Retrieval Strategies: Vector Space Model and Boolean

(COSC 416)

Nazli Goharian nazli@cs.georgetown.edu



Retrieval Strategies

- Manual Systems
 - Boolean, Fuzzy Set
- Automatic Systems
 - Vector Space Model
 - Language Models
 - Latent Semantic Indexing
- Adaptive
 - Probabilistic, Genetic Algorithms, Neural Networks, Inference Networks



















Pivoted Cosine Normalization

• Comparing likelihood of retrieval and relevance in a collection to identify *pivot* and thus, identify the new *correction factor*.

$$SC(Q, D_i) = \frac{\sum_{j=1}^{i} w_{qj} d_{ij}}{(1.0 - s) + (s) \frac{\sqrt{\sum_{j=1}^{i} (d_{ij})^2}}{avgn}}$$

Avgn: average document normalization factor over entire collection *s*: can be obtained empirically





		_	
Q: "gold	l silver truck"		
D ₁ : "Shi	pment of gold damaged	l in a fire"	
D ₂ : "De	livery of silver arrived i	n a silver truck"	
$D_3^2 \cdot Sh$	inment of gold arrived i	n a truck"	
D <i>J</i> . DI	ipinent of gold unived i	in a track	
Id	Term	df	idf
1	а	3	0
2	arrived	2	0.176
3	damaged	1	0.477
4	delivery	1	0.477
5	fire	1	0.477
6	gold	2	0.176
7	in	3	0
8	of	3	0
9	silver	1	0.477
	shipment	2	0.176
10			





Summary: Vector Space Model

• Pros

- Fairly cheap to compute
- Yields decent effectiveness
- Very popular
- Cons
 - No theoretical foundation
 - Weights in the vectors are arbitrary
 - Assumes term independence



Boolean Retrieval

- *Expression*:=
 - term
 - -(expr)
 - NOT expr (not recommended)
 - expr AND expr
 - expr OR expr
- (cost OR price) AND paper AND NOT article

















- Extended Boolean supports term weight and proximity information.
- Example of incorporating term weight:
 - Ranking by term frequency (Sony Search Engine)
 x AND y: tf_x x tf_y
 x OR y: tf_x + tf_y
 - NOT x: 0 if $tf_x > 0$, 1 if $tf_x = 0$
- User may assign term weights cost and +paper

