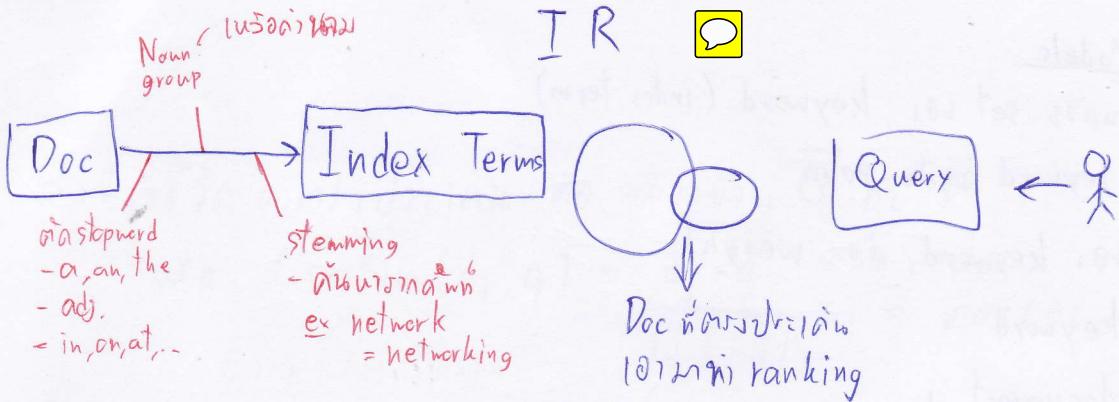
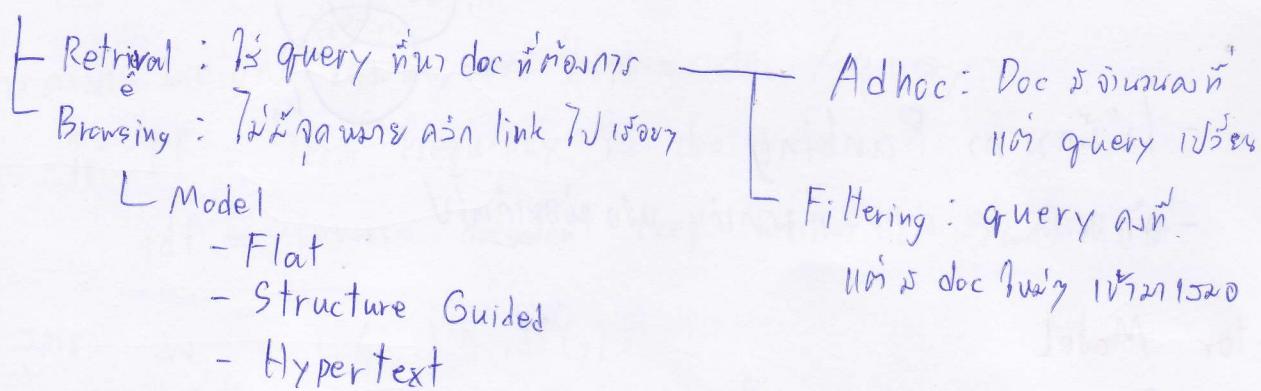


①



- จุดเด่นของ IR, ใช้กันอย่างแพร่หลาย, การจัดการดับเบล relevance (ความซ้ำซ้อน) รำคาญที่สุด
- หา keyword ของหนึ่ง inverted file ของ keyword นั้นๆ ใน Doc ที่เก็บไว้
- Ranking ต้องการเรียงลำดับ Doc ตามความสำคัญ
- User Task



### IR Model ที่ใช้

- Classic Model (บันทึกความถี่ของคำในหัวข้อ)
- boolean
- vector
- probabilistic
- fuzzy
- Structured Model (บันทึกหัวข้อในหน้า, หัวร่อง, หัวเรื่อง)
- Non-Overlapping Lists
- Proximal Nodes

②

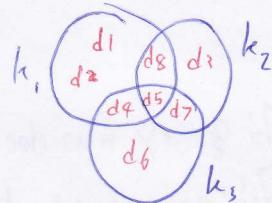
## • Classic IR Models

- Doc contains set var keyword (index term)
- keyword var with noun
- assign var keyword as weight
- $k_i$  = keyword i
- $d_j$  = document j
- $w_{ij}$  = weight var keyword i for document j

## Boolean Model

- no weight like 0 or 1
  - for keyword AND( $\wedge$ ) OR( $\vee$ ) NOT( $\neg$ ) in Query
- $$q = k_a \wedge (k_b \vee \neg k_c) \quad * \text{boolean expression}$$

- AND operation is intersection of documents



- Ranking

- ranking doc according to their weight

## Vector Model

- has w variables = keyword frequencies
- in query
- for doc w.r.t query var the most w.r.t vector in

	$k_1$	$k_2$	$k_3$
$d_1$	2	0	1
$d_2$	1	0	0
$q$	1	2	3

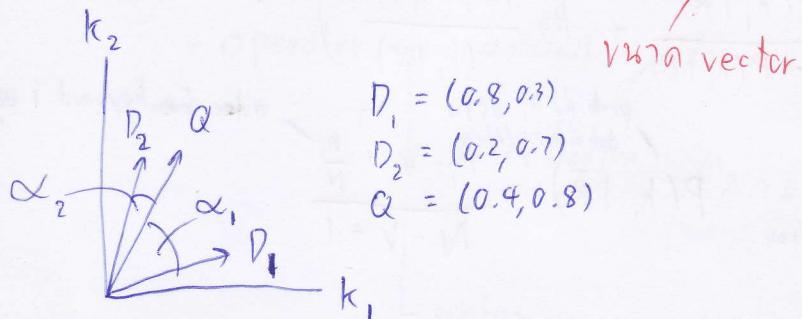
$\Rightarrow$

$$\begin{aligned}\vec{d}_1 &= (2, 0, 1) \\ \vec{d}_2 &= (1, 0, 0) \\ \vec{q} &= (1, 2, 3)\end{aligned}$$

(3)

- ចំណាំបង្កើតអាមេរិកាន់ គោលរូបរាង  $\vec{Q}$  និង  $\vec{D}$

$$\text{ដូច} \quad \text{CosSim}(d_j, q) = \frac{\vec{d}_j \cdot \vec{q}}{|\vec{d}_j| |\vec{q}|} \quad \text{នៅរវយក} \\ \text{សំណាត់សំព័រជាការ} \quad = \cos(\theta)$$



$$\text{sim}(Q, D_1) = \frac{0.8(0.4) + 0.3(0.8)}{(\sqrt{0.8^2 + 0.3^2})(\sqrt{0.4^2 + 0.8^2})} \\ = 0.74$$

- គោលរូបរាង weight រវយក keyword  $\theta_{\text{sim}} = \text{doc} / \text{query}$

-  $tf$  = term frequency នៃ doc ឱ doc

$idf$  = inverse document freq. គោលរូប doc និង keyword

$$\text{ដូច} \quad w_{ij} = tf(i, j) \times idf(i)$$

keyword  $\downarrow$  doc \*  $\uparrow$  រួមទាំង  $\downarrow$  keyword  $i$   $\uparrow$  រួមទាំង  $j$   $\uparrow$  # doc និង keyword  $i$

$$tf(i, j) = \frac{\text{freq}(i, j)}{\max(\text{freq}(i, j))} \quad idf(i) = \log\left(\frac{N}{n_i}\right)$$

# doc និង keyword  $i$  នៅ

- គោលរូប weight រវយក query

$$w_{i,q} = 0.5 + \frac{0.5 \times \text{freq}(i, q)}{\max(\text{freq}(i, q))} \times \log\left(\frac{N}{n_i}\right)$$

$$w_{\text{door}, q} = 0.5 + \frac{0.5 \times 2}{2} \times \log\left(\frac{6}{2}\right)$$

Ex Query = "Visitor at your door or my door"

key terms = {"visitor", "door", "door"}

$$\text{freq}(\text{visitor}, q) = 1, \text{ freq}(\text{door}, q) = 2$$

4

## • Probabilistic Model

- ប្រាក់ទី 1 ដែលមានការងារសម្រាប់រំលែក និងអាជីវកម្ម
  - ក្រុមហ៊ុនសំខាន់សំខាន់ ដែលមានការងារសម្រាប់រំលែក និងអាជីវកម្ម / ជាសម្រាប់
  - សាស្ត្រ Fuhr Model

$$\text{Sim}(d_j, q) = \sum_{i=1}^t w_{i,q} \times w_{i,j} \times \left( \log \frac{P(k_i | R)}{1 - P(k_i | R)} + \log \frac{1 - P(k_i | \bar{R})}{P(k_i | \bar{R})} \right)$$

# doc  $i$  contains keyword  $i$  &  $\bar{R}$

prob  $i$  is  $k_i$  &  $\bar{R}$   
doc  $i$  contains keyword  $i$

# doc  $i$  contains keyword  $i$  &  $R$

prob  $i$  is  $k_i$  &  $R$   
doc  $i$  contains keyword  $i$

$$\text{For } P(k_i | R) = \frac{V_i + \frac{n_i}{N}}{V + 1}$$

# doc  $i$  contains keyword  $i$

# doc  $i$  contains keyword  $i$  &  $\bar{R}$

prob  $i$  is  $k_i$  &  $\bar{R}$   
doc  $i$  contains keyword  $i$

# doc  $i$  contains keyword  $i$  &  $R$

prob  $i$  is  $k_i$  &  $R$   
doc  $i$  contains keyword  $i$

$$P(k_i | \bar{R}) = \frac{n_i - V_i + \frac{n_i}{N}}{N - V + 1}$$

## • Structured Text Model

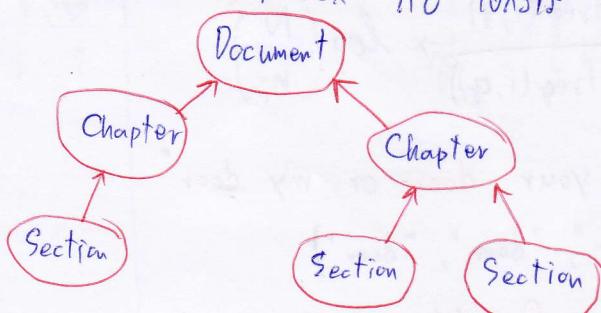
- word នៅលើ title និង weight នៃនំពោះ word ដែលមានអ៊ូរ៉ែ
  - ផ្តល់នូវតម្លៃ នៅលើក្នុង M.R. នៃតម្លៃ
  - នៅ 2 ចំណាំ

## 1) Non-Overlapping Lists

- សៅរ៍ List វារ 
    - Chapter
    - Section
    - SubSection និង text region ដូចជា numbers
  - សំណើនូវ query ទាំងអស់

## 2) Proximal Nodes

- for hierarchical index in 10ms



## Fuzzy Logic

- Fuzzy Set = set memiliki nilai bilangan

- Anggota set dalam skala 0 - 1

- Operator complement  $\mu_A^c(x) = 1 - \mu_A(x)$

intersection  $\mu_{A \cap B}(x) = \mu_A(x) \wedge \mu_B(x)$

$$= \min(\mu_A(x), \mu_B(x))$$

union  $\mu_{A \cup B}(x) = \mu_A(x) \vee \mu_B(x)$

$$= \max(\mu_A(x), \mu_B(x))$$

- A. Simbolisasi index term (keyword)

$$C_{i,j} = \frac{n_{i,j}}{n_i + n_j - n_{i,j}}$$

# doc for i, j = j  
a. simbol  
i ∼ j      # doc for keyword i

- Membuat query i.e. "I would like cat and have dog"  
dari index term i.e. "cat", "dog"

- Dengan menggunakan cat, dog ∼ index term diatas ( $C_{i,j}$ )

ex	Index 1	Index 2	Relate
Dog	Cat	0,33	
~	Bird	0,5	
~	Zebra	~	
~	Zoo	~	
~	Puppy	~	
~	Mouse	~	
~	Robber	~	

(6)

- ឧបតាថ្មីនូវ index term និងនរណា

ex  $\text{dog} = \{d_1 \rightarrow 1, d_2 \rightarrow 0.33, d_3 \rightarrow 1, d_4 \rightarrow 0\}$

$\tilde{x} \text{ dog} \in d_1$

$d_2 \text{ នៃ } \text{dog} \text{ នឹង } \tilde{x} \text{ ក្នុង } d_2$

$\text{cat} \in d_3 \quad c = 0.33$

(នឹង c គឺជាពិនាក់ស្តីពី)

$$d_1 = \{\text{dog, cat, bird, zebra, zool}\}$$

$$d_2 = \{\text{cat, kitty, fish}\}$$

$$d_3 = \{\text{dog, puppy, house, rubber}\}$$

$$d_4 = \{\text{ant, sugar}\}$$

- រាយការណី operator និងសម្រាប់

- union  $\rightarrow \mu_{\text{dog} \cup \text{cat}}(x) = (1 - (1 - \mu_{\text{dog}}(x)) * (1 - \mu_{\text{cat}}(x)))$

- intersect  $\rightarrow \mu_{\text{dog} \cap \text{cat}}(x) = \mu_{\text{dog}}(x) * \mu_{\text{cat}}(x)$

- complement  $\rightarrow \mu_{\overline{\text{dog} \cap \text{cat}}}(x) = (1 - \mu_{\text{dog}}(x)) * (1 - \mu_{\text{cat}}(x))$

- $\text{rank } x$  នៃ  $d_1, d_2, d_3, d_4$  នឹងការបង្ហាញ (a. នូវការបង្ហាញ)

លំដាប់ ranking នៅក្នុងទូទៅ

## Evaluation of IR system

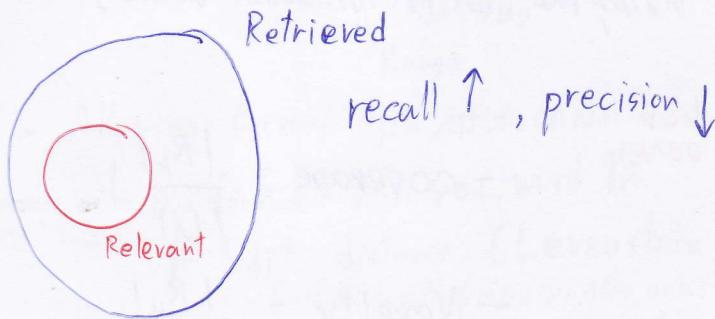
- តាមអាជីវកម្មរបស់វា doc នៃលទ្ធផល Recall និង Precision

- Recall =  $\frac{\# \text{doc} \text{ នៃ } \text{លទ្ធផល}}{\# \text{doc} \text{ នៃ } \text{លទ្ធផល} + \text{ល្អាចិត្ត}}$  = នឹងការបង្ហាញដែលត្រួតពិនិត្យ

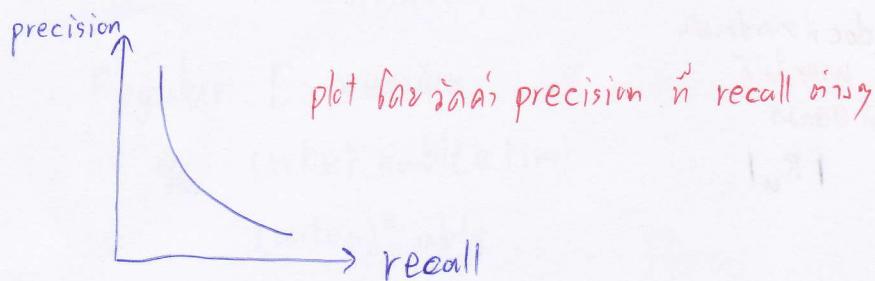
- Precision =  $\frac{\# \text{doc} \text{ នៃ } \text{លទ្ធផល} \text{ ដែលត្រួតពិនិត្យ}}{\# \text{doc} \text{ នៃ } \text{លទ្ធផល}}$    
 $= \frac{\# \text{doc} \text{ នៃ } \text{លទ្ធផល} \text{ ដែលត្រួតពិនិត្យ}}{\# \text{doc} \text{ នៃ } \text{លទ្ធផល} + \text{ល្អាចិត្ត}}$

7

- may recall using memory devices to remember things



- ## - Precision / Recall Curves



- original recall, precision for slide

- ភាគីទំនាក់ស្របតាមរដ្ឋបាល

- Average Precision : ດີຈິກຕ່າງໆ precision (avg) = Mean recall ເພີ້ມ

- R-Precision: ái precision năm nhâi R for  $R = \# \text{ doc nhâi nhâi}$

L Precision Histogram  $\rightarrow$  multivariate algorithm A vs B

$$\text{JMR}_{A/B}(i) = RP_A(i) - RP_B(i)$$

- F-Measure : for harmonic mean  $\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$

$$F = \frac{2}{\frac{1}{R} + \frac{1}{P}} \quad * \text{ in } F \text{ ໃຊ້ວິທີການນັກ recall ເພື່ອສະຫຼຸງ$$

- E-Measure : គឺ F-Measure ក្នុងការសម្រេច weight ស្ថិតិថា precision / recall

$$E = \frac{(1 + \beta^2)}{\frac{\beta^2}{R} + \frac{1}{P}}$$

for  $\beta = 1 \rightarrow$  weight init  
 $\beta > 1 \rightarrow$  ~ precision min  
 $\beta < 1 \rightarrow$  ~ recall min

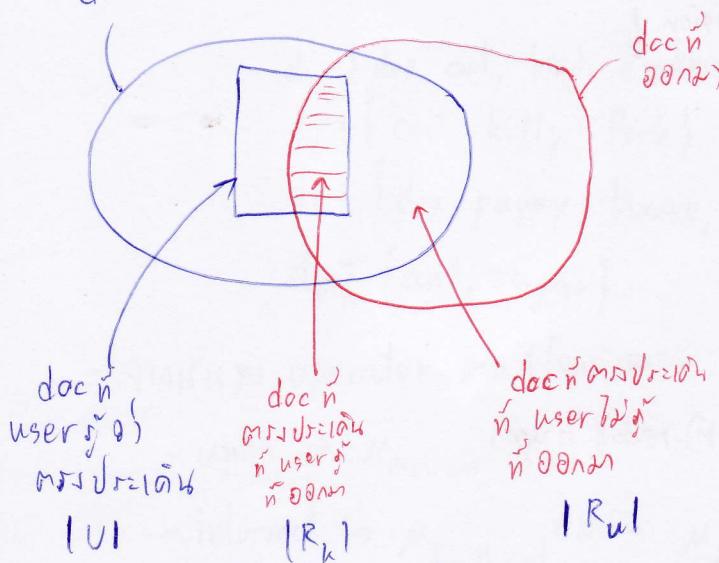
(8)

## User-Oriented Measure

- សិក្សាក្នុង user នៃទូទៅ  $\rightarrow$  អរគុណភាពនៃ user ដែលបាន

ចំណាំប្រើ (ដោយអាជីវកម្មនៃ user ទាំងអស់)

doc in user



$$- \text{coverage} = \frac{|R_kl|}{|UUI|}$$

$$- \text{Novelty} = \frac{|Rkl|}{|Rkl| + |R_kl|}$$

## Query Languages

### Keyword - Base Querying

1) Single-word query ឬ Apple

2) Boolean - query : ឬ single word និងរួចរាល់នៃ Boolean operator  
OR, AND, BUT ឬ Apple AND Banana

3) Natural Language query : ឬនាមឈានឯកសារបាន

4) Context query : ឬ word word

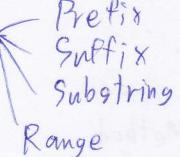
- Phrasal query  $\rightarrow$  ឬនាមឈានឯកសារជាប្រព័ន្ធដែល ex "information theory"  
 $\rightarrow$  ឬនាមឈានឯកសារជាប្រព័ន្ធដែលត្រូវបានត្រួតពិនិត្យដោយ stop words និង use stemming ។  
 $\rightarrow$  ឬនាមឈានឯកសារជាប្រព័ន្ធដែលត្រូវបានត្រួតពិនិត្យដោយ inverted index និងនិងនិង  
និង keywords បានត្រួតពិនិត្យ

- Proximity query  $\rightarrow$  ឬនាមឈានឯកសារជាប្រព័ន្ធដែល  
 $\rightarrow$  ឬនាមឈានឯកសារជាប្រព័ន្ធដែលត្រូវបានត្រួតពិនិត្យដោយ

## ● Pattern Matching

- match ก្នុងចំណាំនាយករដ្ឋមន្ត្រី 100% ទៅ

- Simple Pattern



- Allowing Errors រាល់នូវសម្រាប់ការស្វែងរក

ការស្វែងរកពេលវេលាអំពីរបស់វត្ថុ word បើ

>Edit distance (Levenshtein distance)

គឺជាការស្វែងរកពេលវេលាដែលអាមេរិយាន និងអាមេរិយាន string នៃពីរ

Longest Common Subsequence

គឺជាការស្វែងរកពេលវេលាដែលអាមេរិយាន

- Regular Expression

ex (ule) nable(eling)

(ulen)\* able

## ● Structural Query

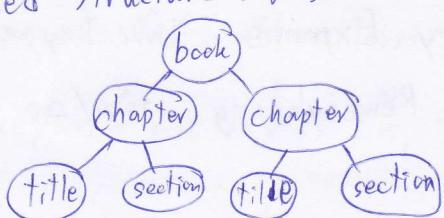
- Its query ដើម្បីរកចំណាំ content ឬ= បានរួចរាល់

word, phrase pattern

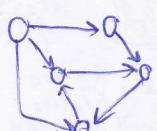
in field this

- បានរួចរាល់ 3 រូបៗ

1) Fixed structure : ជំនួយនៃ field នានា ឬ វិនិស់, email



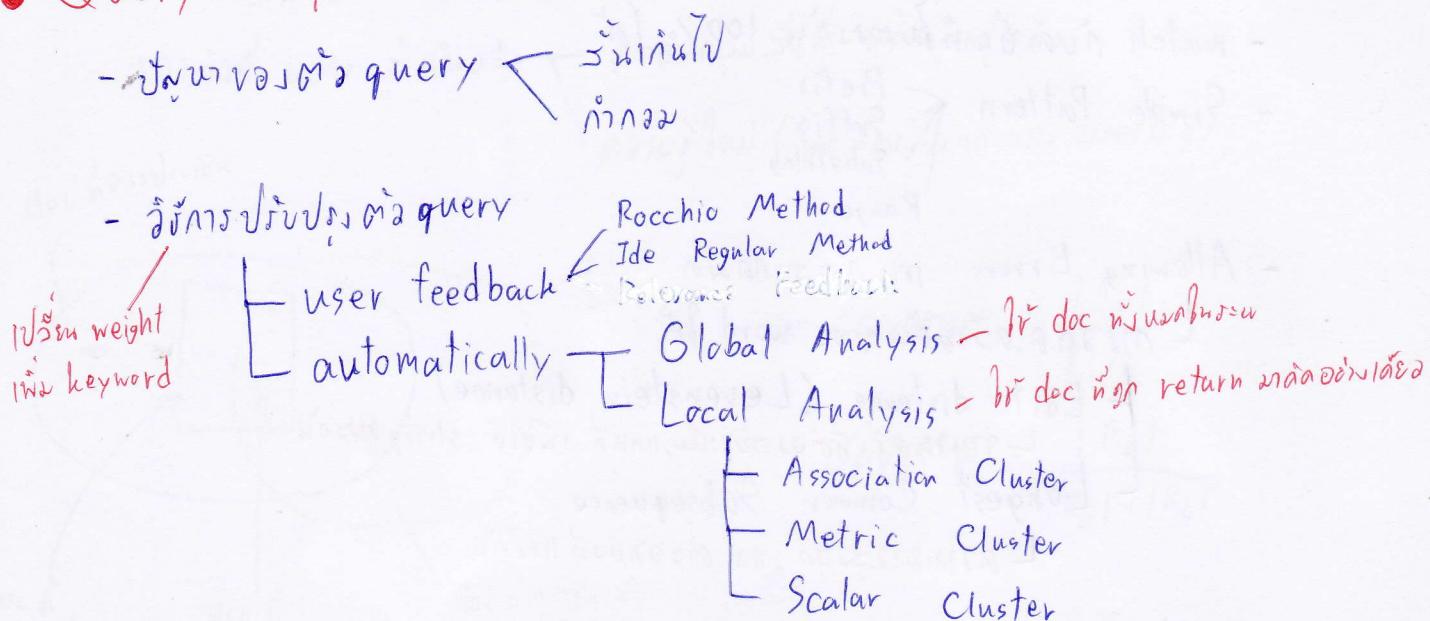
2) Hypertext ; ឬ ជាន់ directed graph បុរាណនៃ text ឬ



3) Hierarchical Structure

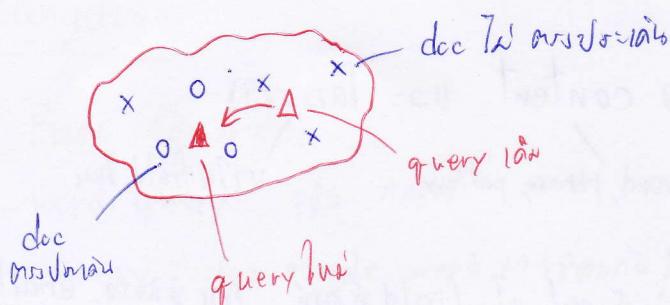
⑩

## Query Operation



## Relevance Feedback

- នៅ query នឹងបានរាយការណ៍ដោយអ្នកប្រើបាន
- ឱ្យ query បានរាយការណ៍ដោយពាក្យសម្រាប់អ្នកប្រើបាន



- និងចំណាំ query
  - ↓ Query Expansion: ពីរ keyword ឬ
  - ↓ Term Reweighting: ពិនិត្យនា W

\* SMF

1) Rocchio Method

$$Q_1 = \alpha Q_0 + \frac{\beta}{n_r} \sum_{\text{adj}_j \in D_r} \vec{d}_j - \frac{\gamma}{n_s} \sum_{\text{adj}_j \in D_s} \vec{d}_j$$

↓ n.r. doc  
 ↓ sum  $\vec{d}$   
 ↓ n.s. doc  
 ↓ sum  $\vec{d}$

$\alpha$  = នរបាលិកិត្តិនុ query នៃវគ្គ

$\beta$  = ~ doc នៃពាក្យសម្រាប់អ្នកប្រើបាន

$\gamma$  = ~ doc នៃពាក្យដែលត្រូវបានស្នើសុំ

## 2) Ide Regular Method

$$\vec{q}_i = \alpha \vec{q}_o + \beta \sum_{\forall j \in D_r} \vec{d}_j - \gamma \sum_{\forall j \in D_n} \vec{d}_j$$

- Խաչող օնդի լուս սահմանը  $n_1, n_2$  առաջ

### 3) Relevance Feedback

$$\vec{q}_m = \vec{q} + \alpha \sum_{\vec{d}_j \in D_r} \vec{d}_j - \beta \sum_{\vec{d}_j \in D_h} \vec{d}_j$$

- ឯកសារនេះត្រូវបានរចនាបានដើម្បីជាអាជីវកម្មសាស្ត្រិយាណីរបស់ខ្លួន។

\* វិជ្ជមាននៃ user feedback (Relevance Feed back) នៅ User mark 7/10

MN / 70000 ຈະ 10w = 100000 ໃນ return ອອນລາຍງານ

will query `getlastBlockNumber` doc in our database which return some data

⑥ Create new query automatically

## Global Analysis

- For document similarity Thesaurus

- គឺជាគារវិភាគនៃវត្ថុការ term (keyword)

- **WWS keyword**  $\vec{w}_i = (w_{i1}, w_{i2}, w_{i3}, w_{iN})$

annular term hi  
in doc j      ↘  
                ↑↑  
W r o s keyword i  
V k dec 2

$$w_{i,j} = \frac{(0.5 + 0.5 \frac{tf_{i,j}}{\max\{tf_{i,j}\}})}{tf_j}$$

$$\text{Term}_i = \sum_{k=1}^N \left[ 0.5 + 0.5 \frac{\text{tf}_{i,k}}{\max\{\text{tf}_{i,k}\}} \right]^2 \text{tf}_k^2$$

$$\text{iff. } = \log \frac{t}{t_0} \quad \text{term nüura}$$

$\text{III}_j = \frac{\text{avg}}{f_j}$  \ given distinct keyword in doc j

(12)

- Term  $C_{u,v} = \vec{k}_u \cdot \vec{k}_v$

$$C_{u,v} = C_{v,u}$$

- Similarity  $\text{sim}(q, k_v) = \vec{q} \cdot \vec{k}_v = \sum_{k_u \in q} w_{u,q} \times C_{u,v}$

- Weighted keyword  $\text{sim}$  function  
from query  $\vec{q}$

$$w_{v,q} = \frac{\text{sim}(q, k_v)}{\sum_{k_u \in q} w_{u,q}}$$

choose  $w_{v,q}$  if  $\text{sim} > 0.53$

ex if sim > query keyword

$$\text{sim}(q, k_1) = 0.53$$

$$\text{sim}(q, k_2) = 0.36$$

$$\text{sim}(q, k_3) = 3.98$$

$$\text{sim}(q, k_4) = 1.87$$

query join  $q = k_1 + k_4$

$$\begin{matrix} / & / \\ 0.53 & 1.87 \end{matrix}$$

if  $\text{sim} > 0.53$   
add keyword

query join  $q' = k_1 + k_3 + k_4$

choose max weight  $w$  for query join

Automatic Local Analysis : If query doc will return 0.0121

- 1 ដែល keyword នឹងរាយនូវ keyword នៃ query តាម
- នូវចំណាំ local clusters
- 1) Association cluster : ប្រាក់នៃ keyword នឹងនូវ keyword ~
- 2) Metric cluster : ពីរសម្រាប់ឱ្យរាយនូវ keyword នឹង
- 3) Scalar cluster : អនុគមន៍ នឹងតាមការបង្ហាញ

## ① Association Cluster

- ឱ្យរាយ matrix រាយការណ៍នូវ keyword នៃ item (សមាជិក)
- ឱ្យរាយ matrix រាយការណ៍នូវ keyword នៃ item : Doc នឹង return 0.0121
- នូវ correlation matrix (C) ដោយ

$$C_{u,v} = \sum_{d_j \in D} f_{s_{u,j}} \times f_{s_{v,j}}$$

- នូវនឹង normalized correlation matrix (S)

$$S_{u,v} = \frac{C_{u,v}}{C_{u,u} + C_{v,v} - C_{u,v}}$$

	A	B	C	D
A	1	0.7	0.18	0.44
B	0.7	1	0.85	0.63
C	0.18	0.85	1	0.63
D	0.44	0.63	0.63	1

term relation

1. {A, B}
2. {B, C}
3. {C, B}
4. {D, B, C}

⇒ 1. {A, B}  
2. {B, C}  
3. {D, B, C}

query តាម A ទៅនីតិវិធីនូវ B និង query តាម B ទៅនីតិវិធី

$$q = A + B$$

B និង C នឹងបាន បាន weight រាយ C នៃ item S<sub>B,C</sub>

query នូវ

$$q = A + (B + 0.85 C)$$

(14)

## Metric Cluster

- ចំណាំសម្រាប់ការគិតថានៅរវាង keyword គួរលាងនូវវិធី correlation matrix

- នឹងចាប់ normalize

$s_u$	$s_1$	$s_2$	$s_3$
$s_1$	0	0.51	0.50
$s_2$	0.51	0	0.60
$s_3$	0.50	0.60	0

stem Relation

$$\left. \begin{array}{l} 1. \{s_1, s_2\} \\ 2. \{s_2, s_3\} \\ 3. \{s_3, s_1\} \end{array} \right\} \Rightarrow \begin{array}{l} 1. \{s_1, s_2\} \\ 2. \{s_2, s_3\} \end{array}$$

$$\text{ឱ្យ } s_1 = \{A, B, C\}, s_2 = \{D, E\}, s_3 = \{F\}$$

query នៃវិធី

$$q = A + 2D$$

$$A \in s_1$$

$$s_1 \cap s_2 \neq \emptyset$$

$$D \in s_2$$

$$s_2 \cap s_3 \neq \emptyset$$

query បុគ្គលិក

$$q' = (s_1 + 0.5s_2) + 2(s_2 + 0.6s_3)$$

$$= s_1 + 2.5s_2 + 1.2s_3$$

- ឯកសារ

$$s_{u,v} = \frac{c_{u,v}}{|V(s_u)| \times |V(s_v)|}$$

នូវវិធី  
សម្រាប់ការ

## Scalar Cluster

- ប្រាក់បញ្ជីនៃ stem នឹងតើសរុបនូវ Tree, Water, Fertilizer

- ការគិតថានៅរវាងការប្រើប្រាស់  $\vec{s}_u = (c_{u,1}, c_{u,2}, \dots, c_{u,t})$

a. សម្រាប់ការ  
stem u នូវ  
keyword នឹង 1

- នឹងចាប់ scalar association matrix a. សម្រាប់ការប្រើប្រាស់ stem

$$s_{u,v} = \frac{\vec{s}_u \cdot \vec{s}_v}{|\vec{s}_u| \times |\vec{s}_v|}$$

## Text Operation

- ณ full text → index term  
(keyword)
- Lexical Analysis: ใช้ input character stream  
ดู stream word or token  
(like binary)
- ณ hyphen  
state-of-the-art ⇒ state of the art
- ณ stopword  
- ณ article, preposition, conjunction, ...
- ณ Stemming วิธีลดคำ
- ตัดส่วนของคำที่ไม่ใช้
- ขั้นตอน 1) Affix removal: ตัดส่วน -ing, un-, -s  
2) Successor Variety: ตัดส่วนของคำต่อไปนี้ successor variety
- /
 

คำนำหน้า prefix ที่ตัดออก  
อย่างไรก็ได้
- 3) Table Lookup ⇒
 

Term	stem
engineering	engineer
engineered	engineer
engineer	engineer

(16)

4) n-gram: រាយការណ៍ប៉ុន្មានលាក្ខណៈសម្រាប់ការសង្គម និងការសង្គម

- ឯកសារ 2-gram មានចំនួន ( $n=2$ )

- ការគិតអនុវត្តន៍ដោយសារតម្លៃ

$$S = \frac{2C}{A+B}$$

A = # diagram និងបានឱ្យឈានទៅលម្អិត (ឯកសារ)

B = # diagram និងមិនឱ្យឈានទៅលម្អិត

C = # diagram និងមិនឱ្យឈានទៅលម្អិត និងមានសម្រាប់ការប្រើប្រាស់ និង 2

- Q-gram stemmer

ឯកសារក្នុង Q នឹង នូវវឌ្ឍន៍ និង បានគិតថា Q-gram distance  
ឯកសារក្នុង Q នឹង នូវវឌ្ឍន៍ និង បានគិតថា Q-gram distance

- Index Terms Selection

- ពាណិជ្ជកម្ម

- រាយការណ៍ index term ដែលមានការប្រើប្រាស់ និង word និង phrase  
រាយការណ៍ index term ដែលមានការប្រើប្រាស់ និង word និង phrase

## ② Indexing

- Invert file : ឯកសារតាមការប្រើប្រាស់

Vocabulary = ពាណិជ្ជកម្ម (keyword)

occurrence = ការប្រើប្រាស់នៃការប្រើប្រាស់

ន.ន.និងចាប់, និង doc 7 ន.

    └ Suffix array

    └ Signature file

+ នៅក្នុង - Keyword = term និងចាប់ document, សារការណ៍ដែលមាន doc មែនក្នុងរបៀបនេះ

- controlled vocabulary = list នៃគោលដៅនិងចាប់ keyword ។

## Inverted File

- ដាក់បង្ហាញអាមុនវត្ថុ Doc 7us
  - Block addressing តើ នីមួយៗ text ឬ block និងអាជីវកិច្ច block 7us
  - ឱ្យអាជីវកិច្ច Doc 7us ត្រូវបានគ្រប់គ្រង (semi-static)

## Suffix Tree and Suffix Array

- māññ tree iou fai feng wāññ māññ kāññ vāññ
  - māññ suffix array fai feng māññ vāññ māññ kāññ

# Signature Files

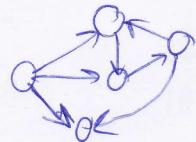
(18)

## • Multimedia IR

- ແອນຂໍ້ມູນສັງເກດ (IRW, IRB)
- Data Modeling
  - RDBMS  $\rightarrow$  ນິບສະແດງ multimedia DBMS
    - inheritance
    - polymorphism
  - OODBMS  $\rightarrow$  ປຳເນົາ object oriented (OO)
    - encapsulation
    - ໂຄສະໝັກ ປະຊາທິປະໄຕ ອາວຸດ, ດີວິຈານ
  - $\rightarrow$  ປຳເນົາ OQL
  - Object-Relational  $\rightarrow$  ໄກສະໝັກພາບມານຸມ  
 $\rightarrow$  ປຳເນົາ SQL
    - ໂຄສະໝັກ ອາວຸດ ລົມລົມ object
- ການຈົບຈັດຂອງ ທີ່ ກົດຫຼັກຂອງ DBMS ສະໜັບ B-Tree
- R-tree (Region tree)
  - ໃນນີ້ spatial object ດີວິຈານ minimum bounding rectangle (MBR)
  - ອັນດາ MBR ພົບທຶນ
  - ຖື່ນໃນ R-tree ຕັ້ງໆ ມີ MBR ທັງໆ ໃນລົມ
  - ບື່ນດູ ກົດ MBR ທີ່ ພົບທຶນ
  - R-tree ສະໜັບພົນກົດ ຢັດ leaf ສະໜັບທຶນ

## ① Searching The Web

- គោរព web site ឬ graph នឹងបាន



- ឬ spider ឬឯកសារណ៍ (page) និង ពេទ័រ index នៅមុខ
- (crawler)
- Search engine

### 1) Centralized Architecture

- crawler-indexer
- Crawler ឬសម្រាប់ឯកសារនិង web និង ឬ index នៅមុខ server នៅរំលែក
- ឬវេជ្ជ crawler = robot, spider, wanderer, knowbot
- ឬស្ថាបី query ដោយមិនត្រូវ server នៃទាំងអស់ ឬ index នៃ crawler ទៀត
- ឬស្ថាបី crawler ដួច Depth-first ឬ Breadth-first Crawling
- ឬស្ថាបី crawler ដែលបានចាត់ចាយនៅលើ server ទៀត

### 2) Distributed Architecture

- ឬ 2 ម៉ោង
- 1. Gatherers: អប់រំកុំរែនសាស្ត្រ index នៃ web server  
នៅមុខ ឬ Broker
- 2. Brokers: ឬប្រភេទសាស្ត្រ ឬទូទៅនៃ user ដែលមានទូទៅនៃ Gatherer នៅរឿង ឬនិង ឬ Broker នៅលើគ្មាននៃ Broker នៅរឿង
- ឬស្ថាបី load នៃទំនួរ
- ស្ថាបី cache ឬបុន្មាននៃនាមូលធម៌

# Information Storage and Retrieval

1 / 2552

## 1. อธิบายความหมาย

- knowbot คือ โปรแกรมที่ไปเขียน web pages ต่างๆ แล้วส่ง index กลับมายัง server กลาง เพื่อใช้ในการทำ search engine (ชื่ออีกหลายชื่อของ knowbot คือ crawler, robot, spider, wanderer )

- n-gram คือ การแบ่งคำออกเป็นชุด ชุดละ n ตัวอักษร เช่น statistics แจกแจงออกมาเป็น st,ta,at,ti,is,st,ti,cs เมื่อแบ่งเป็น 2-gram ( $n = 2$ )

- R-Tree คือ data structure รูปแบบหนึ่งที่ใช้เก็บ spatial object ด้วย minimum bounding rectangle (MBR) ซึ่งเป็นกรอบสี่เหลี่ยมล้อม object ที่เป็นตัวข้อมูล

- Vector Model คือหนึ่งโมเดลในการสืบค้นข้อมูลแบบคลาสสิก โดยแทนที่เอกสาร และ query ด้วยเวกเตอร์ โดยค่าแต่ละมิติในเวกเตอร์คือน้ำหนักแต่มิติ keyword แต่ละคำ ความตรงประเด็นของเอกสารวัดจากมุมที่เอกสารทำกับ query โดยยิ่งมุมเล็กยิ่งตรงประเด็น

- Probabilistic Model คือ หนึ่งโมเดลในการสืบค้นข้อมูลแบบคลาสสิก โดยใช้ความน่าจะเป็นที่ผู้ใช้จะเจอเอกสารที่ตรงประเด็น และมีการรับ feedback จากผู้ใช้ว่าเอกสารไหนตรง/ไม่ตรงประเด็น เพื่อปรับความน่าจะเป็น เพื่อให้ครั้งต่อไปสืบค้นเจอเอกสารที่ตรงประเด็นมากยิ่งขึ้น

- Proximity เป็นวิธีการทำ context query คือการหาคำท้ายคำ โดยคำนับอยู่ห่างกัน ได้โดยจะต้องกำหนดระยะห่างสูงสุดระหว่างคำนับก่อนเขียน (enhance,retrival,4) จะ match “enhance the power of retrieval”

- Supra - index คือ index ที่สร้างขึ้นเหนือ suffix array เพื่อแบ่ง suffix array ออกเป็นช่วงๆ เวลาที่หาคำใหม่ก็ให้ไปหาใน Supra - index ก่อน ว่าคำนั้นน่าจะอยู่ในช่วงไหนของ suffix array แล้วค่อยไปหาคำแทนของคำใน suffix array

- MIMD

- Recall คือ ความสามารถในการดึง doc ที่ตรงประเด็นออกมามีค่าเท่ากับ จำนวนเอกสารที่ออกมารอแล้วตรงประเด็น / จำนวนเอกสารทั้งหมดที่ตรงประเด็น
    - ค่านี้ใช้วัดประสิทธิภาพของ algorithm ในการ query
  - Global Analysis เป็นวิธีการปรับปรุงตัว query แบบอัตโนมัติ โดยใช้ข้อมูลจากเอกสารทั้งหมดในระบบนำมาทำ similarity thesaurus แล้วหาว่าควรเพิ่ม keyword อื่นหรือไม่ และการปรับ weight เป็นเท่าไหร่
2. ให้ออกสารมา 7 เอกสาร พร้อม query กำหนดหน้าหัวนักของ keyword ของแต่ละเอกสารและน้ำหนักของ keyword ใน query พร้อมจัดลำดับเอกสารที่ถูก return โดยการคำนวณให้คำนึงถึงการเป็นรากศัพท์ตัวย
  - ให้ใช้วิธี tf idf ในการหาน้ำหนัก โดยตอนที่หาความถี่ของ keyword ให้คิดถึงรากศัพท์ด้วยเช่นพากเติม -ed -ing -s -es ถือเป็นคำเดียวกัน
    - ส่วนการ query และจัดลำดับให้ทำแบบ vector model จะดีสุด (คิดเป็นเวกเตอร์ ค่าความเหมือนคือ  $\cos(\theta)$ )
3. กำหนดเอกสารที่ตรงประเด็นทั้งหมดในระบบมาให้ มี 2 อัลกอริทึมในการ query พร้อมผลลัพธ์ที่ได้จากการ query ให้วิเคราะห์โดยใช้หลักการใน Evaluation ทุกวิธี
  - ค่าที่ใช้วัดก็มี Average Precision, R-Precision, F-Measure, E-Measure
4. เก็บอัลกอริทึม Proximity query ที่ยอมให้มีระยะห่างระหว่าง keyword ในเอกสารห่างกันได้สูงสุด 6 คำ (ปรับเปลี่ยนจาก phrasal search) พร้อมยกตัวอย่างประกอบ
  - หา doc ที่มี keyword แต่ละคำใน phrase ทุกคำ
  - แต่ละ doc ที่มี keyword อยู่ให้สร้าง array เก็บตำแหน่งของคำแต่ละคำใน phrase
  - เริ่มจากคำที่ขนาด array น้อยที่สุด
  - วนลูปแต่ละตำแหน่งใน array แล้วไปหา array ของคำอื่นๆ ในตำแหน่งที่อยู่ใกล้กันตำแหน่งนั้นที่สุด แล้วดูว่าระยะห่าง keyword ห่างเกิน 6 หรือไม่ ถ้าไม่เกินแสดงว่าเอกสารนั้นเป็นคำตอบ
    - ถ้าเกินให้ดูตำแหน่งต่อไปจนกว่าจะเจอ หรือวนครบทุกตำแหน่งแล้วไม่เจอก็แปลงว่าเอกสารนั้นไม่ใช่คำตอบ
5. มีเอกสารมาให้ 9 เอกสารพร้อม keyword ในเอกสาร กำหนดให้มี 3 processors ให้แสดงวิธีการแบ่งการทำงานแต่ละแบบพร้อมวิเคราะห์