

# Introduction to Supervised Learning

Petchara Pattarakijwanich

Introduction to Data Science, 2 September 2022

# Goal of this week

- Intro to Supervised Learning
  - Types of Variables
  - Types of Methods
  - Test of Accuracy and Methodology
  - Goals of Supervised Learning

# Example of Supervised Learning

## Input

- GPA บ.ปวช. } Quantitative
- คะแนน TCAS } (order correlate?)
- ปว.ปวช.จากโรงเรียนต่าง (พอรู้, ต, ตชช, ตชด) } ordered Qualitative
- จังหวัดที่สอบ } Qualitative
- อาชีพพ่อแม่ } Qualitative
- โรงเรียนที่สอบ } Qualitative

## Output

- GPA บัณฑิต } Quantitative  $\Rightarrow$  Regression
- สอบได้ 1 หรือไม่? } Qualitative  $\Rightarrow$  Classification

$\Rightarrow$  ตัวอย่างของ input  
และ output

# Mathematical Formulation

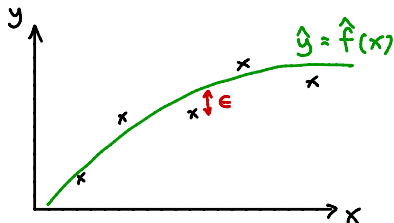
$x$  = input ,  $y$  = output

สมการรวมแล้ว

$$y = f(x) + \epsilon$$

Noise

$f$  คือ ฟังก์ชันที่เราสนใจ



ประมาณค่ารวมแล้ว

$$\hat{y} = \hat{f}(x)$$

$\hat{f}$  เป็น estimator ของ  $f$   
 $\hat{y}$  คือค่าทำนาย (predicted value) ของ  $y$  ที่ตรงกับ  $x$

Irreducible Noise  $\epsilon$

ไม่สามารถลดได้

สาเหตุ - ความไม่แน่นอน

- มี input ที่ไม่แน่นอน

-  $x, y$  ไม่เป็นอิสระ/เกี่ยวข้อง

$$E(\epsilon) = 0$$

$$\text{Var}(\epsilon) = \sigma^2$$

# Types of Variables

By function

- $x =$  input/independent/predictor/feature

↑  
ข้อมูล  
↓

↑  
ML/CS  
↓

- $y =$  output/dependent/response/target

# Types of Variables

By nature

- Quantitative    *ປຶ້ມຕັກເລນ (ຈຳນວນ/ປະມານ)*
  - ປຶ້ມສບເສບສັກໄດ້ ( $>$ ,  $<$  ສຳລານແລນບັດໂຕ)
  - ຕາມກຳໜົດ  $\Leftrightarrow$  ບໍ່ສາມາດຕັກໄດ້
- Qualitative (Categorical)    *ປຶ້ມຕັກ*
  - ປຶ້ມສບເສບສັກໄດ້
  - 2 ຕັກ (Binary) ຫຼື  $> 2$  ຕັກ (multi-class)
- Ordered Quantitative  
Ex. ລາຍຮັບປະຈຳປີ / ພະລັງ / ດີ / ດີແມ່

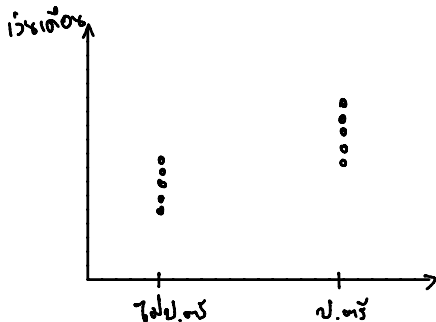
# How to Deal with Qualitative Parameters?

ឧទាហរណ៍ ប្រព័ន្ធគ្រប់គ្រងសំបុត្រស្រាវជ្រាវ

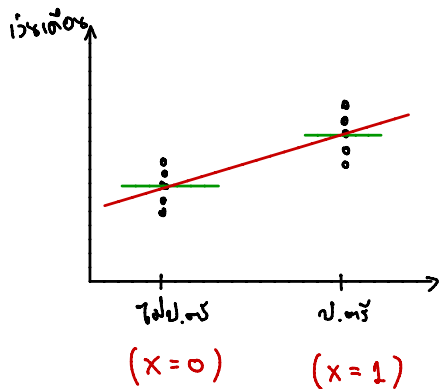
↑  
output  
(Quantitative)

↑  
Input  
(Qualitative (ស្រាវជ្រាវ))

ល.ស្រាវជ្រាវ	លេខសំបុត្រ
Y	...
Z	...
...	...



# How to Deal with Qualitative Parameters?



สมการเส้นตรง

$$\hat{f}(x) = \begin{cases} y_1, & \text{ไม่พบ.นศ} \\ y_2, & \text{พบ.นศ} \end{cases}$$

สมการเส้นตรง

⇒ หาสมการ  $x$  ที่ทำให้ค่า  $y$  เท่ากัน

⇒ Regression

$$y = mx + c$$

ตัวอย่าง?

ไม่พบ.นศ  $y = c$

พบ.นศ  $y = m + c$

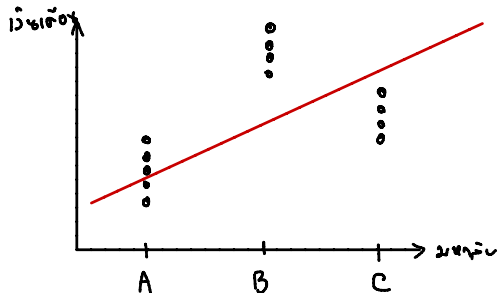
$y_{\text{พบ.นศ}} - y_{\text{ไม่พบ.นศ}} = m$



# How to Deal with Qualitative Parameters?

ข้อสงสัย วิเคราะห์ vs ขาด้านที่สนใจ (A,B,C)

ขาด้าน	วิเคราะห์
A	.
B	.
C	.
...	.



( $x=0$ ) ( $x=1$ ) ( $x=2$ )  
วิเคราะห์ไม่ได้!

# How to Deal with Qualitative Parameters?

โจทย์ โจทย์ 04 vs ขาดตัวแปร (A,B,C)

$$x_1 = \begin{cases} 0 & , A \\ 1 & , \text{ไม่ } A \end{cases}$$

$$x_2 = \begin{cases} 0 & , B \\ 1 & , \text{ไม่ } B \end{cases}$$

ขาดตัวแปร	$x_1$	$x_2$
A	0	1
B	1	0
C	1	1

ฟังก์ชัน

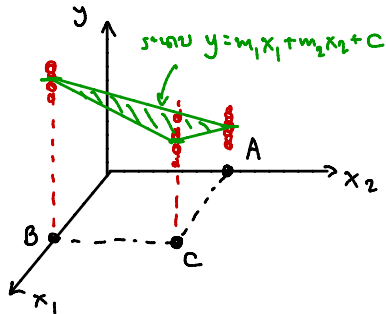
$$y = m_1 x_1 + m_2 x_2 + C$$

ตัวอย่าง

ขาดตัวแปร A :  $y = m_2 + C$

ขาดตัวแปร B :  $y = m_1 + C$

ขาดตัวแปร C :  $y = m_1 + m_2 + C$



# Types of Methods

By output

- Regression    Output = Quantitative  
Ex. น้ำหนัก
- Classification    Output = Qualitative  
Ex. ฤดูฝน

# Types of Methods

By technique  $y = f(x) + \epsilon$

- Parametric

- สมมติตั้งฟังก์ชัน  $f$  ไว้
- ทดสอบ  $\text{function fitting}$
- ข้อดี : ง่าย รวดเร็ว ประหยัด
- ข้อเสีย : ฟังก์ชันที่เรา  $f$  อาจผิดได้

- Non-parametric

- ยึดหลักการเรียนรู้  $f$
- $f$  ถูกปรับให้เหมาะกับข้อมูล
- ข้อดี : ไม่จำเป็นต้องสมมติฟังก์ชัน
- ข้อเสีย : ง่าย รวดเร็ว ประหยัด

# How to Measure Accuracy?

- Regression

MSE = Mean Squared Error

"Loss Function"  $= \frac{1}{n} \sum (y - \hat{y})^2 \Rightarrow \text{loss} = \text{ผิดพลาด}$

loss กับ test data

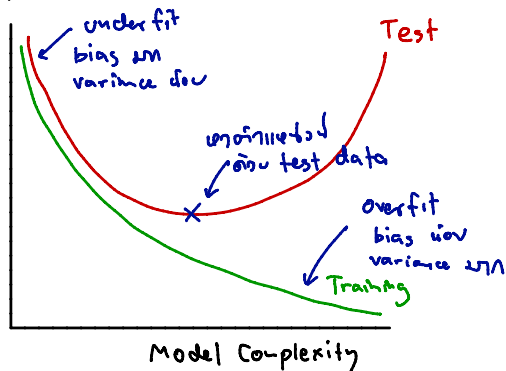
- Classification

Error Rate =  $\frac{\text{จำนวนครั้งที่ทายผิด}}{\text{จำนวนครั้งที่ถาม}}$       loss = ทายไม่ถูก

Confusion Matrix  
(ชุดข้อมูล)

# Bias-Variance Tradeoff

Loss function



# Methodology

- Training

ใช้ input/output ฝึกฝน  $\hat{f}$

- Testing (validation)

ใช้ input/output ฝึกสอน  $\hat{f}$  [  $\hat{y} = \hat{f}(x) \Rightarrow$  เปรียบ  $y, \hat{y}$  ]

- Application

ใช้ input ใส่ output ฝึก  $\hat{f}$  ทำนาย  $\hat{y}$

# Goals of Supervised Learning

- Prediction

- កំណត់  $y$  ចេញពី  $x$
- ភាពត្រឹមត្រូវ (Accuracy ខ្ពស់)
- ក្រៅពី  $\hat{f}$  យើង Black box (មិនដឹងពីអ្វីក្នុងខ្លួន)

- Inference

- ដើម្បី បញ្ជាក់ លទ្ធផលដែលបានមកពីការគណនា  $x, y$   
(Ex. តើការប្រើប្រាស់ថ្នាំបំបាត់ជំងឺបានឬទេ? តើការប្រើប្រាស់ថ្នាំបំបាត់ជំងឺបានឬទេ?)
- ការប្រើប្រាស់លទ្ធផលដែលបានមកពីការគណនា