

Step 0

Answer

2.1 เลือกใช้ Extend Boolean Model เนื่องจากลักษณะของ Query เป็นแบบ Boolean และโจทย์กำหนดให้ Keyword ไม่สัมพันธ์กัน

เอกสาร 10 เอกสารมีการแจกแจง Keyword ดังนี้

D1: {bird, cat, bird, cat, dog, dog, bird}
D2: {cat, tiger, cat, dog}
D3: {dog, bird, bird}
D4: {cat, tiger}
D5: {tiger, tiger, dog, tiger, cat}
D6: {bird, cat, bird, cat, tiger, tiger, bird}
D7: {bird, tiger, cat, dog}
D8: {dog, cat, bird}
D9: {cat, dog, tiger}
D10: {tiger, tiger, tiger}

	Bird	Cat	Dog	Tiger	Max
Doc1	3	2	2	0	3
Doc2	0	2	1	1	2
Doc3	2	0	1	0	2
Doc4	0	1	0	1	1
Doc5	0	1	1	3	3
Doc6	3	2	0	2	3
Doc7	1	1	1	1	1
Doc8	1	1	1	0	1
Doc9	0	1	1	1	1
Doc10	0	0	0	3	3
n	5	8	7	7	

Query = รักแมวและสุนัข แต่ไม่รักเสือ

Query = (Cat AND Dog) AND NOT Tiger

Boolean Model

- เนื่องจาก Boolean Model มีข้อดีคือการไม่สนใจน้ำหนักของ Keyword
- Vector Space Model มีข้อดีคือการเชื่อมโยงความสัมพันธ์ระหว่าง Keyword

จึงใช้มีความเหมาะสมที่จะใช้ของทั้งสองมารวมกัน ทำให้เป็น Model ใหม่ขึ้นมา เรียกว่า **Extended Boolean Model**

Step 1

หา ท.ค. แต่ละ key ใน เอกสาร Doc

tf norm normalize
idf norm normalize

น้ำหนักของ Keyword ในเอกสาร

น้ำหนักของ Keyword "i" ในเอกสาร "j"

$$w_{ij} = tfnorm_{ij} \times idf_{ij}$$

tfnorm_{ij} = $\frac{tf_{ij}}{tfmax_{ij}}$

idfnorm_{ij} = $\frac{idf_{ij}}{idfmax_{ij}}$

normalized TF ของ Keyword "i" ในเอกสาร "j"

normalized IDF ของ Keyword "i" ในเอกสารทั้งหมด

Doc1

tf _{bird} = $\frac{3}{3} = 1.000$	* ทุกตัวให้ size ส่วนไหน, ช่วงไหนต่อ ให้ max = 1
tf _{cat} = $\frac{2}{2} = 0.667$	
tf _{dog} = $\frac{2}{2} = 0.667$	
tf _{tiger} = $\frac{0}{3} = 0.000$	

idf _{bird} = $\log(\frac{10}{3}) = 0.301$	idf _{norm, bird} = $\frac{0.301}{0.301} = 1.000$
idf _{cat} = $\log(\frac{10}{8}) = 0.097$	idf _{norm, cat} = $\frac{0.097}{0.301} = 0.322$
idf _{dog} = $\log(\frac{10}{7}) = 0.155$	idf _{norm, dog} = $\frac{0.155}{0.301} = 0.515$
idf _{tiger} = $\log(\frac{10}{7}) = 0.155$	idf _{norm, tiger} = $\frac{0.155}{0.301} = 0.515$

$$\begin{aligned} w_{bird} &= 1.000 \times 1.000 = 1.000 \\ w_{cat} &= 0.667 \times 0.322 = 0.215 \\ w_{dog} &= 0.667 \times 0.515 = 0.343 \\ w_{tiger} &= 0.000 \times 0.515 = 0.000 \end{aligned}$$

น้ำหนักของแต่ละ Keyword ในแต่ละเอกสาร

	Bird	Cat	Dog	Tiger
Doc1	1.000	0.215	0.343	0.000
Doc2	0.000	0.322	0.257	0.257
Doc3	1.000	0.000	0.257	0.000
Doc4	0.000	0.322	0.000	0.515
Doc5	0.000	0.107	0.172	0.515
Doc6	1.000	0.215	0.000	0.343
Doc7	1.000	0.322	0.515	0.515
Doc8	1.000	0.322	0.515	0.000
Doc9	0.000	0.322	0.515	0.515
Doc10	0.000	0.000	0.000	0.515

Step 2

หา ท.ค.

Query = รักแมวและสุนัข แต่ไม่รักเสือ

Query = (Cat AND Dog) AND NOT Tiger

น้ำหนักของแต่ละ Keyword ในแต่ละเอกสาร

	Bird	Cat	Dog	Tiger
Doc1	1.000	0.215	0.343	0.000
Doc2	0.000	0.322	0.257	0.257
Doc3	1.000	0.000	0.257	0.000
Doc4	0.000	0.322	0.000	0.515
Doc5	0.000	0.107	0.172	0.515
Doc6	1.000	0.215	0.000	0.343
Doc7	1.000	0.322	0.515	0.515
Doc8	1.000	0.322	0.515	0.000
Doc9	0.000	0.322	0.515	0.515
Doc10	0.000	0.000	0.000	0.515

AND

$$sim(q_{and}, d_j) = 1 - \sqrt{\frac{(1 - w_{cat,j})^2 + (1 - w_{dog,j})^2}{2}}$$

$$sim(q_{and}, d_j) = 1 - \sqrt{\frac{\left(1 - \left(1 - \sqrt{\frac{(1 - w_{cat,j})^2 + (1 - w_{dog,j})^2}{2}}\right)^2 + (1 - (1 - w_{tiger,j}))^2\right)}{2}}$$

$$sim(q_{and}, d_1) = 1 - \sqrt{\frac{\left(1 - \left(1 - \sqrt{\frac{(1 - 0.215)^2 + (1 - 0.343)^2}{2}}\right)^2 + (1 - (1 - 0.000))^2\right)}{2}}$$

$$sim(q_{and}, d_1) = 0.488$$

ทุก Doc

จัด 3 ตัว

$$\sqrt{\frac{(1 - w_{cat,j})^2 + (1 - w_{dog,j})^2 + (1 - w_{tiger,j})^2}{2}}$$

degree

	Sim
Doc1	0.488
Doc2	0.465
Doc3	0.377
Doc4	0.295
Doc5	0.291
Doc6	0.320
Doc7	0.447
Doc8	0.583
Doc9	0.447
Doc10	0.205

Step 3

จัด rank

Rank → Doc8, Doc1, Doc2, Doc7, Doc9, Doc3, Doc6, Doc4, Doc5, Doc10