







## How do we represent text?

#### • Remember: computers don't "understand" anything!

# • "Bag of words"

- Treat all the words in a document as index terms for that document
- Assign a "weight" to each term based on "importance"
- Disregard order, structure, meaning, etc. of the words
- Simple, yet effective!

#### Assumptions

- Term occurrence is independent
- Document relevance is independent
- "Words" are well-defined

The iSch University of Maryl



Выступая в Мещанском суде Москвы экс-глава ЮКОСа заявил не совершал ничего противозаконного, в чем обвиняет его генпрокуратура России.

भारत सरकार ने आर्थिक सर्वेक्षण में वित्तीय वर्ष 2005-06 में सात फ़ीसदी विकास दर हासिल करने का आकलन किया है और कर सुधार पर ज़ोर दिया है

#### 日米連合で台頭中国に対処…アーミテージ前副長官提言

조재영 기자= 서울시는 25일 이명박 시장이 '행정중심복합도시" 건설안 에 대해 '군대라도 동원해 막고싶은 심정"이라고 말했다는 일부 언론의 보도를 부인했다. The iSch of Maryl

Univ

The iSchool University of Maryland



## **Boolean Retrieval**

- Users express queries as a Boolean expression
  - AND, OR, NOT
  - Can be arbitrarily nested
- Retrieval is based on the notion of sets
  - Any given query divides the collection into two sets: retrieved, not-retrieved
  - Pure Boolean systems do not define an ordering of the results







## Extensions

- Implementing proximity operators
  - Store word offset in postings
- Handling term variations
  - Stem words: love, loving, loves  $\ldots \rightarrow lov$

#### **Strengths and Weaknesses**

#### Strengths

- Precise, if you know the right strategies
- Precise, if you have an idea of what you're looking for
- Implementations are fast and efficient
- Weaknesses
  - Users must learn Boolean logic
  - Boolean logic insufficient to capture the richness of language
  - No control over size of result set: either too many hits or none
  - When do you stop reading? All documents in the result set are considered "equally good"
  - What about partial matches? Documents that "don't quite match" the query may be useful also

## **Ranked Retrieval**

• Order documents by how likely they are to be relevant to the information need

Uni

- Estimate relevance(q, d<sub>i</sub>)
- Sort documents by relevance
- Display sorted results
- User model
  - Present hits one screen at a time, best results first
  - At any point, users can decide to stop looking
- How do we estimate relevance?
  - Assume document is relevant if it has a lot of query terms
  - Replace relevance  $(q, d_i)$  with sim $(q, d_i)$
  - Compute similarity of vector representations



The iSch

Univ



#### **Similarity Metric**

- How about  $|d_1 d_2|$ ?
- Instead of Euclidean distance, use "angle" between the vectors
  - It all boils down to the inner product (dot product) of vectors



• Term weights consist of two components

- Local: how important is the term in this document?
- Global: how important is the term in the collection?
- Here's the intuition:
  - Terms that appear often in a document should get high weights
  - Terms that appear in many documents should get low weights

The iS University of Mary

- How do we capture this mathematically?
  - Term frequency (local)
  - Inverse document frequency (global)





## **Sketch: Scoring Algorithm**

- Initialize accumulators to hold document scores
- For each query term *t* in the user's query
  - Fetch *t*'s postings
  - For each document,  $score_{doc} += w_{t,d} \times w_{t,q}$
- Apply length normalization to the scores at end
- Return top *N* documents

# MapReduce it?

#### • The indexing problem

- Must be relatively fast, but need not be real time
- For Web, incremental updates are important
- The retrieval problem
  - Must have sub-second response
  - For Web, only need relatively few results

# Indexing: Performance Analysis

- Inverted indexing is fundamental to all IR models
- Fundamentally, a large sorting problem
  - Terms usually fit in memoryPostings usually don't
- How is it done on a single machine?
- How large is the inverted index?
  - Size of vocabulary
  - Size of vocabular
    Size of postings



The iSchool University of Maryland



The iSchool University of Maryland



the	1130021	from	96900	or	54958
of	547311	he	94585	about	53713
to	516635	million	93515	market	52110
а	464736	year	90104	they	51359
in	390819	its	86774	this	50933
and	387703	be	85588	would	50828
that	204351	was	83398	you	49281
for	199340	company	83070	which	48273
is	152483	an	76974	bank	47940
said	148302	has	74405	stock	47401
it	134323	are	74097	trade	47310
on	121173	have	73132	his	47116
by	118863	but	71887	more	46244
as	109135	will	71494	who	42142
at	101779	say	66807	one	41635
mr	101679	new	64456	their	40910
with	101210	share	63925		

the	59	from	92	or	101
of	58	he	95	about	102
to	82	million	98	market	101
а	98	year	100	they	103
in	103	its	100	this	105
and	122	be	104	would	107
that	75	was	105	you	106
for	84	company	109	which	107
is	72	an	105	bank	109
said	78	has	106	stock	110
it	78	are	109	trade	112
on	77	have	112	his	114
by	81	but	114	more	114
as	80	will	117	who	106
at	80	say	113	one	107
mr	86	new	112	their	108
with	91	share	114		

# MapReduce: Index Construction

#### Map over all documents

- Emit term as key, (docid, tf) as value
- Emit other information as necessary (e.g., term position)
- Reduce
  - Trivial: each value represents a posting!
  - Might want to sort the postings (e.g., by docid or tf)
- MapReduce does all the heavy lifting!

#### **Query Execution**

- MapReduce is meant for large-data batch processing
  Not suitable for lots of real time operations requiring low latency
- The solution: "the secret sauce"
  - Most likely involves document partitioning
  - Lots of system engineering: e.g., caching, load balancing, etc.

The iSchool University of Maryland



• High-throughput batch query execution:

- Instead of sequentially accumulating scores per query term:
- Have mappers traverse postings in parallel, emitting partial score components

The iSchool University of Maryland

The iSch University of Maryla

- Reducers serve as the accumulators, summing contributions for each query term
- MapReduce does all the heavy lifting
  - Replace random access with sequential reads
  - Amortize over lots of queries
  - Examine multiple postings in parallel

# **Questions?**