

Why Detect Exact Duplicates

- Smarter crawling**
 - Avoid returning many duplicate results to a query
 - Allow fetching from the fastest or freshest server
- Better connectivity analysis**
 - By combining in-links from the multiple mirror sites to get an accurate PageRank
 - Avoid double counting out-links
- Add redundancy in result listings**
 - If that fails you can try: <mirror>/samepath"
- Clearly**: Crawlers need not crawl pages that are identical or near identical
- Ideally**: given the web's scale and complexity, priority must be given to content that has not already been seen before
 - Saves resources (on the crawler end, as well as the remote host)
 - Increases crawler politeness
 - Reduces the analysis that a crawler will have to do later

Solving the Duplicate/Near-Duplicate Detection Problem

- Duplicate: Exact match** eg: compare password (use hashing)
 - Solution: compute fingerprints using cryptographic hashing
 - SHA-1 and MD5 are/were the two most popular cryptographic hashing methods
- Near-Duplicate: Approximate match**
 - Solution: compute the syntactic similarity with an edit-distance measure
 - Use a similarity threshold to detect near-duplicates
 - e.g., Similarity > 80% => Documents are "near duplicates" or document

Identifying Near Identical Web Pages - Two Approaches

- Test for similar fingerprints** - Treat web documents as a set of features, constituting an n -dimensional vector, and transform this vector into an f -bit fingerprint of a small size
 - Use Simhash and Hamming Distance to compute the fingerprint (more about this later)
 - Compare fingerprints and look for a difference in at most k bits
 - E.g. see Manku et al., WWW 2007, Detecting Near-Duplicates for Web Crawling, <http://www2007.org/papers/paper215.pdf>
- Compute syntactic similarity** TF-IDF: values words that are not common
 - Shingling / n-grams
 - Compute w-shingling (more about shingling later)
 - Broder et al., WWW 1997, Finding Near Duplicate Documents

General Properties of Distance Measures

- Distance measure must satisfy 4 properties**
 1. Non negative distances
 2. $d(x,y) = 0 \Leftrightarrow x=y$ only if
 3. $d(x,y) = d(y,x)$ symmetric
 4. $d(x,y) \leq d(x,z) + d(z,y)$ triangle inequality
- There are several distance measures that can play a role in locating duplicate and near-duplicate documents
 - Euclidean distance** - $d([x_1, \dots, x_n], [y_1, \dots, y_n]) = \sqrt{\sum (x_i - y_i)^2}$ $i=1 \dots n$
 - Jaccard distance** - $d(x,y) = 1 - SIM(x,y)$ or 1 minus the ratio of the sizes of the intersection and union of sets x and y
 - Cosine distance** - the cosine distance between two points (two n element vectors) is the angle that the vectors to those points make; in the range 0 to 180 degrees
 - Edit distance** - the distance between two strings is the smallest number of insertions and deletions of single characters that will convert one string into the other
 - Hamming distance** - between two vectors is the number of principal components in which they differ (usually used on boolean vectors)

Why Detect Near Duplicates

cluster similar things together

- Clustering – BIG USE CASE FOR CONTENT DETECTION AND ANALYSIS**
 - Given a news article some people might wish to see "related articles" describing the same event
- Data extraction**
 - Given a collection of similar pages, e.g. movie reviews, try to identify the schema underlying the collection so one can extract and categorize the information
- Plagiarism**
 - Identify pairs that seem to have significantly borrowed from each other
- Spam detection**
 - Spammers typically send similar emails en masse, so one can use near-similarity techniques to identify the spam
- Duplicates in domain-specific corpora**
 - To identify near-duplicates arising out of revisions, modifications, copying or merging of documents

Identifying Identical Web Pages – Four Approaches

- Compare character by character** two documents to see if they are identical
 - very time consuming !!
- Hash just the first few characters and compare** only those documents that hash to the same bucket
 - But what about web pages where every page begins with <HTML> ??
- Use a hash function** that examines the entire document
 - But this requires lots of buckets different file to same hash is possible
- Better approach** - pick some fixed random positions for all documents and make the hash function depend only on these;
 - This avoids the problem of a common prefix for all or most documents, yet we need not examine entire documents unless they fall into a bucket with another document
 - But we still need a lot of buckets

General Paradigm

- translate documents to number**
- Define a function / that captures the contents of each document in a number
 - E.g. hash function, signature, fingerprint
- Create the pair $\langle f(doc), ID \text{ of } doc \rangle$ for all doc
- Sort the pairs and compare
- Documents that have the same f-value or an f-value within a small threshold are believed to be duplicates

Computing Near Similarity

- Definition of Shingle:**
 - a contiguous subsequence of words in a document is called a *shingle*;
 - The 4-shingling of the phrase below produces a bag of 5 items: "a rose is a rose" => a set $S(D,w)$ is defined as $\{(a_rose_is_a), (rose_is_a_rose), (is_a_rose_is), (a_rose_is_a), (rose_is_a_rose)\}$
 - $S(D,w)$ is the set of shingles of a document D of width w
- Similarity Measures**
 - Resemblance**(A,B) is defined as $\frac{\text{size of } (S(A,w) \cap S(B,w))}{\text{size of } (S(A,w) \cup S(B,w))}$
 - Containment**(A,B) is defined as $\frac{\text{size of } (S(A,w) \cap S(B,w))}{\text{size of } (S(A,w))}$
 - $0 \leq \text{Resemblance} \leq 1$
 - $0 \leq \text{Containment} \leq 1$
 - See *On the resemblance and containment of documents*, Conf. on Compression and Complexity, DEC Research Center, 1997

Shingling Example

- Original text
 - "Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species"
 - All 3-grams (there are 14 of them)
 - 13 of the 14 3-grams include fish (include fish found in tropical environments around the world, including both freshwater and saltwater)
 - Only 1 3-gram does not include fish (salt water species)
- Select only those hash values that are divisible by some number, e.g. here are selected hash values for 3-grams
 - 938, 664, 463, 822, 492, 798, 78, 669, 143, 236, 913, 908, 604, 553, 870, 779
 - Select only those hash values that are divisible by some number, e.g. here are selected hash values for 3-grams
 - 664, 463, 236, 908
- Near duplicates are found by comparing fingerprints and finding pairs with a high overlap

Example taken from Search Engines: Information Retrieval in Practice

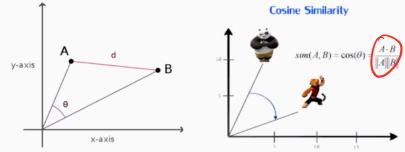
Jaccard Similarity and Shingli

- The Jaccard similarity of sets A and B, $J(A,B)$, is defined as $\frac{|A \cap B|}{|A \cup B|}$
- A k -shingle of a document is any substring of length k found in the document
 - We ignore multiple blanks, tabs, newlines
- Choosing an appropriate value for k
 - If k is too small, then most sequences of length k will occur in most documents
 - k should be picked large enough that the probability of any given shingle appearing in any document is low
 - E.g. suppose in email only letters and white space occur (this is a simplified example); then there are $27^5 = 14,348,907$ possible 5-shingles; since emails are generally small, the value $k=5$ should work well
- In the next few slides we will show how computing shingles is equivalent to computing the Jaccard Similarity

Cosine similarity

Important

- Don't treat features as sets per-se, instead graph them
- Then compute the cosine of the angle between the two graphs
 - 1-d shown
 - Can you imagine for multiple features?



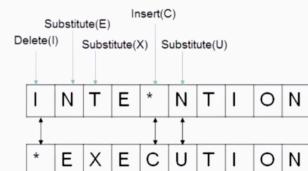
Edit distance for short documents

- Consider the feature set of two documents A, B as a set of characters

- How many changes ^{edit distance} to go from set A to set B?

- Tuning parameters

- Cost of each operation



With a distance metric you can cluster
many cluster algorithms are based on distance

- Basic Clustering

- Hierarchical



General Paradigm

- Define a function f that captures the contents of each document in a number
 - E.g. hash function, signature, fingerprint
- Create the pair $\langle f(\text{doc}_i), \text{ID of doc}_i \rangle$ for all doc;
- Sort the pairs
- Documents that have the same f -value or an f -value within a small threshold are believed to be duplicates

- With $\langle \text{jaccard/edit/cosine} \rangle$ distance as $f(\text{doc}_i)$

- Define a threshold that defines the max accepted delta between $f(\text{doc}_i)$ and $f(\text{doc}_i+1)$

