

# Search Model

- Extended Boolean Model

## Contents

1. Definition
2. Boolean OR
3. Boolean AND
4. Normalized
5. P-Norm
6. AND OR 조합
7. Features

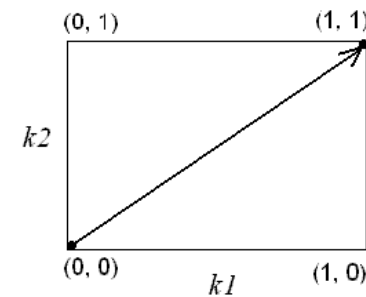
## Definitions

- is based on Boolean Model and added Vector Space Model.
- The weak point of Boolean Model is not give Weight.
- Query used Boolean Model and Results used Vector Space Model.
- find similarity how calculate Euclid distance between Document Term Weight

## Boolean OR

- Assume Query has K1, K2
- So, We can find two extreme points in Term Space.
  - Document include two terms(k1, k2), and complete similarity  $\rightarrow (1,1)$
  - Document hasn't two terms(k1, k2), and similarity is 0  $\rightarrow (0,0)$
- The max distance (dmax) of two point is 1.41
- Therefore, the similarity of documents is between 0 and 1.41

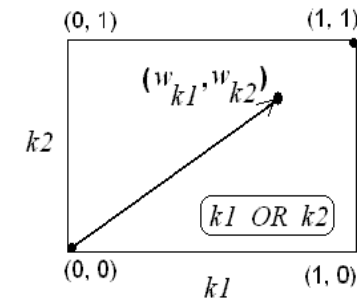
$$d_{\max} = \sqrt{(1-0)^2 + (1-0)^2} = \sqrt{2} = 1.41$$



- We know that Euclid distance of OR Query ( $W_{k1}, W_{k2}$ ) is less than 1.41
- $W_{k1} = TF_{(k1, d)} \times (IDF_{k1} / (\max_d \times IDF_d))$

$$d_{\text{OR}} = \sqrt{(w_{k1} - 0)^2 + (w_{k2} - 0)^2} = \sqrt{w_{k1}^2 + w_{k2}^2}$$

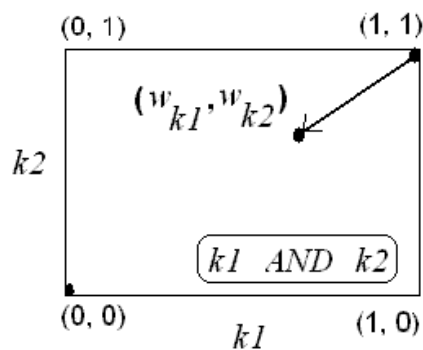
- $TF_{(k1, d)}$  is TF of k1 in document 'd'
- $\max_d$  is TF of term has max TF in document 'd'
- $IDF_d$  is inverse DF(IDF) of term has max TF in document 'd'
- $\max_d \times IDF_d$  is TF\*IDF of term has max TF in document 'd'



## Boolean AND

- Assume Query has K1, K2
- We calculate Euclid distance between  $d_{max}(1,1)$  and  $(W_{k1}, W_{k2})$

$$d_{AND} = \sqrt{2} - \sqrt{(1 - w_{k1})^2 + (1 - w_{k2})^2}$$



## Normalized : Similarity Score

- We need normalized for similarity
- So, divide into term's number

$$Sim(Q_{k1 \text{ OR } k2}, D) = \sqrt{\frac{w_{k1}^2 + w_{k2}^2}{2}}$$

$$Sim(Q_{k1 \text{ AND } k2}, D) = 1 - \sqrt{\frac{(1 - w_{k1})^2 + (1 - w_{k2})^2}{2}}$$

## P-Norm

- Add independent p-parameter(p-norm) into Normalized Similarity

$$Sim(Q_{k1 \text{ OR } k2}, D) = \left( \frac{w_{k1}^p + w_{k2}^p + \dots + w_{km}^p}{m} \right)^{1/p}$$

$$Sim(Q_{k1 \text{ AND } k2}, D) = 1 - \left( \frac{(1 - w_{k1})^p + (1 - w_{k2})^p + \dots + (1 - w_{km})^p}{m} \right)^{1/p}$$

- if p is 1, then get effectiveness of **Vector Space Model**

$$- Sim(Q_{\text{OR}}, D) = Sim(Q_{\text{AND}}, D) = (W_{k1} + W_{k2} + \dots + W_m) / m$$

- if p is  $\infty$ , then get effectiveness of **Boolean Model**

## AND OR 조합

- p-norm has between 1 and  $\infty$

$$Sim(Q_{k1 \text{ AND } k2 \text{ OR } k3}, D) = \left( \frac{\left( 1 - \left( \frac{(1 - w_{k1})^p + (1 - w_{k2})^p}{2} \right)^{1/p} \right)^p + w_{k3}^p}{2} \right)^{1/p}$$

- Upper Similarity's formula can apply using recursive method unrelated Operator's number



## Features

- We can change p-norm 1 to  $\infty$ , so can see Vector Ranking Result to Boolean Ranking Result
- We can use different p-norm
  - for example

$$Sim(Q_{k1 \text{ AND } k2 \text{ OR } k3}, D) = \left( \frac{\left( 1 - \left( \frac{(1 - w_{k1})^p + (1 - w_{k2})^p}{2} \right)^{1/p} + w_{k3} \right)^p}{2} \right)^{1/p}$$

$P = 2$ 
 $P = \infty$

so, can get Vector Space Model's effectiveness at specified Operator ( $p=2$ )  
 also, get Boolean Model's effectiveness at specified Operator ( $p = \infty$ )