

Search Engine Architecture

The goal of this class is to implement a highly-scalable and highly-available search engine by leveraging off-the-shelf tools. We will use elements of information retrieval, NLP, machine learning, and large-scale cluster computing with a focus on practical implementation. Prerequisites are an understanding of basic data structures and algorithms, and working knowledge of Python.

Textbook: Introduction to Information Retrieval, Christopher D. Manning.

- Introduction
 - o Goals of search engines and their architectures
 - o Component and dataflow overview
- Basic modeling
 - o Jaccard coefficient, vector space, cosine similarity, TF-IDF weighting, evaluation
 - o sklearn, NLTK
- Processing
 - o Crawling, filtering noise, tokenization, stemming, phonetic algorithms
 - o Scalability: Hadoop, HDFS, HBase
- Indexing
 - o Inverted indexes, locality-sensitive hashing: Redis
 - o Distributed hash tables: Chord
- Searching
 - o Reverse proxies: HAProxy
- Collaborative filtering
 - o Nonnegative matrix factorization, nearest-neighbor methods
- Link Analysis
 - o PageRank
- Text categorization
 - o Discriminative methods
- Text clustering
 - o Latent semantic indexing, dimensionality reduction, latent Dirichlet allocation