

Lecture 12: Locality-Sensitive Hashing and MinHash

Sam McCauley

April 26, 2020

Williams College

Introduction: Finding Similar Items

Back to Normal Inputs

- Today: no more streaming! Have all data available to us.

Back to Normal Inputs

- Today: no more streaming! Have all data available to us.
- But data is still big!

Back to Normal Inputs

- Today: no more streaming! Have all data available to us.
- But data is still big!
 - In particular: high-dimensional

Back to Normal Inputs

- Today: no more streaming! Have all data available to us.
- But data is still big!
 - In particular: high-dimensional
 - Table with many columns

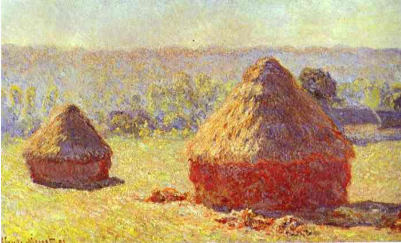
Back to Normal Inputs

- Today: no more streaming! Have all data available to us.
- But data is still big!
 - In particular: high-dimensional
 - Table with many columns
 - For each netflix user, what movies have they seen

Back to Normal Inputs

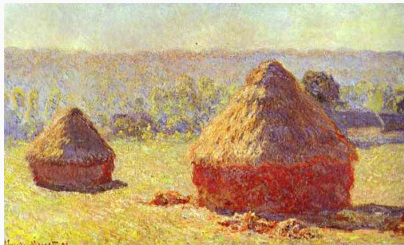
- Today: no more streaming! Have all data available to us.
- But data is still big!
 - In particular: high-dimensional
 - Table with many columns
 - For each netflix user, what movies have they seen
- Goal: solve a difficult, but important, problem

Finding Similar Pair



- Given a set of objects

Finding Similar Pair



- Given a set of objects
- Find the most similar pair of objects in the set

Why Find Similar Objects?

- Find similar news articles for user suggestions.

Why Find Similar Objects?

- Find similar news articles for user suggestions.
- Similar music: Spotify suggests music by finding similar users, and selecting what they listen to

Why Find Similar Objects?

- Find similar news articles for user suggestions.
- Similar music: Spotify suggests music by finding similar users, and selecting what they listen to
- Machine learning in general (training, evaluation, actual algorithms, etc.)

Why Find Similar Objects?

- Find similar news articles for user suggestions.
- Similar music: Spotify suggests music by finding similar users, and selecting what they listen to
- Machine learning in general (training, evaluation, actual algorithms, etc.)
- Data deduplication, etc.

Why Find Similar Objects?

- Find similar news articles for user suggestions.
- Similar music: Spotify suggests music by finding similar users, and selecting what they listen to
- Machine learning in general (training, evaluation, actual algorithms, etc.)
- Data deduplication, etc.
- “Give me a similar pair in this dataset” is a common query!

Strategies for Similarity Search

First attempt: 1-dimensional data

92
44
7
65
60
23
80
67

- Given a list of numbers

First attempt: 1-dimensional data

92
44
7
65
60
23
80
67

- Given a list of numbers
- “Similarity” is the difference between them

First attempt: 1-dimensional data

92
44
7
65
60
23
80
67

- Given a list of numbers
- “Similarity” is the difference between them
- How can we find the closest numbers (i.e. ones with smallest difference)?

First attempt: 1-dimensional data

- How efficiently can we do this?

First attempt: 1-dimensional data

7
23
44
60
65
67
80
92

- How efficiently can we do this?
- Step 1: Sort!

First attempt: 1-dimensional data

7
23
44
60
65
67
80
92

- How efficiently can we do this?
- Step 1: Sort!
- Step 2: Scan through list, find most similar adjacent elements.

First attempt: 1-dimensional data

7
23
44
60
65
67
80
92

- How efficiently can we do this?
- Step 1: Sort!
- Step 2: Scan through list, find most similar adjacent elements.
- $O(n \log n)$ time

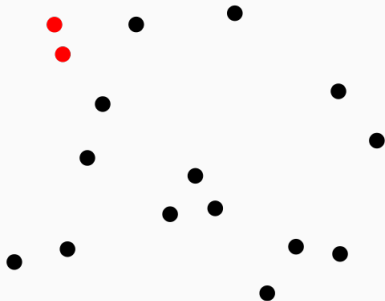
First attempt: 1-dimensional data

Aside: can we do better? Yes, there's a clever $O(n)$ algorithm based on sampling.

60
65
67
80
92

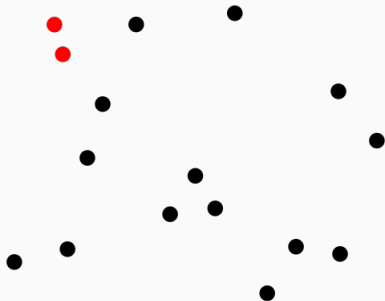
- Step 2: Scan through list, find most similar adjacent elements.
- $O(n \log n)$ time

Two-dimensional Data?



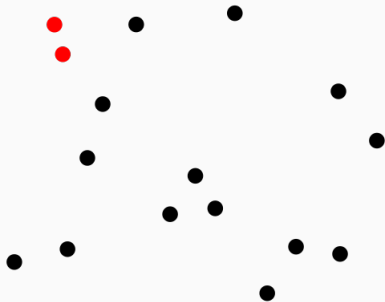
- You likely saw this in CS 256.

Two-dimensional Data?



- You likely saw this in CS 256.
- Divide and conquer, $O(n \log n)$ time.

Two-dimensional Data?



- You likely saw this in CS 256.
- Divide and conquer, $O(n \log n)$ time.
- (Again, possible in $O(n)$)

What About Higher Dimensions?

- We want VERY high dimensions (millions)

What About Higher Dimensions?

- We want VERY high dimensions (millions)
- Songs listened to, movies watched, image tags, etc.

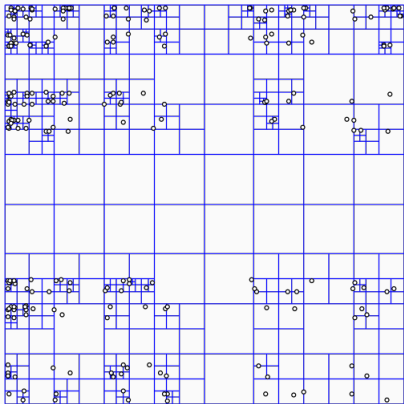
What About Higher Dimensions?

- We want VERY high dimensions (millions)
- Songs listened to, movies watched, image tags, etc.
- Words that appear in a book, k-grams that appear in a DNA sequence

What About Higher Dimensions?

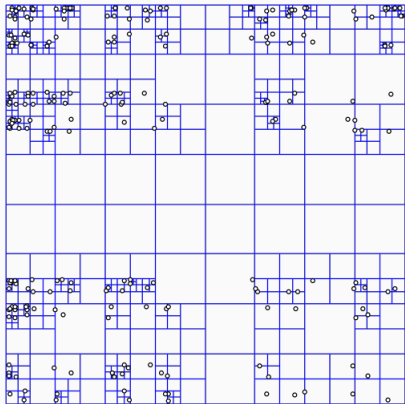
- We want VERY high dimensions (millions)
- Songs listened to, movies watched, image tags, etc.
- Words that appear in a book, k-grams that appear in a DNA sequence
- Classic options: quad trees, kd trees

How Efficient are High-dimensional Algorithms?



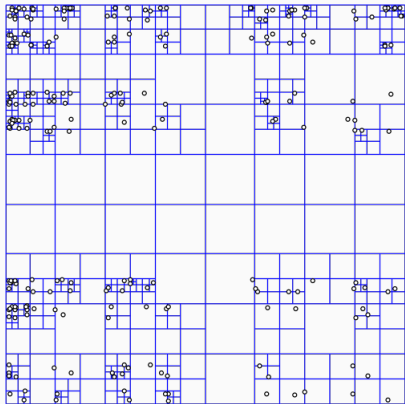
- $O(n \log n)$ for constant dimensions

How Efficient are High-dimensional Algorithms?



- $O(n \log n)$ for constant dimensions
- But: exponential in dimension!

How Efficient are High-dimensional Algorithms?



- $O(n \log n)$ for constant dimensions
- But: exponential in dimension!
- Worse than trying all pairs if $> \log n$ dimensions

Curse of Dimensionality

- Many problems have running time exponential in the dimension of the objects.

Curse of Dimensionality

- Many problems have running time exponential in the dimension of the objects.
- Well-known phenomenon

Curse of Dimensionality

- Many problems have running time exponential in the dimension of the objects.
- Well-known phenomenon
- Applies to similarity search, machine learning, combinatorics

Avoiding the Curse of Dimensionality

- Today we're talking about how to get efficient algorithms for *arbitrarily* large dimensions.

Avoiding the Curse of Dimensionality

- Today we're talking about how to get efficient algorithms for *arbitrarily* large dimensions.
- Linear cost in terms of dimension (but expensive in terms of the problem size).

Avoiding the Curse of Dimensionality

- Today we're talking about how to get efficient algorithms for *arbitrarily* large dimensions.
- Linear cost in terms of dimension (but expensive in terms of the problem size).
- Two tools to get us there:

Avoiding the Curse of Dimensionality

- Today we're talking about how to get efficient algorithms for *arbitrarily* large dimensions.
- Linear cost in terms of dimension (but expensive in terms of the problem size).
- Two tools to get us there:
 - Assume that the close pair is much closer than any other (*approximate* closest pair)

Avoiding the Curse of Dimensionality

- Today we're talking about how to get efficient algorithms for *arbitrarily* large dimensions.
- Linear cost in terms of dimension (but expensive in terms of the problem size).
- Two tools to get us there:
 - Assume that the close pair is much closer than any other (*approximate* closest pair)

We'll come back to this later

Avoiding the Curse of Dimensionality

- Today we're talking about how to get efficient algorithms for *arbitrarily* large dimensions.
- Linear cost in terms of dimension (but expensive in terms of the problem size).
- Two tools to get us there:
 - Assume that the close pair is much closer than any other (*approximate* closest pair)
 - Use hashing! ...A special kind of hashing

We'll come back to this later

Locality-Sensitive Hashing

- Normally, hashing *spreads out* elements.

Locality-Sensitive Hashing

- Normally, hashing *spreads out* elements.
- This is key to hashing: no matter how clustered my data begins, I wind up with a nicely-distributed hash table

Locality-Sensitive Hashing

- Normally, hashing *spreads out* elements.
- This is key to hashing: no matter how clustered my data begins, I wind up with a nicely-distributed hash table
- Locality-sensitive hashing tries to hash similar items together

Locality-Sensitive Hashing: Formal Definition

- Needs a similarity threshold r , an approximation factor $c < 1$

Locality-Sensitive Hashing: Formal Definition

- Needs a similarity threshold r , an approximation factor $c < 1$
- Two guarantees:

Locality-Sensitive Hashing: Formal Definition

- Needs a similarity threshold r , an approximation factor $c < 1$
- Two guarantees:
 - If two items x and y have similarity $\geq r$, $h(x) = h(y)$ with probability at least p_1 .

Locality-Sensitive Hashing: Formal Definition

- Needs a similarity threshold r , an approximation factor $c < 1$
- Two guarantees:
 - If two items x and y have similarity $\geq r$, $h(x) = h(y)$ with probability at least p_1 .
 - If two items x and y have similarity $\leq cr$, $h(x) = h(y)$ with probability at most p_2 .

Locality-Sensitive Hashing: Formal Definition

- Needs a similarity threshold r , an approximation factor $c < 1$
- Two guarantees:
 - If two items x and y have similarity $\geq r$, $h(x) = h(y)$ with probability at least p_1 .
 - If two items x and y have similarity $\leq cr$, $h(x) = h(y)$ with probability at most p_2 .
- High level: close items are likely to collide. Far items are unlikely to collide.

Locality-Sensitive Hashing: Formal Definition

- Needs a similarity threshold r , an approximation factor $c < 1$
- Two guarantees:
 - If two items x and y have similarity $\geq r$, $h(x) = h(y)$ with probability at least p_1 .
 - If two items x and y have similarity $\leq cr$, $h(x) = h(y)$ with probability at most p_2 .
- High level: close items are likely to collide. Far items are unlikely to collide.
- Generally want p_2 to be about $1/n$; then we get a normal hash table for far (i.e. distance $\geq cr$) elements.

Why Locality-Sensitive Hashing Helps

	(101, 37, 65) (103,37,64)	(91,84,3)		(100,18,79)
0	1	2	3	4

Ideally, close items hash to the same bucket.

Issue: Low probability of success!

- If we have $p_2 = 1/n$, then p_1 is usually very small.

Issue: Low probability of success!

We'll put numbers
on this later

- If we have $p_2 = 1/n$, then p_1 is usually very small.
- How can we increase this probability?

Issue: Low probability of success!

- If we have $p_2 = 1/n$, then p_1 is usually very small.
- How can we increase this probability?
- Repetitions! Maintain many hash tables, each with a different locality-sensitive hash function, and try all of them.

LSH with Repetitions

(101, 37, 65)	(103,37,64)	(91,84,3)	(100,18,79)	
0	1	2	3	4

	(101,37,65) (103,37,64)	(91,84,3)		(100,18,79)
0	1	2	3	4

(101, 37, 65)		(103,37,64)		(91,84,3) (100,18,79)
0	1	2	3	4

Similarity

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?
- Images in machine learning: often Euclidean distance (the distance we're familiar with on a day-to-day basis)

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?
- Images in machine learning: often Euclidean distance (the distance we're familiar with on a day-to-day basis)
- What about *sets*?

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?
- Images in machine learning: often Euclidean distance (the distance we're familiar with on a day-to-day basis)
- What about *sets*?
 - Songs listened to by a user

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?
- Images in machine learning: often Euclidean distance (the distance we're familiar with on a day-to-day basis)
- What about *sets*?
 - Songs listened to by a user
 - Movies watched by a user

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?
- Images in machine learning: often Euclidean distance (the distance we're familiar with on a day-to-day basis)
- What about *sets*?
 - Songs listened to by a user
 - Movies watched by a user
 - Human-generated tags given to an image

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?
- Images in machine learning: often Euclidean distance (the distance we're familiar with on a day-to-day basis)
- What about *sets*?
 - Songs listened to by a user
 - Movies watched by a user
 - Human-generated tags given to an image
 - Words that appear in a document

What Do We Mean by “Similar”?

- How can we measure the similarity of objects?
- Images in machine learning: often Euclidean distance (the distance we're familiar with on a day-to-day basis)
- What about *sets*?
 - Songs listened to by a user
 - Movies watched by a user
 - Human-generated tags given to an image
 - Words that appear in a document
- Need a way to measure set similarity

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Khalid
Khalid	Drake
Drake	Travis Scott
Travis Scott	

- When are two sets similar?

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Khalid
Khalid	Drake
Drake	Travis Scott
Travis Scott	

- When are two sets similar?
- Let's look at our two sets. Similar if they have a lot of *overlap*

Set Similarity

Very similar!

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Khalid
Khalid	Drake
Drake	Travis Scott
Travis Scott	

- When are two sets similar?
- Let's look at our two sets. Similar if they have a lot of *overlap*
- I.e. : lots of artists in common, compared to total artists in either list

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	
Khalid	
Drake	
Travis Scott	

- When are two sets similar?
- Let's look at our two sets. Similar if they have a lot of *overlap*
- I.e. : lots of artists in common, compared to total artists in either list

Set Similarity

Not very similar!

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	
Khalid	
Drake	
Travis Scott	

- When are two sets similar?
- Let's look at our two sets. Similar if they have a lot of *overlap*
- I.e. : lots of artists in common, compared to total artists in either list

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Ed Sheeran
Khalid	Drake
Drake	Travis Scott
Travis Scott	Taylor Swift

- When are two sets similar?
- Let's look at our two sets. Similar if they have a lot of *overlap*
- I.e. : lots of artists in common, compared to total artists in either list

Set Similarity

Moderately similar

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Ed Sheeran
Khalid	Drake
Drake	Travis Scott
Travis Scott	Taylor Swift

- When are two sets similar?
- Let's look at our two sets. Similar if they have a lot of *overlap*
- I.e. : lots of artists in common, compared to total artists in either list

Jaccard Similarity

Jaccard Similarity

- Similarity measure for sets A and B

Jaccard Similarity

- Similarity measure for sets A and B
- Defined as:

$$\frac{|A \cap B|}{|A \cup B|}$$

Jaccard Similarity

- Similarity measure for sets A and B

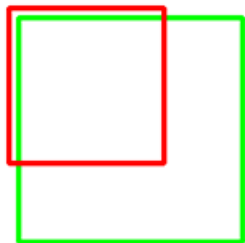
- Defined as:

$$\frac{|A \cap B|}{|A \cup B|}$$

- Intuitively: what fraction of these sets overlaps?

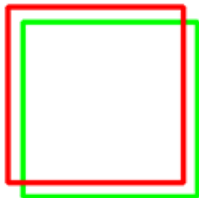
Jaccard Similarity Intuition 1

IoU: 0.4034



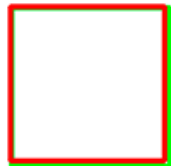
Poor

IoU: 0.7330



Good

IoU: 0.9264



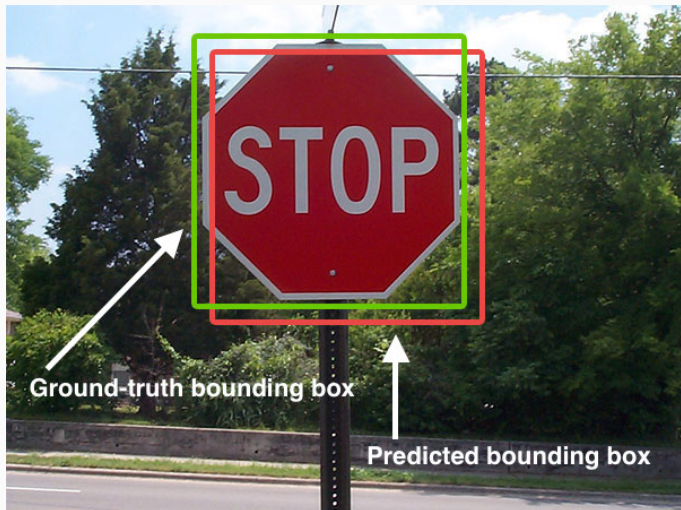
Excellent

Jaccard Similarity Intuition 2

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



Image Search Example



Jaccard Example 1

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Khalid
Khalid	Drake
Drake	Travis Scott
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.

Jaccard Example 1

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Khalid
Khalid	Drake
Drake	Travis Scott
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 4$

Jaccard Example 1

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Khalid
Khalid	Drake
Drake	Travis Scott
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 4$
- $|A \cup B| = 5$

Jaccard Example 1

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Khalid
Khalid	Drake
Drake	Travis Scott
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 4$
- $|A \cup B| = 5$
- Jaccard Similarity: $4/5 = .8$

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	
Khalid	
Drake	
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	
Khalid	
Drake	
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 1$

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	
Khalid	
Drake	
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 1$
- $|A \cup B| = 5$

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	
Khalid	
Drake	
Travis Scott	

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 1$
- $|A \cup B| = 5$
- Jaccard Similarity: $1/5 = .2$

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Ed Sheeran
Khalid	Drake
Drake	Travis Scott
Travis Scott	Taylor Swift

- Similarity: $|A \cap B| / |A \cup B|$.

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Ed Sheeran
Khalid	Drake
Drake	Travis Scott
Travis Scott	Taylor Swift

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 3$

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Ed Sheeran
Khalid	Drake
Drake	Travis Scott
Travis Scott	Taylor Swift

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 3$
- $|A \cup B| = 7$

Set Similarity

User 1	User 2
Post Malone	Ariana Grande
Ariana Grande	Ed Sheeran
Khalid	Drake
Drake	Travis Scott
Travis Scott	Taylor Swift

- Similarity: $|A \cap B| / |A \cup B|$.
- $|A \cap B| = 3$
- $|A \cup B| = 7$
- Jaccard Similarity:
 $3/7 = 0.428$

Jaccard Similarity: Properties

- Works on sets (each dimension is *binary*—an item is in the set, or not in the set)

Jaccard Similarity: Properties

- Works on sets (each dimension is *binary*—an item is in the set, or not in the set)
- Always gives a number between 0 and 1

Jaccard Similarity: Properties

- Works on sets (each dimension is *binary*—an item is in the set, or not in the set)
- Always gives a number between 0 and 1
- 1 means identical, 0 means no items in common

Jaccard Similarity: Properties

- Works on sets (each dimension is *binary*—an item is in the set, or not in the set)
- Always gives a number between 0 and 1
- 1 means identical, 0 means no items in common
- Jaccard similarity ignores items not in either set. So we learn nothing if neither of us like an artist. (Is this good?)

Jaccard Similarity: Properties

- Works on sets (each dimension is *binary*—an item is in the set, or not in the set)
- Always gives a number between 0 and 1
- 1 means identical, 0 means no items in common
- Jaccard similarity ignores items not in either set. So we learn nothing if neither of us like an artist. (Is this good?)
- Still works if one list is much longer than the other. Generally, they'll have small overlap

Locality-Sensitive Hash for Jaccard Similarity

- Want: items with high Jaccard Similarity are likely to hash together

Locality-Sensitive Hash for Jaccard Similarity

- Want: items with high Jaccard Similarity are likely to hash together
- Items with low Jaccard Similarity are UNlikely to hash together

Locality-Sensitive Hash for Jaccard Similarity

- Want: items with high Jaccard Similarity are likely to hash together
- Items with low Jaccard Similarity are UNlikely to hash together
- Classic method: MinHash

MinHash

- Developed by Andrei Broder in 1997 while working at AltaVista

- Developed by Andrei Broder in 1997 while working at AltaVista
- (AltaVista was probably the most popular search engine before Google, they wanted to detect similar web pages to eliminate them from search results)

- Developed by Andrei Broder in 1997 while working at AltaVista
- (AltaVista was probably the most popular search engine before Google, they wanted to detect similar web pages to eliminate them from search results)
- Now used for similarity search, database joins, clustering—LOTS of things.

Bit Vectors

- Can represent any set as a vector of bits

Bit Vectors

- Can represent any set as a vector of bits
- Each bit is an item. "1" means that that item is in the set, "0" means it's not

Bit Vectors

- Can represent any set as a vector of bits
- Each bit is an item. "1" means that that item is in the set, "0" means it's not
- So if I'm keeping track of different people's favorite colors, my universe may be {red, yellow, blue, green, purple, orange}

Bit Vectors

- Can represent any set as a vector of bits
- Each bit is an item. "1" means that that item is in the set, "0" means it's not
- So if I'm keeping track of different people's favorite colors, my universe may be {red, yellow, blue, green, purple, orange}
- If someone likes red and blue, we can store that information as 101000.

Bit Vectors

- Can represent any set as a vector of bits
- Each bit is an item. "1" means that that item is in the set, "0" means it's not
- So if I'm keeping track of different people's favorite colors, my universe may be {red, yellow, blue, green, purple, orange}
- If someone likes red and blue, we can store that information as 101000.
- Effective if universe is smallish; use a list for larger universe

Bit Vectors: Jaccard Similarity

- How can we determine $A \cap B$?

Bit Vectors: Jaccard Similarity

- How can we determine $A \cap B$?
 - This is exactly $A \& B$ in C-style notation

Bit Vectors: Jaccard Similarity

- How can we determine $A \cap B$?
 - This is exactly $A \& B$ in C-style notation
- What about $A \cup B$?

Bit Vectors: Jaccard Similarity

- How can we determine $A \cap B$?
 - This is exactly $A \& B$ in C-style notation
- What about $A \cup B$?
 - This is exactly $A | B$ in C-style notation

Bit Vectors: Jaccard Similarity

- How can we determine $A \cap B$?
 - This is exactly $A \& B$ in C-style notation
- What about $A \cup B$?
 - This is exactly $A | B$ in C-style notation
- We want the size of these sets—need to count the number of 1s in $A \& B$, or $A | B$.

- The hash consists of an *permutation* of all possible items in the universe

128 in the
assignment

- The hash consists of an *permutation* of all possible items in the universe
- To hash a set A : find the first item of A in the order given by the permutation. That item is the hash value!

MinHash example

- Let's stick with favorite colors, out of {red, yellow, blue, green, purple, orange}

MinHash example

- Let's stick with favorite colors, out of {red, yellow, blue, green, purple, orange}
- To hash, we randomly permute them. Let's say our current hash is given by the permutation (blue, orange, green, purple, red, yellow)

MinHash example

- Let's stick with favorite colors, out of {red, yellow, blue, green, purple, orange}
- To hash, we randomly permute them. Let's say our current hash is given by the permutation (blue, orange, green, purple, red, yellow)
- First set is 101000 (same as {red, blue}). blue is in the set, so the hash value is blue.

MinHash example

- Let's stick with favorite colors, out of {red, yellow, blue, green, purple, orange}
- To hash, we randomly permute them. Let's say our current hash is given by the permutation (blue, orange, green, purple, red, yellow)
- First set is 101000 (same as {red, blue}). blue is in the set, so the hash value is blue.
- Second set is 110010 (we could also write {red, yellow, purple}). blue is not in the set; nor is orange; nor is green. purple is, so purple is the hash value

MinHash for Bit Vectors

- On the assignment, have bit vectors of length 128

MinHash for Bit Vectors

- On the assignment, have bit vectors of length 128
- To get a hash function, we need a random permutation of the indices of these bits. That is to say, a random permutation of $\{0, 1, 2, \dots, 127\}$

MinHash for Bit Vectors

- On the assignment, have bit vectors of length 128
- To get a hash function, we need a random permutation of the indices of these bits. That is to say, a random permutation of $\{0, 1, 2, \dots, 127\}$
- To hash an item x , go through the random permutation. Find the first index i in the list such that the i th bit of x is 1.

MinHash for Bit Vectors

- On the assignment, have bit vectors of length 128
- To get a hash function, we need a random permutation of the indices of these bits. That is to say, a random permutation of $\{0, 1, 2, \dots, 127\}$
- To hash an item x , go through the random permutation. Find the first index i in the list such that the i th bit of x is 1.
- Let's say $x = 10011001$, and the permutation is $(1, 5, 2, 0, 7, 6, 4, 3)$.

MinHash for Bit Vectors

- On the assignment, have bit vectors of length 128
- To get a hash function, we need a random permutation of the indices of these bits. That is to say, a random permutation of $\{0, 1, 2, \dots, 127\}$
- To hash an item x , go through the random permutation. Find the first index i in the list such that the i th bit of x is 1.
- Let's say $x = 10011001$, and the permutation is $(1, 5, 2, 0, 7, 6, 4, 3)$.

For the sake
of space let's
do 8 bits

MinHash for Bit Vectors

- On the assignment, have bit vectors of length 128
- To get a hash function, we need a random permutation of the indices of these bits. That is to say, a random permutation of $\{0, 1, 2, \dots, 127\}$
- To hash an item x , go through the random permutation. Find the first index i in the list such that the i th bit of x is 1.
- Let's say $x = 10011001$, and the permutation is $(1, 5, 2, 0, 7, 6, 4, 3)$.
- Then the hash of x is 5.

Analysis of Basic MinHash

- What is the probability that $h(A) = h(B)$?

Analysis

- What is the probability that $h(A) = h(B)$?
- Let's look at the permutation that defines h . We can ignore any item that is not in A or B .

Analysis

- What is the probability that $h(A) = h(B)$?
- Let's look at the permutation that defines h . We can ignore any item that is not in A or B .
- Look at the first index in the permutation that is in A or B (i.e. it is in $A \cup B$)

Analysis

- What is the probability that $h(A) = h(B)$?
- Let's look at the permutation that defines h . We can ignore any item that is not in A or B .
- Look at the first index in the permutation that is in A or B (i.e. it is in $A \cup B$)
 - If this index is in *both* A and B , then $h(A) = h(B)$

Analysis

- What is the probability that $h(A) = h(B)$?
- Let's look at the permutation that defines h . We can ignore any item that is not in A or B .
- Look at the first index in the permutation that is in A or B (i.e. it is in $A \cup B$)
 - If this index is in *both* A and B , then $h(A) = h(B)$
 - If this index is in only one of A or B , then $h(A) \neq h(B)$

Analysis

- What is the probability that $h(A) = h(B)$?
- Let's look at the permutation that defines h . We can ignore any item that is not in A or B .
- Look at the first index in the permutation that is in A or B (i.e. it is in $A \cup B$)
 - If this index is in *both* A and B , then $h(A) = h(B)$
 - If this index is in only one of A or B , then $h(A) \neq h(B)$
- Any index in $A \cup B$ is equally likely to be first. If the index is in $A \cap B$, they hash together; otherwise they do not

Analysis

- What is the probability that $h(A) = h(B)$?
- Let's look at the permutation that defines h . We can ignore any item that is not in A or B .
- Look at the first index in the permutation that is in A or B (i.e. it is in $A \cup B$)
 - If this index is in *both* A and B , then $h(A) = h(B)$
 - If this index is in only one of A or B , then $h(A) \neq h(B)$
- Any index in $A \cup B$ is equally likely to be first. If the index is in $A \cap B$, they hash together; otherwise they do not
- Therefore: probability of hashing together is $|A \cap B|/|A \cup B|$.

MinHash as an LSH

- This means MinHash is an LSH!

MinHash as an LSH

- This means MinHash is an LSH!
- If two items have similarity at least r , they collide with probability at least r

MinHash as an LSH

- This means MinHash is an LSH!
- If two items have similarity at least r , they collide with probability at least r
- If two items have similarity at most cr , they collide with probability at most cr

Analysis: Phrased as bit vectors

- What is the probability that $h(A) = h(B)$?
- Let's look at the permutation that defines h . We can ignore any index that is 0 in both A and B .
- Look at the first index in the permutation that is 1 in A or B
 - If this index is in *both* A and B , then $h(A) = h(B)$
 - If this index is in only one of A or B , then $h(A) \neq h(B)$
- Any index that is 1 in $A|B$ is equally likely to be first. If the index is in $A\&B$, they hash together; otherwise they do not
- Therefore: probability of hashing together is
(number of 1s in $A\&B$)/(number of 1s in $A|B$).

Analysis Example

- Let's say we have $A = \{\text{red, blue, green}\}$ and $B = \{\text{red, orange, purple, green}\}$.

Analysis Example

- Let's say we have $A = \{\text{red, blue, green}\}$ and $B = \{\text{red, orange, purple, green}\}$.
- When do A and B hash together?

Analysis Example

- Let's say we have $A = \{\text{red, blue, green}\}$ and $B = \{\text{red, orange, purple, green}\}$.
- When do A and B hash together?
- If red or green appears before blue, orange, and purple then they hash together

Analysis Example

- Let's say we have $A = \{\text{red, blue, green}\}$ and $B = \{\text{red, orange, purple, green}\}$.
- When do A and B hash together?
- If red or green appears before blue, orange, and purple then they hash together
- If blue or orange or purple appear before red and green, then they hash together

Analysis Example

- Let's say we have $A = \{\text{red, blue, green}\}$ and $B = \{\text{red, orange, purple, green}\}$.
- When do A and B hash together?
- If red or green appears before blue, orange, and purple then they hash together
- If blue or orange or purple appear before red and green, then they hash together
- Probability that red or green is first out of $\{\text{red, blue, green, orange, purple}\}$ is $2/5$.

Analysis Example

- Let's say we have $A = \{\text{red, blue, green}\}$ and $B = \{\text{red, orange, purple, green}\}$.
- When do A and B hash together?
- If red or green appears before blue, orange, and purple then they hash together
- If blue or orange or purple appear before red and green, then they hash together
- Probability that red or green is first out of $\{\text{red, blue, green, orange, purple}\}$ is $2/5$.
- Therefore, A and B hash together with probability $2/5$.

Making Sure We Find the Close Pair

- To find the close pair, compare all pairs of items that hash to the same value

Making Sure We Find the Close Pair

- To find the close pair, compare all pairs of items that hash to the same value
 - (We'll talk about how to do this later—it's similar to MiniMidterm 1)

Making Sure We Find the Close Pair

- To find the close pair, compare all pairs of items that hash to the same value
 - (We'll talk about how to do this later—it's similar to MiniMidterm 1)
- Let's say our close pair has similarity .5. How many times do we need to repeat?

Making Sure We Find the Close Pair

- To find the close pair, compare all pairs of items that hash to the same value
 - (We'll talk about how to do this later—it's similar to MiniMidterm 1)
- Let's say our close pair has similarity $.5$. How many times do we need to repeat?
- Each repetition has the close pair in the same bucket with probability $.5$. So need 2 repetitions in expectation.

An Aside on Expectation

Lemma 1

If a random process succeeds with probability p , then in expectation it takes $1/p$ iterations of the process before success.

Examples:

An Aside on Expectation

Lemma 1

If a random process succeeds with probability p , then in expectation it takes $1/p$ iterations of the process before success.

Examples:

- It takes two coin flips in expectation before we see a heads

An Aside on Expectation

Lemma 1

If a random process succeeds with probability p , then in expectation it takes $1/p$ iterations of the process before success.

Examples:

- It takes two coin flips in expectation before we see a heads
- We need to roll a 6-sided die 6 times before we see (say) a three.

An Aside on Expectation

Lemma 1

If a random process succeeds with probability p , then in expectation it takes $1/p$ iterations of the process before success.

Examples:

- It takes two coin flips in expectation before we see a heads
- We need to roll a 6-sided die 6 times before we see (say) a three.

Proof:

$$\sum_{i=1}^{\infty} ip(1-p)^{i-1} = \frac{p}{(1-(1-p))^2} = \frac{1}{p}$$

Concatenations and Repetitions

Problems with this Approach

- Buckets are really big!! (After all, lots of items are pretty likely to have a given bit set.)

Problems with this Approach

- Buckets are really big!! (After all, lots of items are pretty likely to have a given bit set.)
- How can we decrease the probability that items hash together?

Problems with this Approach

- Buckets are really big!! (After all, lots of items are pretty likely to have a given bit set.)
- How can we decrease the probability that items hash together?
- Answer: concatenate multiple hashes together.

Concatenating Hashes

- Rather than one hash h , concatenate k independent hashes h_1, h_2, \dots, h_k , each with its own permutation P_1, P_2, \dots, P_k .

Concatenating Hashes

- Rather than one hash h , concatenate k independent hashes h_1, h_2, \dots, h_k , each with its own permutation P_1, P_2, \dots, P_k .
- To hash an item: repeat the process of searching through the permutation for each hash. Then concatenate the results together (can just use string concatenation)

Concatenating Hashes

- Rather than one hash h , concatenate k independent hashes h_1, h_2, \dots, h_k , each with its own permutation P_1, P_2, \dots, P_k .
- To hash an item: repeat the process of searching through the permutation for each hash. Then concatenate the results together (can just use string concatenation)
- How does this affect the probability for sets A and B ?

Concatenating Hashes

- Rather than one hash h , concatenate k independent hashes h_1, h_2, \dots, h_k , each with its own permutation P_1, P_2, \dots, P_k .
- To hash an item: repeat the process of searching through the permutation for each hash. Then concatenate the results together (can just use string concatenation)
- How does this affect the probability for sets A and B ?
 - For each of the k independent hashes, A and B collide with probability $|A \cap B|/|A \cup B|$.

Concatenating Hashes

- Rather than one hash h , concatenate k independent hashes h_1, h_2, \dots, h_k , each with its own permutation P_1, P_2, \dots, P_k .
- To hash an item: repeat the process of searching through the permutation for each hash. Then concatenate the results together (can just use string concatenation)
- How does this affect the probability for sets A and B ?
 - For each of the k independent hashes, A and B collide with probability $|A \cap B|/|A \cup B|$.
 - We only obtain the same concatenated hashes if *all* of the hashes are the same.

Concatenating Hashes

- Rather than one hash h , concatenate k independent hashes h_1, h_2, \dots, h_k , each with its own permutation P_1, P_2, \dots, P_k .
- To hash an item: repeat the process of searching through the permutation for each hash. Then concatenate the results together (can just use string concatenation)
- How does this affect the probability for sets A and B ?
 - For each of the k independent hashes, A and B collide with probability $|A \cap B|/|A \cup B|$.
 - We only obtain the same concatenated hashes if *all* of the hashes are the same.
 - They are independent, so we can multiply to obtain probability $(|A \cap B|/|A \cup B|)^k$ of A and B colliding.

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash A .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash A .
 - First hash: red is in A .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash A .
 - First hash: red is in A .
 - Second hash: orange not in A , nor is green. Blue is in A .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash A .
 - First hash: red is in A .
 - Second hash: orange not in A , nor is green. Blue is in A .
 - Third hash: red is in A .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash A .
 - First hash: red is in A .
 - Second hash: orange not in A , nor is green. Blue is in A .
 - Third hash: red is in A .
- Concatenating, we have $h(A) = \text{redbluered}$

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash B .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash B .
 - First hash: red is in B .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash B .
 - First hash: red is in B .
 - Second hash: orange is in B .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash B .
 - First hash: red is in B .
 - Second hash: orange is in B .
 - Third hash: red is in B .

Concatenation Example

- Let's say we have $A = \{\text{red, blue}\}$ and $B = \{\text{red, orange}\}$, and $k = 3$.
- $P_1 = \{\text{red, green, blue, orange}\}, P_2 = \{\text{orange, green, blue, red}\}, P_3 = \{\text{red, green, blue, orange}\}$
- Let's hash B .
 - First hash: red is in B .
 - Second hash: orange is in B .
 - Third hash: red is in B .
- Concatenating, we have $h(B) = \text{redorangered}$

Putting it all Together

- For each hash table, we concatenate k hashes.

Putting it all Together

- For each hash table, we concatenate k hashes.
- Need to repeat all of that multiple times until we find the close pair (let's say we repeat R times)

Putting it all Together

- For each hash table, we concatenate k hashes.
- Need to repeat all of that multiple times until we find the close pair (let's say we repeat R times)
- So: overall need kR permutations

Putting it all Together

- For each hash table, we concatenate k hashes.
- Need to repeat all of that multiple times until we find the close pair (let's say we repeat R times)
- So: overall need kR permutations
- What kind of values work for k and R ?

Putting it Together: Analysis

- Let's say we have a set of n items x_1, \dots, x_n

Putting it Together: Analysis

- Let's say we have a set of n items x_1, \dots, x_n
- The close pair of items has Jaccard similarity $3/4$

Putting it Together: Analysis

- Let's say we have a set of n items x_1, \dots, x_n
- The close pair of items has Jaccard similarity $3/4$
- Every other pair of items has similarity $1/3$

Putting it Together: Analysis

- Let's say we have a set of n items x_1, \dots, x_n
- The close pair of items has Jaccard similarity $3/4$
- Every other pair of items has similarity $1/3$
- How should we set k ? How many repetitions R is it likely to take?

Putting it Together: Analysis (Finding k)

- Non-similar pairs have similarity $1/3$

Putting it Together: Analysis (Finding k)

- Non-similar pairs have similarity $1/3$
- We want buckets to be small (have $O(1)$ size)

Putting it Together: Analysis (Finding k)

- Non-similar pairs have similarity $1/3$
- We want buckets to be small (have $O(1)$ size)
- Look at an element x_i . The expected size of the bucket of x_i is $\sum_{j \neq i} (1/3)^k$ (since x_i and any x_j with $j \neq i$ share a hash value with probability $1/3^k$)

Putting it Together: Analysis (Finding k)

- Non-similar pairs have similarity $1/3$
- We want buckets to be small (have $O(1)$ size)
- Look at an element x_i . The expected size of the bucket of x_i is $\sum_{j \neq i} (1/3)^k$ (since x_i and any x_j with $j \neq i$ share a hash value with probability $1/3^k$)
- We can then solve $(n - 1)(1/3)^k = 1$ to get $k = \log_3 n - 1$.

Putting it Together: Analysis (Predicting R)

- The similar pair has Jaccard similarity .75

Putting it Together: Analysis (Predicting R)

- The similar pair has Jaccard similarity $.75$
- So they are in the same bucket with probability $(.75)^k$

Putting it Together: Analysis (Predicting R)

- The similar pair has Jaccard similarity .75
- So they are in the same bucket with probability $(.75)^k$
- We have $k = \log_3 n - 1$. So...we need to do some algebra

Putting it Together: Analysis (Predicting R)

- The similar pair has Jaccard similarity .75
- So they are in the same bucket with probability $(.75)^k$
- We have $k = \log_3 n - 1$. So...we need to do some algebra
- $(.75)^{\log_3 n - 1} = 2^{\log(n-1) \log(3/4) / \log(3)} = (n - 1)^{\log(3/4) / \log(3)} \approx 1/n^{.26}$

Putting it Together: Analysis (Predicting R)

- The similar pair has Jaccard similarity .75
- So they are in the same bucket with probability $(.75)^k$
- We have $k = \log_3 n - 1$. So...we need to do some algebra
- $(.75)^{\log_3 n - 1} = 2^{\log(n-1) \log(3/4) / \log(3)} = (n - 1)^{\log(3/4) / \log(3)} \approx 1/n^{.26}$
- So we expect about $R = n^{.26}$ repetitions. That's a lot!

Putting it Together: Analysis (Predicting R)

- The similar pair has Jaccard similarity .75
- So they are in the same bucket with probability $(.75)^k$
- We have $k = \log_3 n - 1$. So...we need to do some algebra
- $(.75)^{\log_3 n - 1} = 2^{\log(n-1) \log(3/4) / \log(3)} = (n - 1)^{\log(3/4) / \log(3)} \approx 1/n^{.26}$
- So we expect about $R = n^{.26}$ repetitions. That's a lot!
- But it's essentially the best we know how to do.

Practical MinHash Considerations

So many Permutations!

- OK, so kR repetitions is a LOT of preprocessing, and a lot of random number generation

So many Permutations!

- OK, so kR repetitions is a LOT of preprocessing, and a lot of random number generation
- And most of this won't ever be used! Most of the time, when we hash, we don't make it more than a few indices into the permutation.

So many Permutations!

- OK, so kR repetitions is a LOT of preprocessing, and a lot of random number generation
- And most of this won't ever be used! Most of the time, when we hash, we don't make it more than a few indices into the permutation.
- Idea: Instead of taking just the first hash item that appears in the permutation, take the first (say) 3. Concatenate them together. Then we just need $k/3$ permutations per hash table to get similar bounds.

So many Permutations!

- OK, so kR repetitions is a LOT of preprocessing, and a lot of random number generation
- And most of this won't ever be used! Most of the time, when we hash, we don't make it more than a few indices into the permutation.
- Idea: Instead of taking just the first hash item that appears in the permutation, take the first (say) 3. Concatenate them together. Then we just need $k/3$ permutations per hash table to get similar bounds.
- So let's say we have $A = \{\text{black, red, green, blue, orange}\}$, and we're looking at a permutation $P = \{\text{purple, red, white, orange, yellow, blue, green, black}\}$.

So many Permutations!

- OK, so kR repetitions is a LOT of preprocessing, and a lot of random number generation
- And most of this won't ever be used! Most of the time, when we hash, we don't make it more than a few indices into the permutation.
- Idea: Instead of taking just the first hash item that appears in the permutation, take the first (say) 3. Concatenate them together. Then we just need $k/3$ permutations per hash table to get similar bounds.
- So let's say we have $A = \{\text{black, red, green, blue, orange}\}$, and we're looking at a permutation $P = \{\text{purple, red, white, orange, yellow, blue, green, black}\}$.
- Then A hashes to redorangeblue

Reducing Permutations

- If you take the \hat{k} first items when hashing, rather than just taking the first one, we only need kR/\hat{k} total permutations.

Reducing Permutations

- If you take the \hat{k} first items when hashing, rather than just taking the first one, we only need kR/\hat{k} total permutations.
- Does this affect the analysis?

Reducing Permutations

- If you take the \hat{k} first items when hashing, rather than just taking the first one, we only need kR/\hat{k} total permutations.
- Does this affect the analysis?
 - Yes; the k we're concatenating for each hash table are no longer independent!

Reducing Permutations

- If you take the \hat{k} first items when hashing, rather than just taking the first one, we only need kR/\hat{k} total permutations.
- Does this affect the analysis?
 - Yes; the k we're concatenating for each hash table are no longer independent!
 - But this works fine in practice (and is used all the time)

Problems with Expectation

- We chose parameters so that buckets are small in expectation (i.e. on average)

Problems with Expectation

- We chose parameters so that buckets are small in expectation (i.e. on average)
- But: time to process a bucket is *quadratic*.

Problems with Expectation

- We chose parameters so that buckets are small in expectation (i.e. on average)
- But: time to process a bucket is *quadratic*.
- So getting unlucky is super costly!

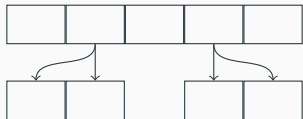
Problems with Expectation

- We chose parameters so that buckets are small in expectation (i.e. on average)
- But: time to process a bucket is *quadratic*.
- So getting unlucky is super costly!
- What can we do if we happen to get a big bucket?

Handling Big Buckets

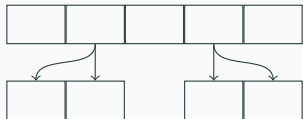
- Recurse!

Handling Big Buckets



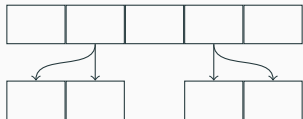
- Recurse!
- Take all items in any really large bucket, rehash them into subbuckets

Handling Big Buckets



- Recurse!
- Take all items in any really large bucket, rehash them into subbuckets
- Might need to repeat

Handling Big Buckets



- Recurse!
- Take all items in any really large bucket, rehash them into subbuckets
- Might need to repeat
- This option can shave off small but significant running time

- Practical datasets have some nice structure

- Practical datasets have some nice structure
 - Similar items tend to be really similar

- Practical datasets have some nice structure
 - Similar items tend to be really similar
 - Tend to appear in clusters

Heuristics in Practice

- Practical datasets have some nice structure
 - Similar items tend to be really similar
 - Tend to appear in clusters
- In practice, can use some simple heuristics in conjunction with MinHash to perform better

Heuristics in Practice

- Practical datasets have some nice structure
 - Similar items tend to be really similar
 - Tend to appear in clusters
- In practice, can use some simple heuristics in conjunction with MinHash to perform better
- If there's enough structure, can use entirely different methods

- Second video on some ideas for implementing MinHash

Conclusion

- Second video on some ideas for implementing MinHash
- You should watch!