# Tree-based Methods

Petchara Pattarakijwanich

Introduction to Data Science, 4 November 2022

# Goal of this week

- Binary Tree
  - Regression Tree
  - Classification Tree
- Purity Metrics
  - Classification Error Rate
  - Gini Index
  - Entropy
- Methods to Improve Binary Tree
  - Bagging
  - Random Forest
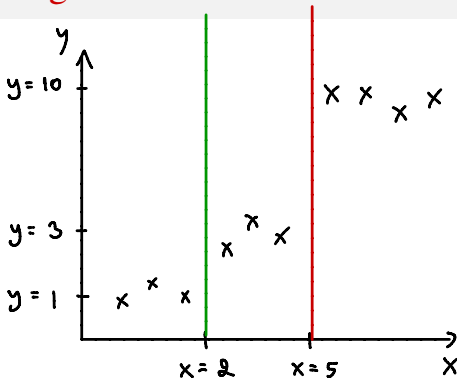  - Boosting

# About the Final Project

Approaches

- Answer specific questions from data, using data science techniques
- Implement, invent, or test some data science techniques or algorithms

Data Source
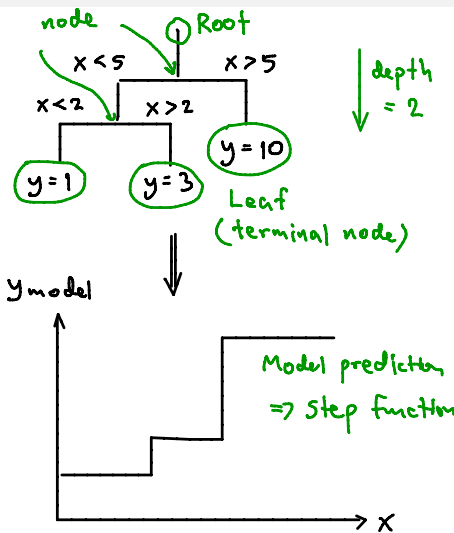
- Your own research
- Publicly available datasets
  https://geekflare.com/open-datasets-for-data-science/
  https://www.dataquest.io/blog/free-datasets-for-projects/
- Simulated data (?)

Presentation (some time around final week or a bit later)

# Regression Tree

# Regression Tree   กรณี 2 มิติ



$x_2$
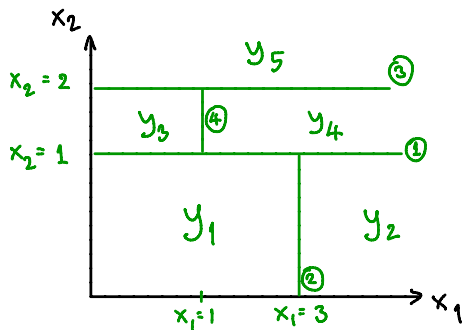
$x_2 = 2$

$x_2 = 1$

$y_5$  ③

$y_3$  ④   $y_4$   ①

$y_1$        $y_2$   ②

$x_1 = 1$   $x_1 = 3$   $x_1$

$x_2 < 1$        $x_2 > 1$

$x_1 < 3$   $x_1 > 3$   $x_2 < 2$   $x_2 > 2$

$x_1 < 1$   $x_1 > 1$

$y_1$   $y_2$   $y_3$   $y_4$   $y_5$

$y_i =$ ค่าเฉลี่ยของจุดข้อมูล
ในช่องนั้นๆ

Model Prediction
Step function ในกรณีนี้

# Classification Tree



ตัดที่ตัวแปรไหน & ค่าไหน ? ตัดยังไงถึง "ดี"

Purity metric
- Error Rate  ← ไม่ค่อยใช้
- Gini Index  } ที่นิยม
- Entropy

# How to Grow Trees

Input $\{x_1, x_2, \dots, x_n\}$  output $\{y\}$

- ตัวเลข → regression
- กลุ่ม → classification

① เริ่มจาก Root

② หาว่าตัด ที่ตัวแปรไหน & ค่าเท่าไหร่

- ลองตัด $x_1$ → หาค่าที่ตัดดีที่สุด
- ลองตัด $x_2$ → หาค่าที่ตัดดีที่สุด
- ⋮
- รูปแบบตัดที่ $x_i$ ที่ค่า $x_i$ = ?

ทำซ้ำ
จนพอใจ

ดีสุด
- Regression ⇒ MSE ลดลง มากสุด
- Classification
⇒ Purity เพิ่มขึ้น มากสุด

③ หยุดเมื่อไหร่?
- ไม่หยุดเลย ⇒ หนึ่ง node มี 1 จุดข้อมูล ⇒ over fit
- หยุดเมื่อความ "ดี" เพิ่มไม่มาก
- หยุดที่ depth ที่ต้องการ

# Purity Metrics



แดง 10

เขียว 5

น้ำเงิน 5

วัด ความ "pure" ของ ก้อน นี้ ยังไง?

- Error Rate

- Gini Index

- Shannon Entropy

node นึง ของ tree

# Classification Error Rate



$$E = 1 - P_{k,max}$$

ทำนายว่าทุกอย่าง เป็น mode

⇒ ทบผิด ณ เปอร์เซ็นต์

แดง 10
เขียว 5     ⇒
ดำน้ำ 5

$P_{แดง} = \frac{1}{2}$

$P_{เขียว} = \frac{1}{4}$

$P_{ดำน้ำ} = \frac{1}{4}$

ทายว่า ทุกอย่าง = แดง

⇒ ทบถูก = $\frac{10}{20}$

ทบผิด = $\frac{10}{20}$ = **0.5**

↑
Error Rate

$P_k = \begin{cases} P_{แดง} \\ P_{เขียว} \\ P_{ดำน้ำ} \end{cases}$

$P_{k,max} = P_{แดง} = 0.5$

$E = 1 - P_{k,max} =$ **0.5**

# Gini Index



$$G = \sum_k P_k(1-P_k) = 1 - \sum_k P_k^2$$

รูปลูก 1 จุด + รูปรวม 1 รูป

⇒ โอกาสผิด เวร ที่ เปอร์เซ็นต์

$$P_{แดง} = \frac{1}{2}$$

$$P_{เขียว} = \frac{1}{4}$$

$$P_{น้ำเงิน} = \frac{1}{4}$$

$$\begin{aligned} G = \ & P_{แดง}(1-P_{แดง}) \\ & + P_{เขียว}(1-P_{เขียว}) \\ & + P_{น้ำเงิน}(1-P_{น้ำเงิน}) \end{aligned}$$

prob หยิบ ได้
สีแดง ๆ

prob หยิบ ได้
label ครบ ๓ สีย

# Gini Index

$$G = P_{เมฆ}(1 - P_{เมฆ}) + P_{ฝน}(1 - P_{ฝน}) + P_{ดีเยี่ยม}(1 - P_{ดีเยี่ยม})$$

$$= P_{เมฆ} - P_{เมฆ}^2 + P_{ฝน} - P_{ฝน}^2 + P_{ดีเยี่ยม} - P_{ดีเยี่ยม}^2$$

$$= 1 - P_{เมฆ}^2 - P_{ฝน}^2 - P_{ดีเยี่ยม}^2$$

$$G = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = \ldots$$

# Shannon Entropy

$$D = -\sum_k P_k \log_2 P_k$$

"Information" ของ Event ที่มี prob $p$

$$I = \log_2\left(\frac{1}{p}\right) = -\log_2 p$$
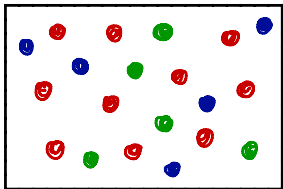
"Information" $p$ ต่ำ $\to$ $I$ สูง

$P_{แดง} = \frac{1}{2}$

$P_{เขียว} = \frac{1}{4}$

$P_{น้ำเงิน} = \frac{1}{4}$

$$I_A + I_B = \log_2\left(\frac{1}{P_A}\right) + \log_2\left(\frac{1}{P_B}\right)$$

Information ของ 2 event

$$= \log_2\left(\frac{1}{P_A P_B}\right)$$

$$= \log_2\left(\frac{1}{P(A \cup B)}\right)$$

# Shannon Entropy



$P_{red} = \frac{1}{2}$

$P_{blue} = \frac{1}{4}$

$P_{green} = \frac{1}{4}$

$$\text{Entropy} = \text{Expected value of Information}$$

$$= \langle I \rangle$$

$$= \sum_k P_k I_k$$
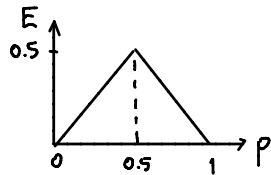
$$= \sum_k P_k \log_2 \left( \frac{1}{P_k} \right)$$

$$= - \sum_k P_k \log_2 P_k$$

$$D = - P_{red} \log_2 P_{red} - P_{blue} \log_2 P_{blue} - P_{green} \log_2 P_{green} = \ldots$$
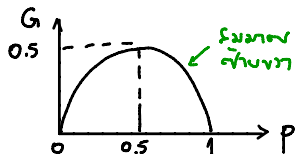
# Purity Metrics

สมมติมี 2 class
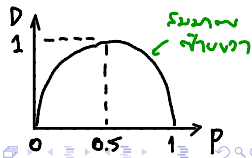$$\begin{cases} + & , \text{ prob } p \\ - & , \text{ prob } 1-p \end{cases}$$

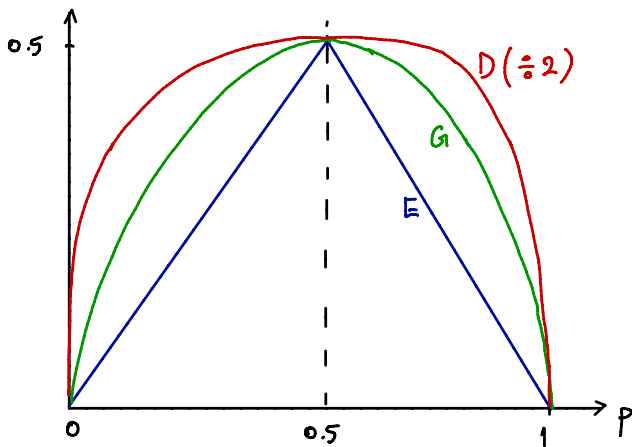$$E = 1 - P_{k,max} = \begin{cases} 1-p, & p \geq 0.5 \\ p, & p < 0.5 \end{cases}$$



$$G = 1 - \sum_k P_k^2 = 1 - p^2 - (1-p)^2$$



สมมาตร สำหรับ

$$D = -\sum_k P_k \log_2 P_k = -p \log_2 p - (1-p) \log_2 (1-p)$$



สมมาตร สำหรับ

# Purity Metrics

# Information Gain

$I_p, N_p$

$\left(\text{Parent}\right)$

$$\Delta I = I_p - \frac{N_L}{N_p} I_L - \frac{N_R}{N_p} I_R$$

ค่าเฉลี่ย ถ่วงน้ำหนัก ของ $I_L, I_R$

$I_L, N_L$

Left child

$I_R, N_R$

Right child

$\Delta I > 0 \Rightarrow$ แปงแล้ว ปรุริตี น้อยลง

$\Rightarrow$ Purity มากขึ้น

$\Delta I$ ยิ่งมากยิ่งดี

แปงด้วยแปรใหน ? ด้านที่ไหน? $\Rightarrow$ ใช้ $\Delta I$ มากสุด

# Information Gain



ตัวอย่างบวก

$x_0$ ตัวอย่างนี้คือ มาไหน?

Parent $\quad N_P = 9 \quad , \quad P_{บวก} = \dfrac{7}{9} \quad , \quad P_{ลบ} = \dfrac{2}{9}$

$$G_P = 1 - \sum_k P_k^2 = 1 - P_{บวก}^2 - P_{ลบ}^2 = 1 - \left(\dfrac{7}{9}\right)^2 - \left(\dfrac{2}{9}\right)^2 = 0.35$$

Left child $\quad N_L = 4 \quad , \quad P_{บวก} = 1 \quad , \quad P_{ลบ} = 0$

$$G_L = 1 - P_{บวก}^2 - P_{ลบ}^2 = 0$$

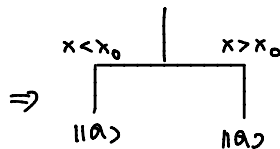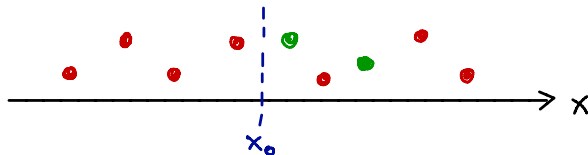Right child $\quad N_R = 5 \quad , \quad P_{บวก} = \dfrac{3}{5} \quad , \quad P_{ลบ} = \dfrac{2}{5}$

$$G_R = 1 - P_{บวก}^2 - P_{ลบ}^2 = 1 - \left(\dfrac{3}{5}\right)^2 - \left(\dfrac{2}{5}\right)^2 = 0.48$$

# Information Gain

$$\Delta G = G_p - \frac{N_L}{N_p} G_L - \frac{N_R}{N_p} G_R$$

$$= 0.35 - \frac{4}{9} \times 0 - \frac{5}{9} \times 0.48$$

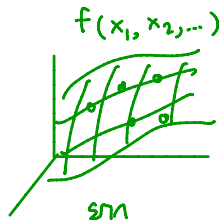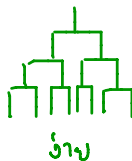$$= 0.08 > 0 \quad \Rightarrow \quad \text{แบ่ง ดีกว่าไม่แบ่ง}$$



$\Rightarrow$

```
      x<x₀  |  x>x₀
        |___|___|
        |       |
      แดง      เกว
```

แปงทำไม
- เพิ่ม node purity
- แปง เริ่ม node ของ
  ตัวใต้

# Pros and Cons of Classification Tree

ข้อดี

- ใช้กับ Qualitative variable ง่าย

- เข้าใจง่าย อธิบายง่าย

  (ใกล้เคียงกับ มนุษย์ตัด
  สินใจ หรือ ครับ?)

$f(x_1, x_2, ...)$



ง่าย

ยาก

ข้อเสีย

- Variance สูง ⎫ overfit ง่าย
- Prediction Accuracy ต่ำ ⎭ - แก้โดย

  Bagging
  Random Forest
  Boosting