figsize=(15,2))
60 sns.boxplot(train['signal_to_noise'],color='yellow',x='signal_to_noise')
61
62 plt.show()
63
64 Run Cell | Run Above | Debug cell
65 # %%
66 avg_reactivity = np.array(list(map(np.array,train.reactivity))).mean(axis=0)
67 avg_deg_50C = np.array(list(map(np.array,train.deg_50C))).mean(axis=0)
68 avg_deg_pH10 = np.array(list(map(np.array,train.deg_pH10))).mean(axis=0)
69
70 avg_deg_Mg_50C = np.array(list(map(np.array,train.deg_Mg_50C))).mean(axis=0)
71 avg_deg_Mg_pH10 = np.array(list(map(np.array,train.deg_Mg_pH10))).mean(axis=0)

```

```

74 fig, ax = plt.subplots(1,3,figsize=(15,4))
75
76 # Distribution of Reactivity Averaged over position
77 sns.kdeplot(avg_reactivity,color='red',ax=ax[0])
78 ax[0].set_title('Average Reactivity',size=10)
79
80 # Distribution of deg_50C Averaged over position
81 sns.kdeplot(avg_deg_50C,color='blue',ax=ax[1])
82 ax[1].set_title('Average deg_50C',size=10)
83
84 # Distribution of deg_pH10 Averaged over position
85 sns.kdeplot(avg_deg_pH10,color='green',ax=ax[2])
86 ax[2].set_title('Average deg_pH10',size=10)
87
88
89 plt.show()
90
91
Run Cell | Run Above | Debug cell
92 # %%
93 plt.figure(figsize=(15,8))
94
95 sns.lineplot(x=range(68),color='red',y=avg_reactivity,label='reactivity')
96 sns.lineplot(x=range(68),color='blue',y=avg_deg_50C,label='deg_50C')
97 sns.lineplot(x=range(68),color='green',y=avg_deg_pH10,label='deg_pH10')
98
99 plt.xlabel('Positions')
100 plt.xticks(range(0,68,2))
101 plt.ylabel('Values')
102 plt.title('Average Target Values (without Mg) and Positions')
103
104 plt.show()
108 plt.figure(figsize=(15,5))
109
110 sns.lineplot(x=range(68),y=avg_deg_Mg_50C,label='deg_Mg_50C')
111 sns.lineplot(x=range(68),y=avg_deg_Mg_pH10,label='deg_Mg_pH10')
112
113 plt.xlabel('Positions')
114 plt.xticks(range(0,68,2))
115 plt.ylabel('Values')
116 plt.title('Average Target Values (w Mg) and Positions')
117
118 plt.show()
119
120
Run Cell | Run Above | Debug cell
121 # %%
122 One_sample = train.iloc[0]
123 One_sequence = One_sample["sequence"]
124 One_structure = One_sample["structure"]
125 One_predicted_loop_type = One_sample["predicted_loop_type"]
126 One_reactivity = One_sample["reactivity"]
127
128
Run Cell | Run Above | Debug cell
129 # %%
130 print("One_sequence :",One_sequence)
131 print("One_structure :",One_structure)
132 print("One_predicted_loop_type :",One_predicted_loop_type)
133 print("One_reactivity :", One_reactivity)
134

```

```
141 Counter(One_sequence)
142
143
Run Cell | Run Above | Debug cell
144 # %%
145 sum_one_sequence = sum(Counter(One_sequence).values())
146 sum_one_sequence
147
148
Run Cell | Run Above | Debug cell
149 # %%
150 Counter(One_structure)
151
152
Run Cell | Run Above | Debug cell
153 # %%
154 sum_one_structure = sum(Counter(One_structure).values())
155 sum_one_structure
156
157
Run Cell | Run Above | Debug cell
158 # %%
159 Counter(One_predicted_loop_type)
160
161
Run Cell | Run Above | Debug cell
162 # %%
163 sum_one_predicted_loop_type = sum(Counter(One_predicted_loop_type).values())
164 sum_one_predicted_loop_type
165
166
Run Cell | Run Above | Debug cell
167 # %%
168 print(f"Samples with signal_to_noise greater than 1: {len(train.loc[(train['signal_to_noise'] > 1]))}")
169 print(f"Samples with SN_filter = 1: {len(train.loc[(train['SN_filter'] == 1]))}")
170 print(f"Samples with signal_to_noise greater than 1, but SN_filter == 0: {len(train.loc[(train['signal_to_noise'] > 1) &
```



## รายละเอียดการทำแบบจำลอง XG boost

```

1  from google.colab import drive
2  drive.mount('/content/drive')
3  |
4  import json
5
6  import numpy as np
7  import pandas as pd
8
9  from sklearn.multioutput import MultiOutputRegressor
10 from sklearn.model_selection import train_test_split, GridSearchCV
11 from sklearn.pipeline import Pipeline
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.metrics import make_scorer
14
15 from xgboost import XGBRegressor
16
17
18 Run Cell | Run Below | Debug cell
19 # %%
20 train_df = pd.read_json('/content/drive/My Drive/IS/train.json', lines=True)
21 test_df = pd.read_json('/content/drive/My Drive/IS/test.json', lines=True)
22 sample_sub_df = pd.read_csv('/content/drive/My Drive/IS/sample_submission.csv')
23
24 Run Cell | Run Above | Debug cell
25 # %%
26 print(train_df.shape)
27 print(test_df.shape)
28 print(sample_sub_df.shape)
29
30 Run Cell | Run Above | Debug cell
31 # %%
32 train_df.head(3)
33
34 Run Cell | Run Above | Debug cell
35 # %%
36 test_df.head(3)
37
38 Run Cell | Run Above | Debug cell
39 # %%
40 sample_sub_df.head(3)
41
42 Run Cell | Run Above | Debug cell
43 # %%
44 # Calculate Means of targets
45 train_df['reactivity'] = train_df['reactivity'].apply(lambda x: np.mean(x))
46 train_df['deg_Mg_pH10'] = train_df['deg_Mg_pH10'].apply(lambda x: np.mean(x))
47 train_df['deg_pH10'] = train_df['deg_pH10'].apply(lambda x: np.mean(x))
48 train_df['deg_Mg_50C'] = train_df['deg_Mg_50C'].apply(lambda x: np.mean(x))
49 train_df['deg_50C'] = train_df['deg_50C'].apply(lambda x: np.mean(x))
50
51 Run Cell | Run Above | Debug cell
52 # %%
53 train_df.head()
54
55 Run Cell | Run Above | Debug cell
56 # %%
57 # Drop unnecessary columns for now
58 train_df = train_df.drop(['id', 'index', 'reactivity_error', 'deg_error_Mg_pH10', 'deg_error_pH10', 'deg_error_Mg_50C', 'deg_error_50C'],
59 train_df.head()
60
61 Run Cell | Run Above | Debug cell
62 # %%
63 # Split data in features and labels
64 X_train = train_df.drop(['reactivity', 'deg_Mg_pH10', 'deg_Mg_50C'], axis=1)
65 Y_train = train_df[['reactivity', 'deg_Mg_pH10', 'deg_Mg_50C']]

```

```

68 X_train, X_test, Y_train, Y_test = train_test_split(X_train, Y_train, test_size=0.15)
69 X_train.shape, X_test.shape, Y_train.shape, Y_test.shape
70
71
72 Run Cell | Run Above | Debug cell
73 # %%
74 def featurize(df):
75     df['total_A_count'] = df['sequence'].apply(lambda s: s.count('A'))
76     df['total_G_count'] = df['sequence'].apply(lambda s: s.count('G'))
77     df['total_U_count'] = df['sequence'].apply(lambda s: s.count('U'))
78     df['total_C_count'] = df['sequence'].apply(lambda s: s.count('C'))
79
80     df['total_dot_count'] = df['structure'].apply(lambda s: s.count('.'))
81     df['total_ob_count'] = df['structure'].apply(lambda s: s.count('('))
82     df['total_cb_count'] = df['structure'].apply(lambda s: s.count(')'))
83
84     return df
85
86
87 Run Cell | Run Above | Debug cell
88 # %%
89 X_train = featurize(X_train)
90 X_test = featurize(X_test)
91
92 Run Cell | Run Above | Debug cell
93 # %%
94 X_train = X_train.drop(['sequence', 'structure', 'predicted_loop_type'], axis=1)
95 X_test = X_test.drop(['sequence', 'structure', 'predicted_loop_type'], axis=1)
96
97 Run Cell | Run Above | Debug cell
98 # %%
99 X_train.head()
100 scaler = StandardScaler()
101 scaler.fit(X_train)
102
103 X_train = scaler.transform(X_train)
104 X_test = scaler.transform(X_test)
105
106 Run Cell | Run Above | Debug cell
107 # %%
108 def mcrmse_loss(y_true, y_pred, N=3):
109     """
110     Calculates competition eval metric
111     """
112     assert len(y_true) == len(y_pred)
113     n = len(y_true)
114     return np.sum(np.sqrt(np.sum((y_true - y_pred)**2, axis=0)/n)) / N
115
116 custom_scorer = make_scorer(mcrmse_loss, greater_is_better=False)
117
118 Run Cell | Run Above | Debug cell
119 # %%
120 # Hyperparameter tune multioutput XGBoost Regressor
121 xgb = XGBRegressor(
122     subsample=0.8,
123     colsample_bytree=0.75,
124     reg_lambda=2,
125     reg_alpha=1,
126     random_state=28
127 )
128
129 params = {
130     'estimator__n_estimators': [ 800, 900, 1000],
131     'estimator__learning_rate': [0.1, 0.25, 0.30],
132     'estimator__max_depth': [3, 4, 5]
133 }
134

```



```

137 reg = MultiOutputRegressor(xgb)
138
139 # Perform 5 fold cross validation on set 1 features
140 gs = GridSearchCV(reg, param_grid=params, cv=5, return_train_score=True, n_jobs=-1, sc
141 gs.fit(X_train, Y_train)
142
143
144 Run Cell | Run Above | Debug cell
145 # %%
146 gs.best_params_
147
148 Run Cell | Run Above | Debug cell
149 # %%
150 # Train using best parameters
151 xgb = XGBRegressor(
152     max_depth=gs.best_params_['estimator__max_depth'],
153     subsample=0.8,
154     colsample_bytree=0.75,
155     reg_lambda=2,
156     reg_alpha=1,
157     n_estimators=gs.best_params_['estimator__n_estimators'],
158     learning_rate=gs.best_params_['estimator__learning_rate'],
159     random_state=28
160 )
161 reg = MultiOutputRegressor(xgb)
162 reg.fit(X_train, Y_train)
163
164
165 Run Cell | Run Above | Debug cell
166 # %%
167 # Train score
168 mcrmse_loss(reg.predict(X_train), np.array(Y_train))
169
170 # Validation score
171 mcrmse_loss(reg.predict(X_test), np.array(Y_test))
172
173
174
175 Run Cell | Run Above | Debug cell
176 # %%
177 test = featurize(test_df.drop(['index', 'id'], axis=1))
178 test = test.drop(['sequence', 'structure', 'predicted_loop_type'], axis=1)
179 test = scaler.transform(test)
180
181
182 Run Cell | Run Above | Debug cell
183 # %%
184 # Predict
185 preds = pd.DataFrame(reg.predict(test))
186
187
188 Run Cell | Run Above | Debug cell
189 # %%
190 # Create submission csv
191 submission_df = preds.loc[preds.index.repeat(list(test_df['seq_length']))].reset_index(drop=True)
192 submission_df = submission_df.rename(columns={0: 'reactivity', 1: 'deg_Mg_pH10', 2: 'deg_Mg_50C'})
193 submission_df['id_seqpos'] = sample_sub_df['id_seqpos']
194 submission_df['deg_pH10'] = 0.0
195 submission_df['deg_50C'] = 0.0
196 submission_df = submission_df[['id_seqpos', 'reactivity', 'deg_Mg_pH10', 'deg_pH10', 'deg_Mg_50C', 'deg_50C']]
197
198
199 Run Cell | Run Above | Debug cell
200 # %%
201 submission_df.to_csv('base.csv', index=False)

```

# วิ ำ ย ล ะ เ คี ย ต ก ำ ร ท ำ แ บ บ จ ำ ล อ ง Catboost

```

1 from google.colab import drive
2 drive.mount('/content/drive')
3 pip install catboost
4 import json
5
6 import numpy as np
7 import pandas as pd
8
9 from sklearn.multioutput import MultiOutputRegressor
10 from sklearn.model_selection import train_test_split, GridSearchCV
11 from sklearn.pipeline import Pipeline
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.metrics import make_scorer
14
15 from catboost import CatBoostRegressor
16
17 train_df= pd.read_json('/content/drive/My Drive/IS/train.json', lines=True)
18 test_df= pd.read_json('/content/drive/My Drive/IS/test.json', lines=True)
19 sample_sub_df= pd.read_csv('/content/drive/My Drive/IS/sample_submission.csv')
20 print(train_df.shape)
21 print(test_df.shape)
22 print(sample_sub_df.shape)
23 train_df.head(3)
24 test_df.head(3)
25 sample_sub_df.head(3)
26 # Calculate Means of targets
27 train_df['reactivity'] = train_df['reactivity'].apply(lambda x: np.mean(x))
28 train_df['deg_Mg_pH10'] = train_df['deg_Mg_pH10'].apply(lambda x: np.mean(x))
29 train_df['deg_pH10'] = train_df['deg_pH10'].apply(lambda x: np.mean(x))
30 train_df['deg_Mg_50C'] = train_df['deg_Mg_50C'].apply(lambda x: np.mean(x))
31 train_df['deg_50C'] = train_df['deg_50C'].apply(lambda x: np.mean(x))
32 train_df.head()
33 # Drop unnecessary columns for now
34 train_df = train_df.drop(['id', 'index', 'reactivity_error', 'deg_error_Mg_pH10', 'deg_error_pH10', 'deg_error_Mg_50C', 'deg_error_50C', 'SM_filter', 'signal_to_noise', ''])
35 train_df.head()
36 # Split data in features and labels
37 X_train = train_df.drop(['reactivity', 'deg_Mg_pH10', 'deg_Mg_50C'], axis=1)
38 Y_train = train_df[['reactivity', 'deg_Mg_pH10', 'deg_Mg_50C']]
39 X_test, X_train, Y_train, Y_test = train_test_split(X_train, Y_train, test_size=0.15)
40 X_train.shape, X_test.shape, Y_train.shape, Y_test.shape
41
42 def featurize(df):
43
44     df['total_A_count'] = df['sequence'].apply(lambda s: s.count('A'))
45     df['total_G_count'] = df['sequence'].apply(lambda s: s.count('G'))
46     df['total_U_count'] = df['sequence'].apply(lambda s: s.count('U'))
47     df['total_C_count'] = df['sequence'].apply(lambda s: s.count('C'))

```





## วิ ำ ย ล ะ ใ ้ ย ด ก ำ ร ท ำ ำ แ บ บ ำ ำ ล อ ง Random Forest

```

1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 import json
5
6 import numpy as np
7 import pandas as pd
8
9 from sklearn.multioutput import MultiOutputRegressor
10 from sklearn.model_selection import train_test_split, GridSearchCV
11 from sklearn.pipeline import Pipeline
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.metrics import make_scorer
14
15 from sklearn.ensemble import RandomForestRegressor
16
17 train_df = pd.read_json('/content/drive/My Drive/IS/train.json', lines=True)
18 test_df = pd.read_json('/content/drive/My Drive/IS/test.json', lines=True)
19 sample_sub_df = pd.read_csv('/content/drive/My Drive/IS/sample_submission.csv')
20
21 print(train_df.shape)
22 print(test_df.shape)
23 print(sample_sub_df.shape)
24
25 train_df.head(3)
26
27 test_df.head(3)
28
29 # Calculate Means of targets
30 train_df['reactivity'] = train_df['reactivity'].apply(lambda x: np.mean(x))
31 train_df['deg_Mg_pH10'] = train_df['deg_Mg_pH10'].apply(lambda x: np.mean(x))
32 train_df['deg_pH10'] = train_df['deg_pH10'].apply(lambda x: np.mean(x))
33 train_df['deg_Mg_50C'] = train_df['deg_Mg_50C'].apply(lambda x: np.mean(x))
34 train_df['deg_50C'] = train_df['deg_50C'].apply(lambda x: np.mean(x))
35
36 train_df.head(3)
37
38 # Drop unnecessary columns for now
39 train_df = train_df.drop(['id', 'index', 'reactivity_error', 'deg_error_Mg_pH10', 'deg_error_pH10', 'deg_error_Mg_50C', 'deg_error_50C', 'SN_filter', 'signal_to_noise', 'deg_
40 train_df.head()
41
42 # Split data in features and labels
43 X_train = train_df.drop(['reactivity', 'deg_Mg_pH10', 'deg_Mg_50C'], axis=1)
44 Y_train = train_df[['reactivity', 'deg_Mg_pH10', 'deg_Mg_50C']]
45
46 X_train, X_test, Y_train, Y_test = train_test_split(X_train, Y_train, test_size=0.15)
47 X_train.shape, X_test.shape, Y_train.shape, Y_test.shape

```



```

49 def featurize(df):
50
51     df['total_A_count'] = df['sequence'].apply(lambda s: s.count('A'))
52     df['total_G_count'] = df['sequence'].apply(lambda s: s.count('G'))
53     df['total_U_count'] = df['sequence'].apply(lambda s: s.count('U'))
54     df['total_C_count'] = df['sequence'].apply(lambda s: s.count('C'))
55
56     df['total_dot_count'] = df['structure'].apply(lambda s: s.count('.'))
57     df['total_ob_count'] = df['structure'].apply(lambda s: s.count('('))
58     df['total_cb_count'] = df['structure'].apply(lambda s: s.count('('))
59
60     return df
61
62 X_train = featurize(X_train)
63 X_test = featurize(X_test)
64
65 X_train = X_train.drop(['sequence', 'structure', 'predicted_loop_type'], axis=1)
66 X_test = X_test.drop(['sequence', 'structure', 'predicted_loop_type'], axis=1)
67
68 X_train.head()
69
70 scaler = StandardScaler()
71 scaler.fit(X_train)
72
73 X_train = scaler.transform(X_train)
74 X_test = scaler.transform(X_test)
75
76 def mcrmse_loss(y_true, y_pred, N=3):
77     """
78     Calculates competition eval metric
79     """
80     assert len(y_true) == len(y_pred)
81     n = len(y_true)
82     return np.sum(np.sqrt(np.sum((y_true - y_pred)**2, axis=0)/n)) / N
83
84 custom_scorer = make_scorer(mcrmse_loss, greater_is_better=False)
85
86 rf = RandomForestRegressor(random_state=42, verbose=1)
87 params = {
88     'estimator__n_estimators': [100, 200, 300],
89     'estimator__max_depth': [8, 9, 10],
90 }
91
92 reg = MultiOutputRegressor(rf)
93 # Perform 5 fold cross validation on set 1 features
94 rfr = GridSearchCV(estimator=reg, param_grid = params, cv = 5, n_jobs=-1, scoring=custom_scorer)
95
96 rfr.fit(X_train, Y_train)
97
98 rfr.best_params_
99
100 # Train using best parameters
101 rf = RandomForestRegressor(n_jobs=-1, random_state=28, verbose=1, n_estimators=500, max_depth=13)
102
103 rfr = MultiOutputRegressor(rf)
104 rfr.fit(X_train, Y_train)
105
106 # Train score
107 mcrmse_loss(rfr.predict(X_train), np.array(Y_train))
108
109 # Validation score
110 mcrmse_loss(rfr.predict(X_test), np.array(Y_test))
111
112 test = featurize(test_df.drop(['index', 'id'], axis=1))
113 test = test.drop(['sequence', 'structure', 'predicted_loop_type'], axis=1)
114 test = scaler.transform(test)
115
116 # Predict
117 preds = pd.DataFrame(rfr.predict(test))

```

## รายละเอียดการทำแบบจำลอง Ligth GBM

```

1  from google.colab import drive
2  drive.mount('/content/drive')
3
4  import gc
5  import os
6  import random
7
8  import lightgbm as lgb
9  import numpy as np
10 import pandas as pd
11 import seaborn as sns
12 import itertools
13
14 from matplotlib import pyplot as plt
15 from sklearn.metrics import mean_squared_error
16 from sklearn.preprocessing import LabelEncoder
17 from sklearn.model_selection import StratifiedKFold, KFold, GroupKFold
18 from sklearn.cluster import KMeans
19
20 sns.set(style='darkgrid')
21 SEEDS = 28
22
23 Run Cell | Run Below | Debug cell
24 # %%
25 def rmse(y_true, y_pred):
26     return (mean_squared_error(y_true, y_pred))**.5
27
28 # treemodel_wrapper
29 class TreeModel:
30     def __init__(self, model_type):
31         self.model_type = model_type
32         self.tr_data = None
33         self.vl_data = None
34         self.model = None
35
36     def train(self, params, train_x, train_y, valid_x=None, valid_y=None, num_round=None, early_stopping=None, verbose=None):
37         if self.model_type == 'lgb':
38             self.tr_data = lgb.Dataset(train_x, label=train_y)
39             self.vl_data = lgb.Dataset(valid_x, label=valid_y)
40             self.model = lgb.train(params, [self.tr_data, self.vl_data],
41                                   num_boost_round=num_round, early_stopping_rounds=early_stopping, verbose_eval=verbose)
42
43         if self.model_type == 'rf_reg':
44             self.train_x = train_x
45             self.train_y = train_y
46             self.model = RandomForestRegressor(*params).fit(self.train_x, self.train_y)
47
48         if self.model_type == 'xgb':
49             self.tr_data = xgb.DMatrix(train_x, train_y)
50             self.vl_data = xgb.DMatrix(valid_x, valid_y)
51             self.model = xgb.train(params, [self.tr_data, self.vl_data], num_boost_round=num_round,
52                                   evals=[(self.tr_data, 'train'), (self.vl_data, 'val')],
53                                   verbose_eval=verbose, early_stopping_rounds=early_stopping)
54
55         if self.model_type == 'cat':
56             params['num_boost_round'] = num_round
57             self.cat_cols = list(train_x.select_dtypes(include='object').columns)
58             self.tr_data = Pool(train_x, train_y, cat_features=self.cat_cols)
59             self.vl_data = Pool(valid_x, valid_y, cat_features=self.cat_cols)
60             self.model = CatBoost(params).fit(self.tr_data, eval_set=self.vl_data,
61                                             early_stopping_rounds=early_stopping, verbose=verbose, use_best_model=True)
62
63         return self.model
64
65     def predict(self, X):
66         if self.model_type == 'lgb':
67             return self.model.predict(X, num_iteration=self.model.best_iteration)
68
69         if self.model_type == 'rf_reg':
70             return self.model.predict(X)
71
72         if self.model_type == 'xgb':
73             X_DM = xgb.DMatrix(X)
74             return self.model.predict(X_DM)
75
76         if self.model_type == 'cat':
77             X_pool = Pool(X, cat_features=self.cat_cols)
78             return self.model.predict(X_pool)
79
80 @property
81 def feature_names_(self):
82     if self.model_type == 'lgb':
83         return self.model.feature_name()
84
85     if self.model_type == 'rf_reg':
86         return self.train_x.columns
87
88     if self.model_type == 'xgb':
89         return list(self.model.get_score(importance_type='gain').keys())
90
91     if self.model_type == 'cat':
92         return self.model.feature_names_

```

```

94     def feature_importances_(self):
95         if self.model_type == 'lgb':
96             return self.model.feature_importance(importance_type='gain')
97
98         if self.model_type == 'rf_reg':
99             return self.model.feature_importances_
100
101         if self.model_type == 'xgb':
102             return list(self.model.get_score(importance_type='gain').values())
103
104         if self.model_type == 'cat':
105             return self.model.feature_importances_
106
107
108 train= pd.read_json('/content/drive/My Drive/IS/train.json', lines=True)
109 test= pd.read_json('/content/drive/My Drive/IS/test.json', lines=True)
110 submission= pd.read_csv('/content/drive/My Drive/IS/sample_submission.csv')
111
112 train_data = []
113 for mol_id in train['id'].unique():
114     sample_data = train.loc[train['id'] == mol_id]
115     sample_seq_length = sample_data.seq_length.values[0]
116
117     for i in range(68):
118         sample_dict = {'id' : sample_data['id'].values[0],
119                       'id_seqpos' : sample_data['id'].values[0] + '_' + str(i),
120                       'sequence' : sample_data['sequence'].values[0][i],
121                       'structure' : sample_data['structure'].values[0][i],
122                       'predicted_loop_type' : sample_data['predicted_loop_type'].values[0][i],
123                       'reactivity' : sample_data['reactivity'].values[0][i],
124                       'reactivity_error' : sample_data['reactivity_error'].values[0][i],
125                       'deg_Mg_pH10' : sample_data['deg_Mg_pH10'].values[0][i],
126                       'deg_error_Mg_pH10' : sample_data['deg_error_Mg_pH10'].values[0][i],
127                       'deg_pH10' : sample_data['deg_pH10'].values[0][i],
128                       'deg_error_pH10' : sample_data['deg_error_pH10'].values[0][i],
129                       'deg_Mg_50C' : sample_data['deg_Mg_50C'].values[0][i],
130                       'deg_error_Mg_50C' : sample_data['deg_error_Mg_50C'].values[0][i],
131                       'deg_50C' : sample_data['deg_50C'].values[0][i],
132                       'deg_error_50C' : sample_data['deg_error_50C'].values[0][i]}
133
134         shifts = [1,2,3,4,5]
135         shift_cols = ['sequence', 'structure', 'predicted_loop_type']
136         for shift,col in itertools.product(shifts, shift_cols):
137             if i - shift >= 0:
138                 sample_dict['b'+str(shift)+'_'+col] = sample_data[col].values[0][i-shift]
139             else:
140                 sample_dict['b'+str(shift)+'_'+col] = -1
141
142             if i + shift <= sample_seq_length - 1:
143                 sample_dict['a'+str(shift)+'_'+col] = sample_data[col].values[0][i+shift]
144             else:
145                 sample_dict['a'+str(shift)+'_'+col] = -1
146
147         train_data.append(sample_dict)
148
149 train_data = pd.DataFrame(train_data)
150 train_data.head()
151
152
153 test_data = []
154 for mol_id in test['id'].unique():
155     sample_data = test.loc[test['id'] == mol_id]
156     sample_seq_length = sample_data.seq_length.values[0]
157     for i in range(sample_seq_length):
158         sample_dict = {'id' : sample_data['id'].values[0],
159                       'id_seqpos' : sample_data['id'].values[0] + '_' + str(i),
160                       'sequence' : sample_data['sequence'].values[0][i],
161                       'structure' : sample_data['structure'].values[0][i],
162                       'predicted_loop_type' : sample_data['predicted_loop_type'].values[0][i]}
163
164         shifts = [1,2,3,4,5]
165         shift_cols = ['sequence', 'structure', 'predicted_loop_type']
166         for shift,col in itertools.product(shifts, shift_cols):
167             if i - shift >= 0:
168                 sample_dict['b'+str(shift)+'_'+col] = sample_data[col].values[0][i-shift]
169             else:
170                 sample_dict['b'+str(shift)+'_'+col] = -1
171
172             if i + shift <= sample_seq_length - 1:
173                 sample_dict['a'+str(shift)+'_'+col] = sample_data[col].values[0][i+shift]
174             else:
175                 sample_dict['a'+str(shift)+'_'+col] = -1
176
177         test_data.append(sample_dict)
178 test_data = pd.DataFrame(test_data)
179 test_data.head()

```

```

182 # label_encoding
183 sequence_encmap = {'A': 0, 'G': 1, 'C': 2, 'U': 3}
184 structure_encmap = {'.' : 0, '(' : 1, ')' : 2}
185 looptype_encmap = {'S':0, 'E':1, 'H':2, 'I':3, 'X':4, 'W':5, 'B':6}
186
187 enc_targets = ['sequence', 'structure', 'predicted_loop_type']
188 enc_maps = [sequence_encmap, structure_encmap, looptype_encmap]
189
190 for t,m in zip(enc_targets, enc_maps):
191     for c in [c for c in train_data.columns if t in c]:
192         train_data[c] = train_data[c].replace(m)
193         test_data[c] = test_data[c].replace(m)
194
195 not_use_cols = ['id', 'id_seqpos']
196 features = [c for c in test_data.columns if c not in not_use_cols]
197 targets = ['reactivity', 'deg_Mg_pH10', 'deg_pH10', 'deg_Mg_50C', 'deg_50C']
198
199 FOLD_N = 5
200 gkf = GroupKFold(n_splits=FOLD_N)
201
202 params = {'objective': 'regression',
203          'boosting': 'gbdt',
204          'metric': 'rmse',
205          'learning_rate': 0.05,
206          'seed': SEEDS}
207
208 feature_importances = pd.DataFrame()
209 result = {}
210 oof_df = pd.DataFrame(train_data.id_seqpos)
211
212 for target in targets:
213     oof = pd.DataFrame()
214     preds = np.zeros(len(test_data))
215     scores = 0.0
216
217     for n, (tr_idx, vl_idx) in enumerate(gkf.split(train_data[features], train_data[target])):
218         tr_x, tr_y = train_data[features].iloc[tr_idx], train_data[target].iloc[tr_idx]
219         vl_x, vl_y = train_data[features].iloc[vl_idx], train_data[target].iloc[vl_idx]
220         vl_id = train_data['id_seqpos'].iloc[vl_idx]
221
222         model = TreeModel(model_type='lgb')
223         model.train(params, tr_x, tr_y, vl_x, vl_y,
224                   num_round=20000, early_stopping=100, verbose=0)
225
226         fi_tmp = pd.DataFrame()
227         fi_tmp['feature'] = model.feature_names_
228
229         fi_tmp['importance'] = model.feature_importances_
230         fi_tmp['fold'] = n
231         fi_tmp['target'] = target
232         feature_importances = feature_importances.append(fi_tmp)
233
234         vl_pred = model.predict(vl_x)
235         score = rmse(vl_y, vl_pred)
236         scores += score / FOLD_N
237         print(f'score : {score}')
238
239         oof = oof.append(pd.DataFrame({'id_seqpos': vl_id, target: vl_pred}))
240
241         pred = model.predict(test_data[features])
242         preds += pred / FOLD_N
243
244     oof_df = oof_df.merge(oof, on='id_seqpos', how='inner')
245     submission[target] = preds
246
247     print(f'{target}_rmse : {scores}')
248     result[target] = scores
249
250 display(result)
251 display(f'total : {np.mean(list(result.values()))}')
252
253 # feature_importances
254 for target in targets:
255     tmp = feature_importances[feature_importances.target==target]
256     order = list(tmp.groupby('feature').mean().sort_values('importance', ascending=False).index)
257
258     plt.figure(figsize=(10, 10))
259     sns.barplot(x="importance", y="feature", data=tmp, order=order)
260     plt.title(target)
261     plt.tight_layout()

```





## ประวัติผู้เขียน

ชื่อ-สกุล	ณัฐนรี พอสูงเนิน
วัน เดือน ปี เกิด	01 เมษายน 2533
สถานที่เกิด	กรุงเทพมหานคร
วุฒิการศึกษา	พ.ศ.2555 วิทยาศาสตรบัณฑิต สาขาวิชาเคมี จากมหาวิทยาลัยธรรมศาสตร์

