

CSCI 4022 Spring 2021

Shingling and Minhashing



Example: prove that the edit distance is a proper distance measure.

Announcements and To-Dos

Announcements:

1. HW 1 up!

Minute form comments: ☺

Python!

Last time Recap

Example: prove that the edit distance is a proper distance measure.

Last time Recap

insert/delete string characters
 \nearrow m, n such operations (count) to turn

Example: prove that the edit distance is a proper distance measure.

$x \rightarrow y$
 $y \rightarrow x$

Solution:

1. No negative distances. *Counts are non-negative \checkmark .*
2. Distances are only zero from a point to itself. *1) if points identical \rightarrow dist 0 \checkmark
 2) if $d=0 \rightarrow$ points are identical*
3. Distance is symmetric. *operations are invertible in the same # of steps.*
4. Distances satisfy the **triangle inequality**.

and \rightarrow nan \geq
 and \rightarrow nan \rightarrow nan \geq


must be \geq or mde.

(formal pf: use "minimum" $\frac{1}{\epsilon}$
 proof by contradiction)

Last time Recap

Example: prove that the edit distance is a proper distance measure.

Solution:

1. $d(x, y) \geq 0$: ☒. We're counting operations: obviously can't be negative!
2. $d(x, y) = 0 \leftrightarrow x = y$: ☒. Two statements: if $x = y$, clearly the edit distance is zero. On the other hand, if we don't edit a string x then we get the same string, so no edit distance also implies $x = y$.
3. $d(x, y) = d(y, x)$ ☒ Given a sequence of edits to turn x to y , we can turn this into a same-length sequence of edits to get y to x by changing each insertion to a deletion and vice versa. (in other words: each operation is invertible!)
4. $d(x, y) \leq d(x, z) + d(z, y)$ ☒. If we turn x to y by first turning x to z , then it the number of edits can't be less than going directly to y . 

So edit distance is a proper distance measure.

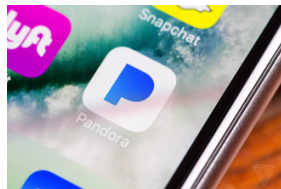
Last time: Similarity

Many problems - both in this course and in general - can be expressed as the task of finding **similar** elements.

We often phrase this as finding “near-neighbors,” which may exist in **high-dimensional** spaces.

Examples:

1. Documents with similar words/code
2. Users who watch similar movies
3. Songs that have similar attributes
4. Images with similar features
5. etc., etc.



Last time: Distance Measures

We first considered distances in **Euclidean** spaces, like the general L_r -norm:

$$d(x, y) = \left(\sum_{i=1}^n (x_i - y_i)^r \right)^{1/r}$$

L_1 : Manhattan

L_2 : Euclidean /
hypotenuse

L_∞ : "max" distance

which worked well for vectors in a continuous space.

For comparing sets, we saw **Jaccard similarity**,

$$\text{sim}(S, T) = \frac{|S \cap T|}{|S \cup T|}$$

and the associated **distance**: $d = 1 - \text{sim}$.


From there we introduced two other non-Euclidean distances meant to describe strings or lists:

edit distance:, # of deletions/insertions required to move from one string to another;

Hamming distance, number of components by which two strings differ

Today's plan

Hi!
lots of huge sets



Today we discuss the plan of how to implement Jaccard Similarity on a *large data* scale. This is a three part-process.

1. First, we **shingle**: we map a document into a set of numbers.
2. Second, we discuss **approximate Jaccard similarity**: by considering the properties of reshuffling/permuting bitstrings representing sets.
3. Finally, we discuss **approximate permutations** via modular transformations.

Document Similarity

Jaccard similarity was a measure for the similarity between *sets*. Consider the task of measure similarity between **documents**. Applications might include:

1. plagiarism detection
2. mirroring web pages
3. finding news articles from the same sources or authors

The first question: how can we represent an entire document as a set?

Document Similarity

Jaccard similarity was a measure for the similarity between *sets*. Consider the task of measure similarity between **documents**. Applications might include:

1. plagiarism detection
2. mirroring web pages
3. finding news articles from the same sources or authors

The first question: how can we represent an entire document as a set?



- A document is just a long string of *characters*. We can
- Represent a document as a collection of *substrings* of those characters.
 - these substrings of characters are called *shingles*

k-shingles

Definition: a *k-shingle* for a document is any ^{substring} ~~sub~~string of length k found in that document.

Example: Suppose we have a document D that consists only of the word "remember" and we pick $k = 2$. What is the entire set of all 2-shingles for D ?

$\{\underline{re}\}$; overlap $\cup \{mb\}$
 $\cup \{\underline{em}\}$ $\cup \{be\}$
 $\cup \{me\}$ $\cup \{er\}$
 $\cup \{em\}$

$"aa" ?$ No
 $"ab" ?$ No
 $"ac" ?$ No

k-shingles

gram: chops into subsequences of words!
 $k=2$: "we love class" $\rightarrow \{(we\ love), (love\ class)\}$

Definition: a k -shingle for a document is any substring of length k found in that document.

Example: Suppose we have a document D that consists only of the word "remember" and we pick $k = 2$. What is the entire set of all 2-shingles for D ?

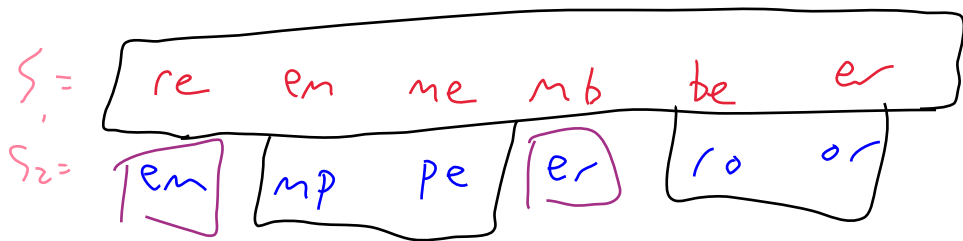
Solution: {re, em, me, mb, be, er}.

Note:

1. "em" would appear twice... but this is a *set*, so we don't include duplicates?
2. A variant would consider a *bag* of 2-shingles, which might include duplicates.
3. We generally want characters to be case-insensitive and not using punctuation. So we pre-process by lowering cases and replace any spaces, tabs, newlines, etc. with a general white-space character. (letters + _ = 27 total)

k-shingles

Suppose document D_1 consists only of the word “remember” and document D_2 consists only of the word “emperor.” Representing the documents as their sets of 2-shingles (call them S_1 and S_2), what is the Jaccard similarity of D_1 and D_2 ?



Total = 10 distinct 2-shingles

2 are shared

$$\text{Sim}(S_1, S_2) = \frac{2}{10}$$

k-shingles

Suppose document D_1 consists only of the word “remember” and document D_2 consists only of the word “emperor.” Representing the documents as their sets of 2-shingles (call them S_1 and S_2), what is the Jaccard similarity of D_1 and D_2 ?

Solution:

$$S_1 = \{\text{re, em, me, mb, be, er}\}.$$

$$S_2 = \{\text{em, mp, pe, er, ro, or}\}.$$

$$|S_1 \cup S_2| = |\{\text{re, em, me, mb, be, er, mp, pe, ro, or}\}| = 10.$$

$$|S_1 \cap S_2| = |\{\text{em, er}\}| = 2.$$

So the similarity is $2/10 = 0.2$.

k-shingles

The next task: how do we choose k to represent a document?

1. If k is too small, most sequences of k chars might appear in most documents. As a result, everything looks similar.
2. If k is too large, many documents might share no sequences as all and everything will look different.

General rule: choose k just large enough such that the probability of any given shingle appearing in any given document is low.

k-shingles

The next task: how do we choose k to represent a document?

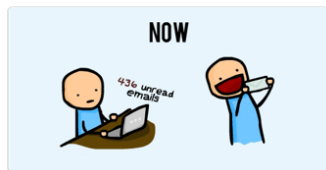
1. If k is too small, most sequences of k chars might appear in most documents. As a result, everything looks similar.
2. If k is too large, many documents might share no sequences as all and everything will look different.

General rule: choose k just large enough such that the probability of any given shingle appearing in any given document is low.

Cool, but what does “low” mean?

k-shingles

General rule: choose k just large enough such that the probability of any given shingle appearing in any given document is low.



Example: suppose our corpus of documents to consider is a bunch of emails.

Suppose we consider only letters, all converted to lower case.

With 27 total characters (white space, too!) we'd have $27^k \approx 14.3$ million total possible shingles.

$$k = 5$$

Since typical e-mails are typically much shorter than this, we expect $k = 5$ to work well... and it does!

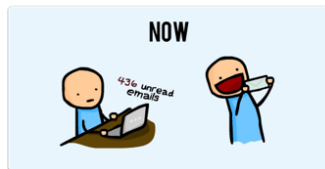
k-shingles

Example: suppose our corpus of documents to consider is a bunch of emails.

But are all characters (letters/spaces) equally likely in an e-mail? **Definitely not!**

Common letters and blanks are far more common!

As a result, we might revise our heuristic: if only 20 of the characters are "common," we'd think of only around 20^k common shingles. Still, $20^5 \approx 3.2$ million, which means $k = 5$ for e-mails is likely ok. Longer documents like essays, code, etc. might require e.g. $k = 9$.



hashing

Example: Suppose document D_1 is the word “banana”; document D_2 is the word “bandit,”; and document D_3 is the work “brand.” Representing the documents only as their **characteristic matrix**, what are the Jaccard similarities of the 3 documents?

	banana	bandit	brand

hashing

Example: Suppose document D_1 is the word “banana”; document D_2 is the word “bandit,”; and document D_3 is the work “brand.” Representing the documents only as their **characteristic matrix**, what are the Jaccard similarities of the 3 documents?

Definition: The *characteristic matrix* of a set of shingles is the matrix where each row is a shingle and each column is a document. Row i and col j is a 1 if shingle i appears in document j .

Solution:

	banana	bandit	brand	
	1	1	0	ba
	1	1	1	an
	1	0	0	nd
	0	1	1	di
	0	1	0	it
	0	0	1	br
	0	0	1	ra

hashing

Example: Suppose document D_1 is the word “banana”; document D_2 is the word “bandit,”; and document D_3 is the work “brand.” Representing the documents only as their **characteristic matrix**, what are the Jaccard similarities of the 3 documents?

Definition: The *characteristic matrix* of a set of shingles is the matrix where each row is a shingle and each column is a document. Row i and col j is a 1 if shingle i appears in document j .

Solution: The Jaccard similarity for columns 1 and 2 is $2/6=1/3$.

	banana	bandit	brand
ba	1	1	0
an	1	1	1
na	1	0	0
nd	0	1	1
di	0	1	0
it	0	1	0
br	0	0	1
ra	0	0	1

hashing

Often, we *hash* the shingles instead of storing the substrings themselves.

A hash:

1. Map each substring of length k to an integer
2. Each document is now represented by a set of integers or bucket numbers/indices
3. We have 8 shingles here, so we need at least 8 buckets. We want more to be safe!

Suppose we choose 19. *← "large" prime*

Definition: A *hash function* is a modular function to map string characters to buckets. It may use e.g. ASCII values of those characters. We could use:

Handwritten: #s ↓

	banana	bandit	brand
ba	1	1	0
an	1	1	1
na	1	0	0
nd	0	1	1
di	0	1	0
it	0	1	0
br	0	0	1
ra	0	0	1

$$h(k) := (\text{sum of ASCII values of characters in shingle}) \bmod 19$$

hashing

Often, we *hash* the shingles instead of storing the substrings themselves. We could use:

$$h(k) := (\text{sum of ASCII values of characters in shingle}) \bmod 19$$

Which results in...

1. $ba \rightarrow \text{ord}(b) + \text{ord}(a) \bmod 19 = 5$
2. $an \rightarrow \text{ord}(a) + \text{ord}(n) \bmod 19 = 17$
3. $na \rightarrow \text{ord}(n) + \text{ord}(a) \bmod 19 = 17$

	banana	bandit	brand
ba	1	1	0
an	1	1	1
na	1	0	0
nd	0	1	1
di	0	1	0
it	0	1	0
br	0	0	1
ra	0	0	1

hashing

Often, we *hash* the shingles instead of storing the substrings themselves. We could use:

$$h(k) := (\text{sum of ASCII values of characters in shingle}) \bmod 19$$

Which results in...

1. $ba \rightarrow \text{ord}(b) + \text{ord}(a) \bmod 19 = 5$
2. $an \rightarrow \text{ord}(a) + \text{ord}(n) \bmod 19 = 17$
3. $na \rightarrow \text{ord}(n) + \text{ord}(a) \bmod 19 = 17$

Ack! A **collision!**.

	banana	bandit	brand
ba	1	1	0
an	1	1	1
na	1	0	0
nd	0	1	1
di	0	1	0
it	0	1	0
br	0	0	1
ra	0	0	1

hashing

To avoid collisions, we often use hash functions that stretch the strings out to much larger numbers.

reality: $(17 \cdot \text{first} + 101 \cdot \text{second}) \bmod 19$

$h(k) := 2 \cdot \text{first ASCII value} + \text{second ASCII value} \bmod 19$

Prime

Which results in less overlaps, especially if those coefficients get further apart!

We can also replace the substrings for each row in the characteristic matrix with the hash bucket numbers.

Then, if we sort by bucket number, we get much easier indexing!

	banana	bandit	brand
8	1	1	0
0	1	1	1
13	1	0	0
16	0	1	1
1	0	1	0
3	0	1	0
6	0	0	1
2	0	0	1

(a)

hashing

The sorted, renumbered characteristic matrix is to the left. The missing numbers correspond to possibly *many* missing shingles.

	banana	bandit	brand
0	1	1	1
1	0	1	0
2	0	0	1
3	0	1	0
6	0	0	1
8	1	1	0
13	1	0	1 0
16	0	1	1

The data science question:

Is there a fast way to efficiently determine which pairs are most similar? Idea:

1. Replace each set/document with a smaller representation called a *signature*.
2. Use the signatures to *estimate* Jaccard similarity.

minhashing

So how do we find “signatures?”

	banana	bandit	brand
0	1	1	1
1	0	1	0
2	0	0	1
3	0	1	0
6	0	0	1
8	1	1	0
13	1	0	0
16	0	1	1

Definition: To *minhash* a set represented by a column of the characteristic matrix, pick a permutation of the rows. The *minhash value* of any column is the number of the first row, in the permuted order, in which the column has a 1.

Definition: A row in the *signature matrix* is built up from the minhash values of all columns under a given permutation

↳ rows: permutation per row
col: that doc's first "/"

Permutations and Jaccard

Why does this make any sense?

	banana	bandit	brand
0	1	1	1
1	0	1	0
2	0	0	1
3	0	1	0
6	0	0	1
8	1	1	0
13	1	0	0
16	0	1	1

Idea: if we pick a row *at random* and look at 2 columns, what happens?

1. Type “X” is both 1s. This should be in the numerator *and* denominator of Jaccard similarity.
2. Type “Y” is exactly 1 1. This should be in the denominator of Jaccard similarity.
3. Type “Z” is both 0s. This should not be included in Jaccard similarity at all.

If counted these types over *all* rows, that's $\text{sim} = \frac{X}{X+Y}$ ⁽¹¹⁾ ₍₁₁₎₊₍₀₁₎₊₍₁₀₎
 We're going to use permutations to *randomly* select a row, and only count the first row that's either X-type or Y-type.

hashing Example

Let's generate some permutations.

	banana	bandit	brand
0	1	1	1
1	0	1	0
2	0	0	1
3	0	1	0
6	0	0	1
8	1	1	0
13	1	0	10
16	0	1	1

1. Consider the permutation of rows: 8, 6, 2, 16, 0, 13, 1, 3
2. The resulting minhash values are:
 - 2.1 banana: 1 (b/c first 1 is in row 1)
 - 2.2 bandit: 1
 - 2.3 brand: 2 (b/c first 1 is in row 2)
3. Save those values into the new, *signature matrix*:

$$\begin{bmatrix} \text{banana} & \text{bandit} & \text{brand} \\ 1 & 1 & 2 \end{bmatrix}$$

first row of permutation *second part of permutation*

hashing Example

Let's generate some permutations.

	banana	bandit	brand
0	1	1	1
1	0	1	0
2	0	0	1
3	0	1	0
6	0	0	1
8	1	1	0
13	1	0	10
16	0	1	1

- Next, consider the permutation of rows: 6, 16, 3, 1, 8, 0, 13, 2
- The resulting minhash values are:
 - banana: 5
 - bandit: 2
 - brand: 1
- Finally let's also add the original permutation of rows: 3, 6, 16, 1, 8, 0, 13, 2, and the resulting minhash values...
 banana: 5, bandit: 1, brand: 2.
- Append those values into the *signature matrix*:

banana	bandit	brand
<u>1</u> 5	<u>1</u> 2	<u>2</u> 1
5	1	2

hashing Example

"Permutations work" *te*

The Jaccard similarity of the document is approximated by the Jaccard similarity of the rows of the *signature matrix*

	banana	bandit	brand
0	1	1	1
1	0	1	0
2	0	0	1
3	0	1	0
6	0	0	1
8	1	1	0
13	1	0	10
16	0	1	1

$$\begin{bmatrix} \text{banana} & \text{bandit} & \text{brand} \\ 1 & 1 & 2 \\ 5 & 2 & 1 \\ 5 & 1 & 2 \end{bmatrix}$$

(Note: In the original image, red and blue lines connect the values 1, 5, and 5 in the first column to the values 1, 2, and 1 in the second column, and 2, 1, and 2 in the third column, illustrating a permutation.)

So we have:

1. $\text{sim}(\text{banana}, \text{bandit}) \approx 1/3$
2. $\text{sim}(\text{banana}, \text{brand}) \approx 0/3$
3. $\text{sim}(\text{bandit}, \text{brand}) \approx 0/3$

which is uh... not great compared to the true values of (2/6, 1/6, 2/7). But in reality we'd do this many many times!

Permuting

So we're at this point:

1. Collapse documents into smaller matrices using shingles, and possible hashing to reduce the total number of rows.
2. Permuting the rows around allows us to save on memory. Instead of loading all the documents at once, we can create their *signatures*, or the first row with a "1" under a given permutation. We then apply the *same permutation* to other documents and compare *signatures*

But did this actually save memory? Now we have to save permutations of huge matrices to use the same permutations on new documents!

Worse... aren't truly random permutations pretty computationally expensive??



Not Permuting

It turns out we don't *actually have to permute*. Instead, another hash function allows us to **approximate** permutations. If we use a hash function we can randomly grab rows in a way that *looks* like a random permutation.

This is called a **universal hash**. For *random* integers a and b and a large prime number p (much greater than the total number of rows N ,

$$h_{a,b}(x) = ((a \cdot x + b) \mod p) \mod N$$

is a function that in practice *scans the rows* of the **characteristic matrix** in an approximately random order.

The full minhash:

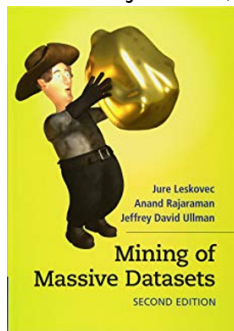
1. Step 1: pick n hash functions.
2. Step 2: initialize the *signature matrix* with all infinities. We'll replace the terms once we find "first 1's".
3. Step 3: For each row r of the characteristic matrix, compute $h_i(r)$ for that row for each and every one of the i hash functions. This represents "pick a random row" for n different permutations.
4. For each document or column x :
 - 4.1 If the char. matrix has a 0 in column c , row r : do nothing.
 - 4.2 If the char. matrix has a 1 in column c , row r , then take each of the hash functions and replace $sig(i, c)$ by $\min(sig(i, c), h_i(r))$

Why this last bit? We're looking for the *first row* of a 1, which will be the minimum time that it occurred.

Acknowledgments

We'll pick up with an example of this next time.. and start discussing clustering!

Some material is adapted/adopted from Mining of Massive Data Sets, by Jure Leskovec, Anand Rajaraman, Jeff Ullman (Stanford University) <http://www.mmds.org>



Special thanks to Tony Wong for sharing his original adaptation and adoption of slide material.