Information Retrieval INFO/CS 4300

- Instructor: Chris Buckley

 Office hours Wednesdays 11am Gates 231
- Piazza will be the main communication tool
 - <u>https://piazza.com/cornell/fall2014/info4300/home</u>
 - Lecture notes will appear there.TA office hours and locations appear there.

Course Admin

- Critique 1, Homework 1 Graded, grades available on-line through CMS, hard copy can be picked up in the homework return room in Gates 216, open Mon-Fri noon-4pm
- Homework 2: Due today at beginning of class (hardcopy plus CMS).
- Warning: class **may** get interrupted late today by guest visitor (Amit Singal, Director of Search at Google).

Previous Lectures

- Overview
- Evaluation 1
- Indexing
- Retrieval
 - Boolean Model
 - Tf*idf weighting
 - Vector Space Model
 - Retrieval Optimization (DAAT, TAAT, safe vs non-safe)
 - Basic Probabilistic Model

Retrieval Models

- Older models
- Boolean retrieval (still used, special applications)
- Vector Space model (still used with tf*idf weighting for general retrieval)
- Basic Probabilistic model
 Newer Probabilistic Models
- BM25
- Language models
- Newer tf*idf variants
- Pivoted unique normalizationCombining evidence (later in course)
- Combining evidence (later in course)
 Inference networks
- Interence networ
 Learning to Rank

Probability Ranking Principle

- Robertson (1977)
 - "If a reference retrieval system's response to each request is a ranking of the documents in the collection in order of decreasing probability of relevance to the user who submitted the request,
 - where the probabilities are estimated as accurately as possible on the basis of whatever data have been made available to the system for this purpose,
 - the overall effectiveness of the system to its user will be the best that is obtainable on the basis of those data."

Binary Independence Model

- Scoring function is
 - $\sum_{i:d_i=1} \log \frac{p_i(1-s_i)}{s_i(1-p_i)}$
- Query provides information about relevant documents.
- If we assume p_i constant, s_i approximated by entire collection, get idflike weight

$$\log \frac{0.5(1-\frac{n_i}{N})}{\frac{n_i}{N}(1-0.5)} = \log \frac{N-n_i}{n_i}$$

Contingency Table Relevant Non-relevant Total $p_i = (r_i + 0.5)/(R+1)$ $s_i = (n_i - r_i + 0.5)/(N - R + 1)$ (the +0.5, +1 factors are *smoothing* values) Gives scoring function: $\sum_{i:d_i=q_i=1} \log \frac{(r_i+0.5)/(R-r_i+0.5)}{(n_i-r_i+0.5)/(N-n_i-R+r_i+0.5)}$

Basic Probabilistic Model

Advantages

- Gives (with enough assumptions) basis for idf weighting model. · Important since gives insight into how to incorporate additional information
- Can be extended
- Disadvantages
 - Poor performance in itself
 - Smoothing needed for any application
 - Natural term dependence extensions don't work in practice.

BM25

- Popular and effective ranking algorithm based on binary independence model
 - adds document and query term weights

 $\sum_{i \in Q} \log \frac{(r_i + 0.5)/(R - r_i + 0.5)}{(n_i - r_i + 0.5)/(N - n_i - R + r_i + 0.5)} \cdot \frac{(k_1 + 1)f_i}{K + f_i} \cdot \frac{(k_2 + 1)qf_i}{k_2 + qf_i}$

- $-k_1$, k_2 and K are parameters whose values are set empirically
- $K = k_1((1-b) + b \cdot \frac{dl}{avdl})$ where *dl* is doc length - Typical TREC value for k_1 is 1.2, k_2 varies from 0 to
- 1000, b = 0.75

- r, is the # of relevant documents containing term i
- (set to 0 if no relevancy info is known)
- n_i is the # of docs containing term i N is the total # of docs in the collection
- R is the number of relevant documents for this query
- (set to 0 if no relevancy info is known)
 f_i is the frequency of term i in the doc under consideration
- t, is the requency of term in the doc under consideration (if is the frequency of term in the query k, determines how the tf component of the term weight changes as f, increases. (if 0, then tf components is ginored.) "typical value for TRECS 1.2; so f, is very non-linear (similar to the use of log in term viso of the vector space model) --- after 3 or 4 occurrences of a term, additional occurrences will have little impact. k, has a similar role for the query term weights. Typical values (see slide) make the equation less sensitive to k, than k, because query term frequencies are much lower and less variable than doc term frequencies.
- K is more complicated. Its role is basically to normalize the tf component by document
- b regulates the impact of length normalization. (0 means none; 1 is full normalization.)

BM25 Example

- Query with two terms, "president lincoln", (qf = 1)
- No relevance information (r and R are zero)
- N = 500,000 documents
- "president" occurs in 40,000 documents (n₁ = 40,000)
- "lincoln" occurs in 300 documents (n₂ = 300)
- "president" occurs 15 times in doc ($f_1 = 15$)
- "lincoln" occurs 25 times (f₂ = 25)
- document length is 90% of the average length (dl/avdl = .9)
- $k_1 = 1.2$, b = 0.75, and $k_2 = 100$
- K = 1.2 · (0.25 + 0.75 · 0.9) = 1.11

BM25 Example BM25(Q, D) = $\log \frac{(0+0.5)/(0-0+0.5)}{(40000-0+0.5)/(500000-40000-0+0+0.5)}$ $\times \frac{(1.2+1)15}{1.11+15} \times \frac{(100+1)1}{100+1} \\ + \log \frac{(0+0.5)/(0-0+0.5)}{(00-0.5)/(0-0+0.5)}$ $+\log\frac{(0+0.5)/(0-0+0.5)}{(300-0+0.5)/(500000-300-0+0+0.5)}$ $\times \frac{(1.2+1)25}{1.11+25} \times \frac{(100+1)1}{100+1}$ $= \log 460000.5/40000.5 \cdot 33/16.11 \cdot 101/101$ $+\log 499700.5/300.5 \cdot 55/26.11 \cdot 101/101$ $2.44 \cdot 2.05 \cdot 1 + 7.42 \cdot 2.11 \cdot 1$ = 5.00 + 15.66 = 20.66

BM25 Example

• Effect of term frequencies





Retrieval Models

- Older models
 - Boolean retrieval (still used, special applications)
 Vector Space model (still used with tf*idf weighting for general retrieval)
- Basic Probabilistic model
- Newer Probabilistic Models
 - BM25
 - Language models
- Newer tf*idf variants
- Pivoted unique normalization
 Combining evidence (Later in course)
- Inference networks
- Learning to Rank

Language Model

- Unigram language model
 - probability distribution over the words in a language
 - generation of text consists of pulling words out of a "bucket" according to the probability distribution (and replacing them)
- N-gram language model
 - some applications use bigram and trigram language models where probabilities depend on previous words

Language Model

- A *topic* in a document or query can be represented as a language model
 - i.e., words that tend to occur often when discussing a topic will have high probabilities in
 - the corresponding language model
- Multinomial distribution over words
- text is modeled as a finite sequence of words, where there are t possible words at each point in the sequence
- commonly used, but not only possibility
- doesn't model burstiness

LMs for Retrieval

- 3 possibilities:
 - probability of generating the query text from a document language model
 - probability of generating the document text from a query language model
 - comparing the language models representing the query and document topics
- · Models of topical relevance

Query-Likelihood Model

- Rank documents by the probability that the query could be generated by the document model (i.e. same topic)
- Start with a query, so calculate P(D|Q) to rank the documents
- Use Bayes' Rule
- $p(D|Q) \stackrel{rank}{=} P(Q|D)P(D)$ Assuming prior is uniform, unigram model

$$P(Q|D) = \prod_{i=1}^{n} P(q_i|D)$$

Estimating Probabilities

· Obvious estimate for unigram probabilities is

$$P(q_i|D) = \frac{f_{q_i,D}}{|D|}$$

- Maximum likelihood estimate – makes the observed value of f_{so} (the frequency of term q_i in document D) most likely
- Problem: If query words are missing from document, score for the document will be 0
 - Missing 1 out of 6 query words (score of 0) is the same as missing 5 out of 6

Example

- D = (tropical, tropical, fish, fish, fish, water, aquarium)
- Q = salt water tropical fish
- P ("tropical" | D) = 2/7
- P ("fish" | D) = 3/7
- P ("water" | D) = 1 / 7
- P ("salt" | D) = 0/7
- P(Q | D) = 2/7 * 3/7 * 1/7 * 0

Smoothing

- Document texts are a *sample* from the language model
 - Missing words should not have zero probability of occurring
- Smoothing is a technique for estimating probabilities for missing (or unseen) words
- lower (or *discount*) the probability estimates for words that are seen in the document text
- assign that "left-over" probability to the estimates for the words that are not seen in the text

Estimating Probabilities

- Estimate for unseen words is $\alpha_D P(q_i | C)$
 - $-P(q_i|C)$ is the probability for query word *i* in the *collection* language model for collection *C* (background probability)
 - $-\alpha_D$ is a parameter
- Estimate for words that occur in a query is
 (1 - α_D) P(q_i|D) + α_DP(q_i|C)
- Different forms of estimation come from different α_D

Jelinek-Mercer Smoothing

- α_D is a constant, λ
 - Gives estimate of $p(q_i|D) = (1-\lambda) \frac{f_{q_i,D}}{|D|} + \lambda \frac{c_{q_i}}{|C|}$
 - Ranking score
 - $P(Q|D) = \prod_{i=1}^{n} ((1-\lambda) \frac{f_{q_i,D}}{|D|} + \lambda \frac{c_{q_i}}{|C|})$
- Use logs for convenience
- accuracy problems multiplying small numbers $\log P(Q|D) = \sum_{i=1}^{n} \log((1-\lambda)\frac{f_{q_i,D}}{|D|} + \lambda \frac{c_{q_i}}{|C|})$







Frequency of	Frequency of	OL.
"president"	"lincoln"	score
15	25	-10.53
15	1	-13.75
15	0	-19.05
1	1 25 -12.99 0 25 -14.40	



5

QL Language Models

- Advantages
 - Principled probabilistic model
 - Intuitively appealing captures notions we think are important
 - Smoothing itself is appealing
 - Extensible, including smoothing
- Works well
- Disadvantages
- Parameter estimation
- Smoothing methods are somewhat ad hoc (but they are targeted!)

Retrieval Models

- Older models
 - Boolean retrieval (still used, special applications)
 - Vector Space model (still used with tf*idf weighting for general retrieval)
 Basic Probabilistic model (Not used)
- Basic Probabilistic mode
 Newer Probabilistic Models
- BM25
- Language models
- Newer tf*idf variants
- Pivoted unique normalization
- Combining evidence (later in course)

 Inference networks
 - Learning to Rank

Improved vector space retrieval

- · Standard weighting scheme is tf*idf and cosine similarity
 - We've discussed tf*idf variants (and you'll see more of this next week with project 1) and possible alternatives
 - Can cosine document length normalization be improved?
 - Recall that cosine normalization put all vectors on a unit hypersphere
 Nice theoretically is it what we want in practice?
 - The answer turns out to be no, but this was discovered in a very roundabout fashion

Case study: Query expansion digression

- I was looking at query expansion techniques in the mid 1990's (early TREC days where the field was doing experimenting with long documents in test collections for the first time)
- · Basic tf*idf weights with cosine normalization
- Discovered a method of expanding (lengthening) queries by adding related terms with appropriate weights
 - Worked very well improving retrieval results, in fact too well!!!
 - Was still getting improvements (small, but significant) by adding 200-300 terms.

Query expansion => length normalization

- Query expansion results yielded very smooth improvement curves as terms were added. At the end of even 100 added terms, the terms being added were common terms that looked visually to be pretty random. Either
 - A. I had discovered something fundamental like an "atomic" unit or measurement of semantic information content, and how to weight it. Or
 - B. Something else was going on
- · I alternated between the two possibilities for a couple of weeks
- The answer turned out to be B.
 - Adding random common terms (with low weights) to the query increases the score of documents randomly (wrt original query), but will tend to increase the score of longer documents more than shorter documents!
 - Cosine length normalization turned out to be biased against long documents

Investigating length normalization

- We had had length normalization on our list of factors to investigate for a while. Steve Robertson and Okapi had just come out with BM25 which approached length normalization differently.
- · How do you investigate length normalization issues?
- First step, figure out what is happening in practice now
- We divided the document collection into bins according to length
- Calculated Prob (retrieval of D | D in bin i)
- Calculated Prob (relevance of D | D in bin i)



Cosine favors shorter documents

- In an ideal world, we would hope the two plots would be much closer.
- If we change our retrieval normalization to match Prob(Relevance), will that improve performance?
- Amit Singhal took over finding methods to do this.
- First step was introducing pivoted normalization.



Pivoted Normalization

- Pivoted_normalization = (1.0 slope) * pivot + slope * old_normalization
- If we take the pivot point as being the average old normalization document, and divide by it (doesn't affect ranking)
- Pivoted_normalization = (1.0 slope) + slope *
 <u>old normalization</u>
 average old normalization
- · A document of length average old normalization will be unchanged
- Slope can be interpreted as our "belief in length"
- Derived very differently, but similar to Okapi's BM25 length normalization factor (done first) $K = k_1((1-b) + b \cdot \frac{d}{avdl})$



Pivoted unique terms

- Still a problem with very long documents.
- Analyzing further, issue is the tf contribution in very long documents

 Even standard taking logs is too much (for the length part)
- Instead of cosine of log (1.0+tf) weighting or something similar, use the number of unique terms.
- Pivoted_unique = (1-slope) * pivot + slope * # of unique terms
- · Again can take the pivot as the average number of unique terms in a doc



Pivoted unique: Final tf*idf weighting

• Lnu.ltc where Lnu weighting in documents is

 $(1 + \log(tf))$ 1+log(average tf)

(1.0 - slope) + slope * (# of unique terms) average # of unique terms

Where slope = 0.20 works well across collections

Retrieval Models

- Older models
 - Boolean retrieval (still used, special applications) - Vector Space model (still used with tf*idf weighting for general retrieval)
- Basic Probabilistic model
- Newer Probabilistic Models BM25
- Language models
- Newer tf*idf variants
- Pivoted unique normalization Combining evidence (later in course)
- Inference networks
- Learning to Rank

Features of these newer models

- All work about the same as far as test collection evaluation goes
- All require estimating parameter(s)
 - The estimations are not motivated by the models
 - But most parameters are insensitive to small changes (should do reasonably on other collections)
- · Precursors, in some ways, to "Learning to Rank" models, covered later