

1.)

$$\text{gini index} = 1 - \left( \left( \frac{10}{20} \right)^2 + \left( \frac{10}{20} \right)^2 \right) = \frac{1}{2}$$

$$\text{Entropy} = - \left( \frac{5}{10} \right) \log_2 \left( \frac{5}{10} \right) - \left( \frac{5}{10} \right) \log_2 \left( \frac{5}{10} \right) = 1$$

ข. 1 - 20

$$\text{Gini index} = 1 - (0 + 1^2) = 0$$

$$\text{Information gain} = 1 - 0 = 1$$

$$\text{Entropy} = -0 \log(0) - 1 \log(1) = 0$$

ค. M

$$\text{Gini index} = 1 - \left( \left( \frac{6}{10} \right)^2 + \left( \frac{4}{10} \right)^2 \right) = 0.52$$

$$\text{Information gain} = 1 - 0.97 = 0.03$$

$$\text{Entropy} = - \left( \frac{6}{10} \right) \log_2 \left( \frac{6}{10} \right) - \left( \frac{4}{10} \right) \log_2 \left( \frac{4}{10} \right) = 0.97$$

F

$$\text{Gini index} = 1 - \left( \left( \frac{4}{10} \right)^2 + \left( \frac{6}{10} \right)^2 \right) = 0.52$$

$$\text{Information gain} = 1 - 0.97 = 0.03$$

$$\text{Entropy} = - \left( \frac{4}{10} \right) \log_2 \left( \frac{4}{10} \right) - \left( \frac{6}{10} \right) \log_2 \left( \frac{6}{10} \right) = 0.97$$

ง. Family

$$\text{Gini index} = 1 - \left( \left( \frac{1}{4} \right)^2 + \left( \frac{3}{4} \right)^2 \right) = 0.625$$

$$\text{Information gain} = 1 - 0.81 = 0.19$$

$$\text{Entropy} = - \left( \frac{1}{4} \right) \log_2 \left( \frac{1}{4} \right) - \left( \frac{3}{4} \right) \log_2 \left( \frac{3}{4} \right) = 0.81$$

Sport

$$\text{Gini index} = 1 - \left( \left( \frac{8}{8} \right)^2 + \left( \frac{0}{8} \right)^2 \right) = 0$$

$$\text{Information gain} = 1 - 0 = 1$$

$$\text{Entropy} = - \left( \frac{8}{8} \right) \log_2 \left( \frac{8}{8} \right) - \left( \frac{0}{8} \right) \log_2 \left( \frac{0}{8} \right) = 0$$

Luxury

$$\text{Gini index} = 1 - \left( \left( \frac{7}{8} \right)^2 + \left( \frac{1}{8} \right)^2 \right) = 0.78$$

$$\text{Information gain} = 1 - 0.54 = 0.46$$

$$\text{Entropy} = - \left( \frac{7}{8} \right) \log_2 \left( \frac{7}{8} \right) - \left( \frac{1}{8} \right) \log_2 \left( \frac{1}{8} \right) = 0.54$$

จ. Small

$$\text{Gini index} = 1 - \left( \left( \frac{3}{5} \right)^2 + \left( \frac{2}{5} \right)^2 \right) = 0.52$$

$$\text{Information gain} = 1 - 0.97 = 0.03$$

$$\text{Entropy} = - \left( \frac{3}{5} \right) \log_2 \left( \frac{3}{5} \right) - \left( \frac{2}{5} \right) \log_2 \left( \frac{2}{5} \right) = 0.97$$

Medium

$$\text{Gini index} = 1 - \left( \left( \frac{3}{7} \right)^2 + \left( \frac{4}{7} \right)^2 \right) = 0.51$$

$$\text{Information gain} = 1 - 0.985 = 0.015$$

$$\text{Entropy} = - \left( \frac{3}{7} \right) \log_2 \left( \frac{3}{7} \right) - \left( \frac{4}{7} \right) \log_2 \left( \frac{4}{7} \right) = 0.985$$

Large

$$\text{Gini index} = 1 - \left( \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right) = 0.5$$

$$\text{Information gain} = 1 - 1 = 0$$

$$\text{Entropy} = - \left( \frac{2}{4} \right) \log_2 \left( \frac{2}{4} \right) - \left( \frac{2}{4} \right) \log_2 \left( \frac{2}{4} \right) = 1$$

Extra large

$$\text{Gini index} = 1 - \left( \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right) = 0.5$$

$$\text{Information gain} = 1 - 1 = 0$$

$$\text{Entropy} = - \left( \frac{2}{4} \right) \log_2 \left( \frac{2}{4} \right) - \left( \frac{2}{4} \right) \log_2 \left( \frac{2}{4} \right) = 1$$

ฉ. Customer ID

ข. เหมาะ เพราะ มีความบริสุทธิ์มากที่สุด

## 2. จงอธิบายความหมายของ underfitting และ overfitting

๓๐๒ **Model underfitting** เกิดเมื่อโมเดลมีความซับซ้อนไม่พอที่จะเรียนรู้รูปแบบที่แฝงอยู่ในข้อมูลได้ สามารถสังเกตได้จากการที่ทั้ง training error และ test error มีค่าสูง

**Model overfitting** เกิดเมื่อโมเดลมีความซับซ้อนมากเกินไปทำให้เรียนรู้รูปแบบปลอม (spurious patterns) แทนที่จะเป็นรูปแบบที่แท้จริงตามธรรมชาติของข้อมูล สามารถสังเกตได้จากการที่ training error มีค่าต่ำหรือลดลงในขณะที่ test error มีค่าสูงหรือเพิ่มขึ้น

3. กำหนดผลการทดสอบประสิทธิภาพของโมเดลการจำแนกประเภทแบบสองคลาสดัง confusion matrix ต่อไปนี้ จงประเมินประสิทธิภาพของโมเดลดังกล่าวโดยคำนวณหาค่า Accuracy, Error rate, Recall, Precision, Specificity, False Positive Rate และ F1-score

		Predicted Class	
		Positive Class	Negative Class
Actual Class	Positive Class	180	20
	Negative Class	60	340

	Accuracy	Error rate	Recall	Precision	Specificity	FPR	F1-Score
Model 1	$\frac{180+340}{180+20+60+340}$	$\frac{20+60}{180+20+60+340}$	$\frac{180}{180+20}$	$\frac{180}{180+60}$	$\frac{340}{60+340}$	$\frac{60}{60+340}$	$\frac{2(0.75)(0.9)}{0.75 + 0.9}$
	0.866	0.133	0.9	0.75	0.85	0.15	0.818

4.  $\text{Cost}(\text{internal}) = \log m$   $\text{Cost}(\text{leaf}) = \log 3$

Tree1

Tree2

$$\text{Cost}(\text{tree1}) = 2 * \log m + 3 * \log 3$$

$$\text{Cost}(\text{tree2}) = 4 * \log m + 5 * \log 3$$

$$= \log m^2 + \log 3^3$$

$$= \log m^4 3^5$$

$$= \log m^2 3^3$$

$$\text{Cost}(D|\text{tree1}) = 7 * \log n$$

$$\text{Cost}(D|\text{tree2}) = 4 * \log n$$

$$\text{Cost}(\text{tree1}, D) = \log m^2 3^3 + 7 * \log n$$

$$\text{Cost}(\text{tree2}, D) = \log m^4 3^5 + 4 * \log n$$

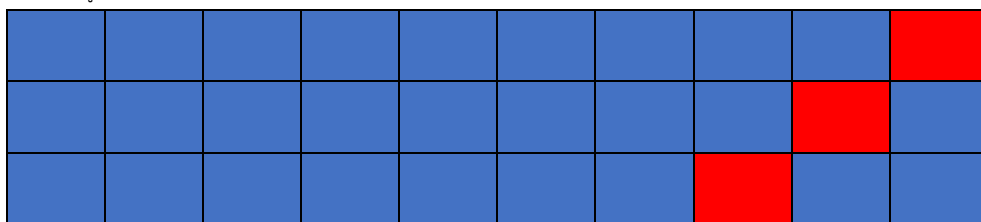
$$= \log m^2 3^3 n^7$$

$$= \log m^4 3^5 n^4$$

$$\text{Cost}(\text{tree1}, D) - \text{Cost}(\text{tree2}, D) = \log m^2 3^3 n^7 - \log m^4 3^5 n^4 = \log m^{-2} 3^{-2} n^3 = 3 \log n - 2 \log 3 - 2 \log m$$

∴ ควรใช้ tree2 เมื่อ  $2n \geq 3 + m$  และควรใช้ tree1 เมื่อ  $2n \leq 3 + m$

5. แบ่งข้อมูลออกเป็น 10 ส่วน แล้วทำการหาความแม่นยำ 10 รอบ โดยแต่ละรอบจะใช้ test set 1 ส่วนที่ไม่ซ้ำกัน

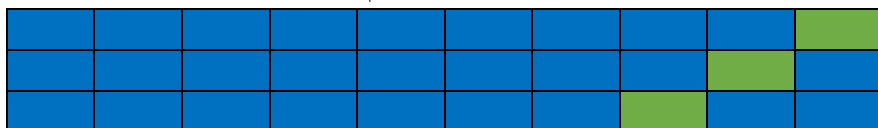


Train set



Test set

ค่าความแม่นยำจะใช้ค่าเฉลี่ยจากทั้ง 10 รอบ และ train set (กล่องสีฟ้า) จะมีการคัดเลือก model จากการแบ่งเป็น 10 ส่วน โดยใช้ 1 ส่วนที่ไม่ซ้ำกัน 10 รอบ เป็น validation set เพื่อหาค่า error และเลือก model ที่มีค่า error รวมต่ำที่สุด



Train set



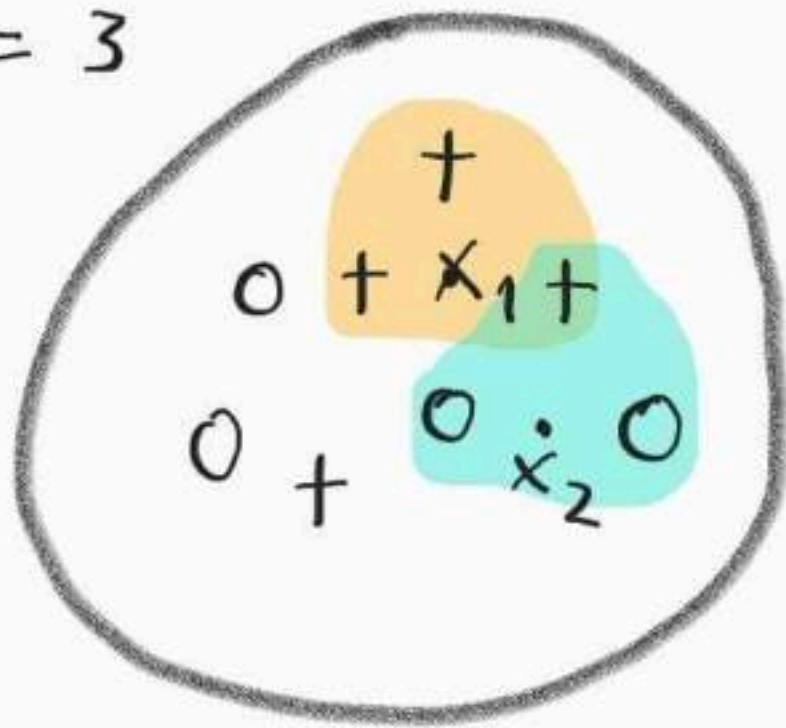
validation set



# K-NN แบ่งกลุ่ม

## 1. Majority voting แบ่งแบบลงคะแนนโดยใครเยอะกว่า

$k=3$



• ดูรอบ  $x_1$  3 ตัว  
จะเห็นว่า + อยู่รอบ  $x_1$  3 ตัว  
กลุ่ม  $x_1$  เป็น +

\* ถ้า  $x_1$  มี 2 ตัว + 2 ตัว  
จะแบ่งกลุ่มไม่ได้

• ดูรอบ  $x_2$  3 ตัว  
จะเห็นว่า o อยู่รอบ  $x_2$  2 ตัว กับ + 1 ตัว  
กลุ่ม  $x_2$  เป็น o

## 2. Weighted distance พิจารณาความถี่ในกลุ่มใกล้เคียงจุดที่กำหนด

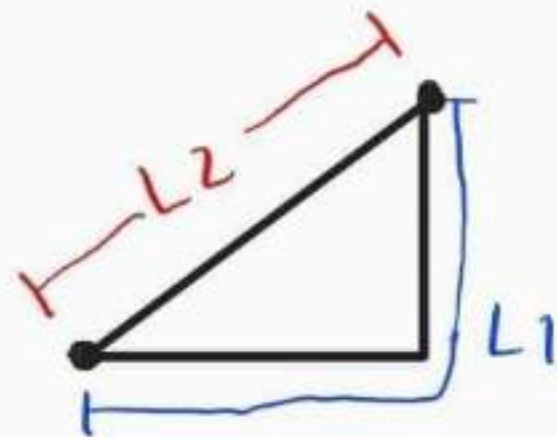
มี 3 แบบ คือ Euclidean, Manhattan และ Cosine

### • Euclidean distance

$$L2 = \left[ \sum_{i=1}^n (x_i - y_i)^2 \right]^{\frac{1}{2}}$$

### • Manhattan

$$L1 = \sum_{i=1}^n |x_i - y_i|$$



แล้วแต่จะเลือกใช้

### • Cosine Similarity เหมาะกับเอกสาร, text, search engine วัดความคล้ายกัน

$$\text{dist}(x, y) = \cos(x, y)$$

$$= \frac{x \cdot y}{\|x\| \|y\|}$$

การกำหนดค่า K ต้องลองทดลองเองว่า K ที่เท่าไรดีที่สุด

Gini Index ,Entropy และ Misclassification Error ต่างกันอย่างไร

Gini Index ,Entropy และ Misclassification Error เป็นการหาความบริสุทธิ์ของข้อมูลโดยจะมีผลลัพธ์ที่ใกล้เคียงกันคือ มีค่าเป็น 0 เมื่อข้อมูลมีเพียงคลาสเดียว และมีความมากที่สุดเมื่อ คลาสทุกคลาสมีปริมาณที่เท่าๆ กัน โดยค่าที่ได้จากทั้ง 3 แบบสามารถเรียงค่าจากผลลัพธ์ ดังนี้  $Entropy \geq Gini\ Index \geq Misclassification\ Error$