

Assignment 6

Training data			Test data		
Doc#	Content	Category	Doc#	Content	Category
D1	He loves her badly	Novel	D7	He walks with her	?
D2	I slept	Diary	D8	I fight a dog	?
D3	A dog walks on a highway	Article			
D4	He slept	Novel			
D5	I and a dog walk on a highway	Diary			
D6	He fights near a highway	Article			

*คำที่ควรจะมีรากศัพท์เดียวกัน เช่น walk และ walks จะนับเป็นคำเดียวกัน

(3 คะแนน) 1 ค่าความน่าจะเป็นของ Category แต่ละตัวมีขนาดเท่าไร

(2 คะแนน) 2 ดูจาก Training set แล้วตอบคำถามต่อไปนี้

2.1 จำนวนคำของทุกเอกสารที่ถูกจัดในแต่ละ Category (รวมตัวซ้ำด้วย)

2.2 จำนวนคำของทุกเอกสารของทุก Category (ไม่รวมตัวซ้ำ)

(3 คะแนน) 3 ค่าความน่าจะเป็นของแต่ละคำทั้งหมดโดยอิงจาก Category มีขนาดเท่าไร

(2 คะแนน) 4 จงหาว่าเอกสาร D7 D8 ควรถูกจัดให้อยู่ในกลุ่มใด

1. คำความน่าจะเป็นของ Category แต่ละตัวขนาดเท่าไร

Ans $P(\text{Novel}) = \frac{2}{6} = \frac{1}{3}$

$$P(\text{Diary}) = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{Article}) = \frac{2}{6} = \frac{1}{3}$$

2. ดูจาก Training set แล้วตอบคำถามต่อไปนี้

2.1 จำนวนคำของทุกเอกสารที่ถูกจัดใน แต่ละ Category (รวมตัวซ้ำด้วย)

Ans คำทั้งหมดใน Novel = 6

คำทั้งหมดใน Diary = 10

คำทั้งหมดใน Article = 11

2.2 จำนวนคำของทุกเอกสารของทุก Category (ไม่รวมตัวซ้ำ)

Ans มีทั้งหมด 14 คำ

3. คำความน่าจะเป็นของ แต่ละคำทั้งหมดโดยอิงจาก Category มีขนาดเท่าไร

Ans ใน Novel

$$P(\text{He} | \text{Novel}) = 0.15$$

$$P(\text{her} | \text{Novel}) = 0.1$$

$$P(\text{loves} | \text{Novel}) = 0.1$$

$$P(\text{badly} | \text{Novel}) = 0.1$$

$$P(\text{Slept} | \text{Novel}) = 0.1$$

$$P(\text{I} | \text{Novel}) = 0.05$$

$$P(\text{A} | \text{Novel}) = 0.05$$

$$P(\text{dog} | \text{Novel}) = 0.05$$

$$P(\text{walk} | \text{Novel}) = 0.05$$

$$P(\text{and} | \text{Novel}) = 0.05$$

$$P(\text{on} | \text{Novel}) = 0.05$$

$$P(\text{highway} | \text{Novel}) = 0.05$$

$$P(\text{fight} | \text{Novel}) = 0.05$$

$$P(\text{near} | \text{Novel}) = 0.05$$

Q6 Diary

$$P(\text{He}|\text{Diary}) = 0.05$$

$$P(\text{loves}|\text{Diary}) = 0.05$$

$$P(\text{Slept}|\text{Diary}) = 0.1$$

$$P(\text{A}|\text{Diary}) = 0.15$$

$$P(\text{walk}|\text{Diary}) = 0.1$$

$$P(\text{on}|\text{Diary}) = 0.1$$

$$P(\text{fight}|\text{Diary}) = 0.05$$

$$P(\text{her}|\text{Diary}) = 0.05$$

$$P(\text{badly}|\text{Diary}) = 0.05$$

$$P(\text{I}|\text{Diary}) = 0.15$$

$$P(\text{dog}|\text{Diary}) = 0.1$$

$$P(\text{and}|\text{Diary}) = 0.1$$

$$P(\text{highway}|\text{Diary}) = 0.1$$

$$P(\text{near}|\text{Diary}) = 0.05$$

Q6 Article

$$P(\text{He}|\text{Article}) = 0.1$$

$$P(\text{loves}|\text{Article}) = 0.05$$

$$P(\text{Slept}|\text{Article}) = 0.05$$

$$P(\text{A}|\text{Article}) = 0.2$$

$$P(\text{walk}|\text{Article}) = 0.1$$

$$P(\text{on}|\text{Article}) = 0.1$$

$$P(\text{fight}|\text{Article}) = 0.1$$

$$P(\text{her}|\text{Article}) = 0.05$$

$$P(\text{badly}|\text{Article}) = 0.05$$

$$P(\text{I}|\text{Article}) = 0.05$$

$$P(\text{dog}|\text{Article}) = 0.1$$

$$P(\text{and}|\text{Article}) = 0.05$$

$$P(\text{highway}|\text{Article}) = 0.15$$

$$P(\text{near}|\text{Article}) = 0.1$$

4. จงหาว่าเอกสาร D7 D8 ควรถูกจัดให้อยู่ในกลุ่มใด

Ans (D7) $P(\text{Novel} | D7) = P(\text{He} | \text{Novel}) \cdot P(\text{walks} | \text{Novel}) \cdot P(\text{with} | \text{Novel}) \cdot P(\text{her} | \text{Novel}) \cdot P(\text{Novel})$

$$= 0.15 \times 0.05 \times 0.05 \times 0.1 \times \frac{1}{3}$$

$$= 0.0000125$$

$$P(\text{Diary} | D7) = P(\text{He} | \text{Diary}) \cdot P(\text{walks} | \text{Diary}) \cdot P(\text{with} | \text{Diary}) \cdot P(\text{her} | \text{Diary}) \cdot P(\text{Diary})$$

$$= 0.05 \times 0.1 \times 0.05 \times 0.05 \times \frac{1}{3}$$

$$= 0.00000416$$

$$P(\text{Article} | D7) = P(\text{He} | \text{Article}) \cdot P(\text{walks} | \text{Article}) \cdot P(\text{with} | \text{Article}) \cdot P(\text{her} | \text{Article}) \cdot P(\text{Article})$$

$$= 0.1 \times 0.1 \times 0.05 \times 0.05 \times \frac{1}{3}$$

$$= 0.00000833$$

\therefore D7 ควรถูกจัดใน Novel

(D8)

$$\begin{aligned}P(\text{Novel} | D8) &= P(I | \text{Novel}) \cdot P(\text{fight} | \text{Novel}) \cdot P(a | \text{Novel}) \cdot P(\text{dog} | \text{Novel}) \cdot P(\text{Novel}) \\&= 0.05 \times 0.05 \times 0.05 \times 0.05 \times \frac{1}{3} \\&= 0.00000208\end{aligned}$$

$$\begin{aligned}P(\text{Diary} | D8) &= P(I | \text{Diary}) \cdot P(\text{fight} | \text{Diary}) \cdot P(a | \text{Diary}) \cdot P(\text{dog} | \text{Diary}) \cdot P(\text{Diary}) \\&= 0.15 \times 0.05 \times 0.15 \times 0.1 \times \frac{1}{3} \\&= 0.0000375\end{aligned}$$

$$\begin{aligned}P(\text{Article} | D8) &= P(I | \text{Article}) \cdot P(\text{fight} | \text{Article}) \cdot P(a | \text{Article}) \cdot P(\text{dog} | \text{Article}) \cdot P(\text{Article}) \\&= 0.1 \times 0.1 \times 0.2 \times 0.1 \times \frac{1}{3} \\&= 0.0000667\end{aligned}$$

$\therefore D8$ ควรจัดใน Article