



**การบ้านปฏิบัติการ11**  
**Artificial Intelligence (ID3)(20 คะแนน)**

- 1) **20คะแนน**(HW11\_5XXXXXXX.py)ให้เขียนโปรแกรมเพื่ออ่านข้อมูลตัวอย่าง (training data) จากไฟล์ input.txt จากนั้นนำข้อมูลที่ได้มาทำการสร้าง Decision tree โดยใช้วิธีคำนวณตามหลักการ ID3 และแสดงผลลัพธ์ดังตัวอย่างที่กำหนดให้

**ตัวอย่างข้อมูล training data**

Film	Country of origin	Big Star	Genre	Success
1	USA	yes	Science	true
2	USA	no	Comedy	false
3	USA	yes	Comedy	true
4	Europe	no	Comedy	true
5	Europe	yes	Science	false
6	Europe	yes	Romance	false
7	Rest of world	yes	Comedy	false
8	Rest of world	no	Science	false
9	Europe	yes	Comedy	true
10	USA	yes	Comedy	true

training dataอยู่ในไฟล์โดยมีรูปแบบดังนี้

1;USA;yes;Science;true  
2;USA;no;Comedy;false  
3;USA;yes;Comedy;true

**ผลลัพธ์**ให้แสดงดังนี้

First factor is Country of origin  
The next factor of USA is Big star  
The next factor of Europe is Genre

- ให้เขียนฟังก์ชัน `information_gain(feature)` เพื่อคืนค่า `information gain` ของ feature แต่ละตัวของข้อมูลตามหลักการของ ID3 โดยสามารถมีการเรียกฟังก์ชันย่อยที่เหมาะสมได้
- ให้เขียนฟังก์ชัน `factor_list()` เพื่อแสดงผลลำดับของปัจจัยของ Decision Tree ที่ได้ ตามตัวอย่างด้านบน

หลักการคำนวณโดยใช้ ID3 มีรายละเอียดดังนี้

$$\text{Entropy} = E(S) = -p_1 \log_2 p_1 - p_0 \log_2 p_0$$

จากข้อมูล training data จะหาปัจจัยตัวแรกได้ดังนี้

กรณีปัจจัย **Country of origin** จะหา Entropy ของค่าตัวแปร 3 ค่าได้ดังนี้

$$\begin{aligned} E(\text{USA}) &= -(3/4) \log_2 (3/4) - (1/4) \log_2 (1/4) \\ &= 0.311 + 0.5 = 0.811 \end{aligned}$$

$$E(\text{Europe}) = -(2/4) \log_2 (2/4) - (2/4) \log_2 (2/4) = 1$$

$$E(\text{Rest of world}) = -0 - (2/2) \log_2 (2/2) = 0$$

จากนั้นนำมาคำนวณหาค่า Information gain ของ Country of origin

$$\begin{aligned} \text{Country of origin} &= 1 - P(\text{USA}) \times E(\text{USA}) - P(\text{Europe}) \times E(\text{Europe}) - P(\text{Rest of world}) \times E(\text{Rest of world}) \\ &= 1 - (0.4 \times 0.811) - (0.4 \times 1) - (0.2 \times 0) \\ &= 0.2756 \end{aligned}$$

กรณีปัจจัย **Big Star** จะหา Entropy ของค่าตัวแปร 2 ค่าได้ดังนี้

$$\begin{aligned} E(\text{yes}) &= -(4/7) \log_2 (4/7) - (3/7) \log_2 (3/7) \\ &= 0.9852 \end{aligned}$$

$$E(\text{no}) = -(1/3) \log_2 (1/3) - (2/3) \log_2 (2/3) = 0.9151$$

จากนั้นนำมาคำนวณหาค่า Information gain ของ Big Star

$$\begin{aligned} \text{Big Star} &= 1 - P(\text{yes}) \times E(\text{yes}) - P(\text{no}) \times E(\text{no}) \\ &= 1 - (0.7 \times 0.9852) - (0.3 \times 0.9151) \\ &= 0.044 \end{aligned}$$

กรณีปัจจัย **Genre** จะหา Entropy ของค่าตัวแปร 3 ค่าได้ดังนี้

$$E(\text{science}) = -(1/3) \log_2 (1/3) - (2/3) \log_2 (2/3) = 0.9151$$

$$E(\text{comedy}) = -(4/6) \log_2 (4/6) - (2/6) \log_2 (2/6) = 0.9151$$

$$E(\text{romance}) = 0 - 1 \log_2 (1) = 0$$

จากนั้นนำมาคำนวณหาค่า Information gain ของ Genre

$$\begin{aligned} \text{Genre} &= 1 - (0.6 \times 0.9151) - (0.3 \times 0.9151) - (0.1 \times 0) \\ &= 0.176 \end{aligned}$$

ดังนั้น ปัจจัยตัวแรกของ Decision Tree คือ Country of origin และได้ tree ดังรูป



ขั้นถัดมา หาปัจจัยตัวที่สองที่ต่อจาก USA พิจารณาเฉพาะข้อมูลที่มี Country เป็น USA

Film	Country of origin	Big Star	Genre	Success
1	USA	yes	Science	true
2	USA	no	Comedy	false
3	USA	yes	Comedy	true
4	USA	yes	Comedy	true

กรณีปัจจัย **Big Star** จะหา Entropy ของค่าตัวแปร 2 ค่าได้ดังนี้

$$E(\text{yes}) = -1 \log_2(1) = 0$$

$$E(\text{no}) = 0 - 1 \log_2(1) = 0$$

จากนั้นนำมาคำนวณหาค่า Information gain ของ Big Star

$$\begin{aligned} \text{Big Star} &= 0.811 - P(\text{yes}) \times E(\text{yes}) - P(\text{no}) \times E(\text{no}) \\ &= 0.811 - (0.75 \times 0) - (0.25 \times 0) = 0.811 \end{aligned}$$

คิดภายใต้เงื่อนไขของ USA เลยต้องใช้ E(USA)

กรณีปัจจัย **Genre** จะหา Entropy ของค่าตัวแปร 3 ค่าได้ดังนี้

$$E(\text{science}) = -1 \log_2(1) = 0$$

$$E(\text{comedy}) = -(2/3) \log_2(2/3) - (1/3) \log_2(1/3) = 0.9151$$

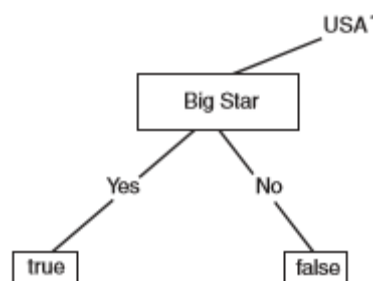
$$E(\text{romance}) = 0 \text{ เพราะไม่มีค่าในส่วนนี้}$$

จากนั้นนำมาคำนวณหาค่า Information gain ของ Genre

$$\begin{aligned} \text{Genre} &= 0.811 - (0.25 \times 0) - (0.75 \times 0.9151) - 0 \\ &= 0.125 \end{aligned}$$

คิดภายใต้เงื่อนไขของ USA เลยต้องใช้ E(USA)

ดังนั้น ปัจจัยตัวถัดมาของ Decision Tree ที่ต่อจาก USA คือ Big Star และได้ tree ดังรูป



ขั้นถัดมา หาปัจจัยตัวที่สองที่ต่อจาก Europe พิจารณาเฉพาะข้อมูลที่มี Country เป็น Europe

Film	Country of origin	Big Star	Genre	Success
1	Europe	no	Comedy	true
2	Europe	yes	Science	false
3	Europe	yes	Romance	false
4	Europe	yes	Comedy	true

กรณีปัจจัย **Big Star** จะหา Entropy ของค่าตัวแปร 2 ค่าได้ดังนี้

$$E(\text{yes}) = - (1/3) \log_2 (1/3) - (2/3) \log_2 (2/3) = 0.9151$$

$$E(\text{no}) = - 1 \log_2 (1) = 0$$

จากนั้นนำมาคำนวณหาค่า Information gain ของ Big Star

$$\text{Big Star} = 1 - P(\text{yes}) \times E(\text{yes}) - P(\text{no}) \times E(\text{no})$$

$$= 1 - (0.75 \times 0.9151) - (0.25 \times 0) = 0.314$$

คิดภายใต้เงื่อนไขของ Europe เลยต้องใช้  
E(Europe)

กรณีปัจจัย **Genre** จะหา Entropy ของค่าตัวแปร 3 ค่าได้ดังนี้

$$E(\text{science}) = 0$$

$$E(\text{comedy}) = 0$$

$$E(\text{romance}) = 0$$

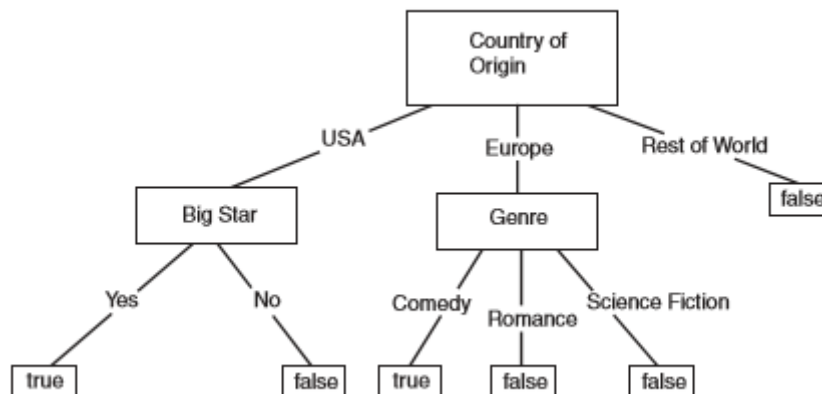
จากนั้นนำมาคำนวณหาค่า Information gain ของ Genre

$$\text{Genre} = 1 - 0$$

$$= 1$$

คิดภายใต้เงื่อนไขของ Europe เลยต้องใช้  
E(Europe)

ดังนั้น ปัจจัยตัวถัดมาของ Decision Tree ที่ต่อจาก Europe คือ Genre และได้ tree ดังรูป



### การส่งงาน

1. ลักษณะ/ลำดับข้อความของการรับค่า/แสดงผลจะต้องเป็นไปตามที่ระบุในตัวอย่างการ run
2. ไฟล์งาน/ใบงานที่ส่ง จะต้องมีการแทรก comment/หัวกระดาษ ตามข้อกำหนดใน website รายวิชา
3. ไฟล์งานโปรแกรมที่ส่ง จะต้องมีการแทรก pseudocode เป็น comment ในแต่ละขั้นตอน
4. Upload ไฟล์ source code ตามที่ระบุในแต่ละข้อ ไปยัง website ที่ใช้ส่งการบ้าน

<http://hw.cs.science.cmu.ac.th> ตาม section ที่นักศึกษาเรียน