

605.744 Information Retrieval: Topics for Midterm Exam

The exam will consist primarily of problem-solving and short-answer questions. The exam will be closed-book. You can bring one 8.5 x 11" sheet of paper with anything written on the front and back (e.g., study notes, equations, sample problems). I recommend that you bring a simple or scientific calculator to aid with arithmetic, but note that no wireless devices or computing devices with substantial storage are allowed; use of smart phones, iPads, e-Readers, etc.. is prohibited. Key topics that you should be familiar with include:

- issues in processing text, such as dealing with punctuation, stemming
- dictionaries, and how they can be represented and compressed
- document frequency, term frequency, IDF, Zipf's law
- Boolean, vector-space (e.g., cosine), binary independence probabilistic model, and statistical language modeling retrieval models
- methods for term weighting (e.g., binary, tf, idf, tf-idf, $1+\log(\text{tf})$)
- indexing process, algorithms for indexing, including inverted file construction when memory is limited or not very limited.
- inverted file data structures and compression techniques (gap-lists; gamma and delta codes)
- how to score and rank documents for a query
- wildcard querying
- evaluation metrics, especially, precision, recall, precision at x docs, interpolated recall-precision graphs, and mean average precision
- test collections, TREC evaluations
- term operations such as stopword removal, stemming, n-gram tokenization; the effect of stopword removal and stemming on lexicon and inverted file sizes
- relevance feedback and Rocchio's method for query modification
- term similarity measures between two terms (e.g., mutual information)
- general familiarity with assigned readings (IIR: 1-9, 11-12; assigned papers, which are available on the course website)