

CSEP 521: Applied Algorithms

Lecture 16 – Document Similarity

Richard Anderson, February 25, 2021



Announcements

Topics

- Document similarity
- MinHash
- String similarity
 - Edit Distance / Longest Common Subsequence
 - Sharding strings
- Dimension reduction

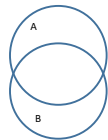
Document Similarity

- Want to be able to identify documents that are “very close” to each other
- Very large number of documents
- Individually pre-process documents
 - Save a small amount of data per document (sketch)
 - Perform similarity tests based on sketch



Jaccard Similarity

$$\text{Jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$



Let X be the characteristic vector for A where x_j is the multiplicity of item j and Y be the characteristic vector for B where y_j is the multiplicity of item j .

$$\text{Jaccard}(A, B) = \frac{\sum_j \min(x_j, y_j)}{\sum_j \max(x_j, y_j)}$$

Representation scheme

- Tokenize document
- Break document into shards
- Hash each shard into a domain of size 2^{64} (unsigned long)
- Treat as a bag of words*
- Use Jaccard similarity measure

far out in the uncharted backwaters of the unfashionable end of the western spiral arm of the galaxy lies a small unregarded yellow sun

1. far out in the uncharted
2. out in the uncharted backwaters
3. in the uncharted backwaters of
4. the uncharted backwaters of the
5. uncharted backwaters of the unfashionable
6. backwaters of the unfashionable end
7. of the unfashionable end of
8. the unfashionable end of the
9. unfashionable end of the western
10. end of the western spiral
11. of the western spiral arm
12. the western spiral arm of

* In this application, we use bag of words without multiplicity

Similarity testing

- Identify document pairs that have high similarity by doing pairwise comparison
- Precompute hashes of shards – n shards for document of n tokens
- Cost of comparison is O(n)

- How to improve this: reduce the amount of information stored per document

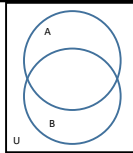
MinHash

- U is the domain (in this case, the hash of the shards, $[0 \dots 2^{64}]$)
- Choose a random permutation π on U
- Let $A \subseteq U$
- MinHash(A) = $\text{argmin}_{x \in A} \pi(x)$
 - MinHash is the smallest element of A under the random permutation

```

U = {a, b, c, d, e, f, g, h, i, j}
A = {b, c, e, g}
B = {c, e, f}
π1 = [a, c, d, i, j, h, b, e, f, g]
π2 = [j, i, g, c, b, h, e, a, d, f]
    
```

An amazing result



$$\Pr[\text{MinHash}(A) = \text{MinHash}(B)] = \frac{|A \cap B|}{|A \cup B|} = \text{Jaccard}(A, B)$$

Using the MinHash

- Identify document pairs where $\text{Jaccard}(A, B) \geq 0.95$
- Run MinHash with k independent random permutations
- Number of times $\text{MinHash}(A) = \text{MinHash}(B)$ is a good estimate of Jaccard Similarity
- Compute the k MinHashes for each documents as a sketch
- Comparison of documents requires k comparisons

Similarity of Strings

- String edit distance – how many edits to convert S_1 into S_2
- Edit operations: Add character, Remove character, (Change character)

BARTHOLEMESIMPSON → KRUSTYTHECLOWN

B	A	R	T	H	O	L	E	M	E	W	S	I	M	P	S	O	N							B-	
A	R	T	H	O	L	E	M	E	W	S	I	M	P	S	O	N									A-
R	T	H	O	L	E	M	E	W	S	I	M	P	S	O	N										K+
K	R	T	H	O	L	E	M	E	W	S	I	M	P	S	O	N									U+
K	R	U	T	H	O	L	E	M	E	W	S	I	M	P	S	O	N								S+
K	R	U	S	T	H	O	L	E	M	E	W	S	I	M	P	S	O	N							Y+
K	R	U	S	T	Y	H	O	L	E	M	E	W	S	I	M	P	S	O	N						T+
K	R	U	S	T	Y	T	H	O	L	E	M	E	W	S	I	M	P	S	O	N					O-
K	R	U	S	T	Y	T	H	L	E	M	E	W	S	I	M	P	S	O	N						L-
K	R	U	S	T	Y	T	H	E	M	E	W	S	I	M	P	S	O	N							M-
K	R	U	S	T	Y	T	H	E	E	W	S	I	M	P	S	O	N								E-
K	R	U	S	T	Y	T	H	E	W	S	I	M	P	S	O	N									C+
K	R	U	S	T	Y	T	H	E	C	W	S	I	M	P	S	O	N								L+

BARTHOLEMESIMPSON → KRUSTYTHECLOWN

K	R	U	S	T	Y	T	H	E	C	L	W	S	I	M	P	S	O	N		O+
K	R	U	S	T	Y	T	H	E	C	L	O	W	S	I	M	P	S	O	N	S-
K	R	U	S	T	Y	T	H	E	C	L	O	W	I	M	P	S	O	N		I-
K	R	U	S	T	Y	T	H	E	C	L	O	W	M	P	S	O	N			M-
K	R	U	S	T	Y	T	H	E	C	L	O	W	P	S	O	N				P-
K	R	U	S	T	Y	T	H	E	C	L	O	W	S	O	N					S-
K	R	U	S	T	Y	T	H	E	C	L	O	W	O	N						O-
K	R	U	S	T	Y	T	H	E	C	L	O	W	N							

Longest Common Subsequence

- $C=c_1...c_g$ is a subsequence of $A=a_1...a_m$ if C can be obtained by removing elements from A (but retaining order)
- $LCS(A, B)$: A maximum length sequence that is a subsequence of both A and B

ocurranecc attacggct bartholemewsimpson
 occurrence tacgacca krustythecrown

Edit Distance and LCS

- String A has length n and B has length m
- Suppose that A is converted to B by removing k characters and adding j characters
 - Number of unchanged characters is $c = n - k = m - j$
 - Edit distance d is $k + j$
 - $n + m = 2c + k + j = 2c + d$
 - $d = n + m - 2c$
- Minimizing the edit distance is maximizing the length of the common sequence

LCS Optimization

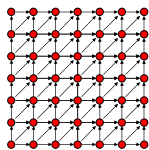
- $A = a_1a_2...a_m$
- $B = b_1b_2...b_n$
- $Opt[j, k]$ is the length of $LCS(a_1a_2...a_j, b_1b_2...b_k)$
- Optimization recurrence

If $a_j = b_k, Opt[j, k] = 1 + Opt[j-1, k-1]$

If $a_j \neq b_k, Opt[j, k] = \max(Opt[j-1, k], Opt[j, k-1])$

$Opt[j, 0] = Opt[0, k] = 0$

Dynamic Programming Computation



Code to compute $Opt[n, m]$

```

for (int i = 0; i < n; i++)
for (int j = 0; j < m; j++)
if (A[i] == B[j])
    Opt[i, j] = Opt[i-1, j-1] + 1;
else if (Opt[i-1, j] >= Opt[i, j-1])
    Opt[i, j] := Opt[i-1, j];
else
    Opt[i, j] := Opt[i, j-1];
    
```


Of course this is going to be hashing



- Choose a good hash function $h: U \rightarrow 2^{32}$
- Let $f_1(x) = h(x) \bmod 2$
- 1 bit representation
 - If $x = y$, then $f_1(x) = f_1(y)$
 - If $x \neq y$, then $\Pr[f_1(x) = f_1(y)] \leq 1/2$
 - Property preserved with probability at least 50%
- Repeat with k independent hash functions h_1, \dots, h_k
 - If $x = y$, then $f_i(x) = f_i(y)$ for all $i = 1, \dots, k$
 - If $x \neq y$, then $\Pr[f_i(x) = f_i(y) \text{ for all } i = 1, \dots, k] \leq 2^{-k}$
- To achieve error of δ , we need to use $k = \lceil \log_2 1/\delta \rceil$