# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.
Lecture 10

- Problem Set 2 is due Monday 10/16 at 11:59pm.
- The midterm is in class on Tuesday 10/24. Midterm study material will be posted shortly.
- We have a quiz this week, but not the next two weeks (due to the problem set and midterm).

-I away next week

## Summary

Last Class:

- Discussion of practical algorithms for distinct items estimation (LogLog/HyperLogLog). $10\,1\,0\,0\,0\,0$
- Introduction of Jaccard similarity and the similarity research problem.

## Summary

### Last Class:

- Discussion of practical algorithms for distinct items estimation (LogLog/HyperLogLog).

- Introduction of Jaccard similarity and the similarity research problem.

### This Class:

- Locality sensitive hashing for fast similarity search.

- MinHash as a locality sensitive hash function for Jaccard similarity

- Balancing false positives and negatives with LSH signatures and repeated hash tables.

3

# Search with Jaccard Similarity

$(A, B, C)$
$(C, D, E)$

$\boxed{\dfrac{1}{5}}$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\text{\# shared elements}}{\text{\# total elements}}.$$

**Want Fast Implementations For:**

- **Near Neighbor Search:** Have a database of $n$ sets/bit strings and given a set $A$, want to find if it has high Jaccard similarity to anything in the database. $\Omega(n)$ time with a linear scan.

  **All-pairs Similarity Search:** Have $n$ different sets/bit strings and want to find all pairs with high Jaccard similarity. $\Omega(n^2)$ time if we check all pairs explicitly.

Will speed up via randomized locality sensitive hashing.

approximate    matching

## Locality Sensitive Hashing

Goal: Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).

## Locality Sensitive Hashing

**Goal:** Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).
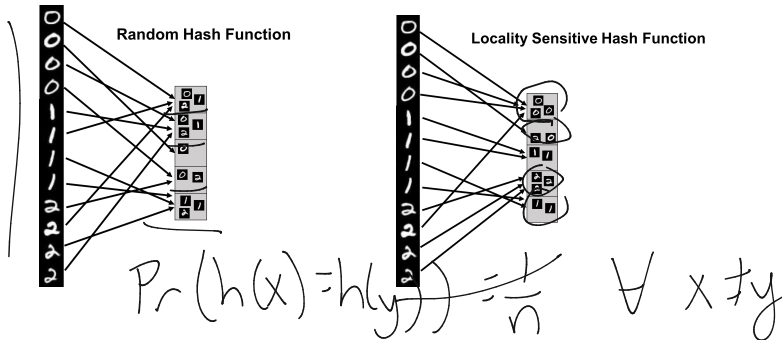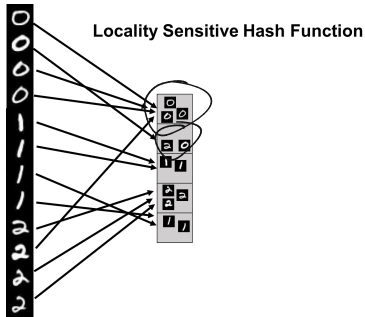
**Strategy:** Locality sensitive hashing (LSH).

- Design a hash function where the collision probability is higher when two inputs are more similar (can design different functions for different similarity metrics.)

# Locality Sensitive Hashing

**Goal:** Speed up Jaccard similarity search (near neighbor and all-pairs similarity search).

**Strategy:** Locality sensitive hashing (LSH).

- Design a hash function where the collision probability is higher when two inputs are more similar (can design different functions for different similarity metrics.)
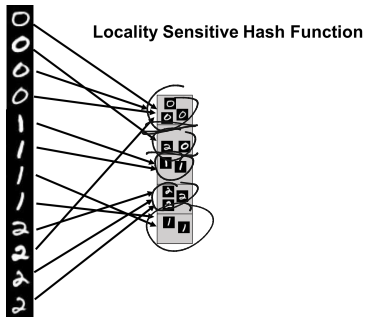


$$Pr(h(x) = h(y)) = \frac{t}{n} \quad \forall \; x \neq y$$

# LSH For Similarity Search

How does locality sensitive hashing (LSH) help with similarity search?



**Locality Sensitive Hash Function**

## LSH For Similarity Search

How does locality sensitive hashing (LSH) help with similarity search?



**Locality Sensitive Hash Function**

- **Near Neighbor Search:** Given item $x$, compute $h(x)$. Only search for similar items in the $h(x)$ bucket of the hash table.
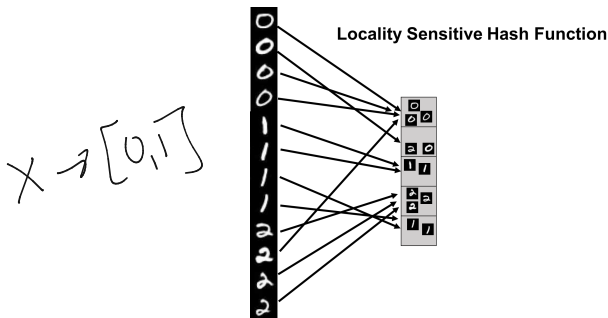
## LSH For Similarity Search

How does locality sensitive hashing (LSH) help with similarity search?



- **Near Neighbor Search:** Given item *x*, compute $h(x)$. Only search for similar items in the $h(x)$ bucket of the hash table.
- **All-pairs Similarity Search:** Scan through all buckets of the hash table and look for similar pairs within each bucket.

# LSH For Similarity Search

How does locality sensitive hashing (LSH) help with similarity search?



**Locality Sensitive Hash Function**

$X \to [0,1]$

- **Near Neighbor Search:** Given item $x$, compute $h(x)$. Only search for similar items in the $h(x)$ bucket of the hash table.
- **All-pairs Similarity Search:** Scan through all buckets of the hash table and look for similar pairs within each bucket.
- We will use $h(x) = g(MinHash(x))$ where $g : [0,1] \to [n]$ is a random hash function. Why?

## MinHashing

An Example: Locality sensitive hashing for Jaccard similarity.

## MinHashing

An Example: Locality sensitive hashing for Jaccard similarity.

Strategy: Use random hashing to map each set to a single hash value. The probably that two sets have colliding hash values will be proportional to their Jaccard similarity.

MinHash(A): [Andrei Broder, 1997 at Altavista]

- Let $h : U \to [0, 1]$ be a random hash function

- $s := 1$

- For $x_1, \ldots, x_{|A|} \in A$

   - $s := \min(s, h(x_k))$

- Return $s$

$$\{A, B, C\} \to .711$$
$$\{C, D, E\} \to .52$$
$$\{A, B, D\} \to .711$$
$$\downarrow \quad \downarrow \quad \downarrow$$
$$.98 \quad .82 \quad (.711)$$

# MinHashing

**An Example:** Locality sensitive hashing for Jaccard similarity.

**Strategy:** Use random hashing to map each set to a single hash value. The probably that two sets have colliding hash values will be proportional to their Jaccard similarity.

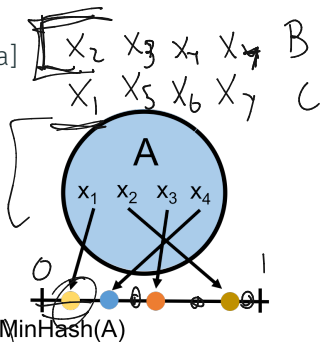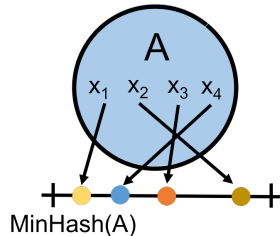**MinHash(A):** [Andrei Broder, 1997 at Altavista]

· Let **h** : $U \rightarrow [0, 1]$ be a random hash function

· **s** := 1

· For $x_1, \ldots, x_{|A|} \in A$

    · **s** := min(**s**, **h**($x_k$))

· Return **s**

$$\begin{array}{cccc} x_2 & x_3 & x_4 & x_7 \\ x_1 & x_5 & x_6 & x_7 \end{array} \quad \begin{array}{c} B \\ C \end{array}$$

$J(A, B) = .6$

$J(A, C) = \frac{1}{7}$



MinHash(A)

$\{ dog, \ cat \}$
$\{ rabbit \ horse \}$

$\{ car, truck \}$

$\mathbb{E}(minHash(A)) = \frac{1}{|A| + 1}$

7

# MinHashing

**An Example:** Locality sensitive hashing for Jaccard similarity.

**Strategy:** Use random hashing to map each set to a single hash value. The probably that two sets have colliding hash values will be proportional to their Jaccard similarity.

**MinHash(A):** [Andrei Broder, 1997 at Altavista]

- Let $h : U \rightarrow [0, 1]$ be a random hash function

- $s := 1$

- For $x_1, \ldots, x_{|A|} \in A$

    - $s := \min(s, h(x_k))$

- Return $s$


MinHash(A)

Identical to our distinct elements sketch!

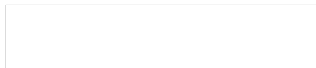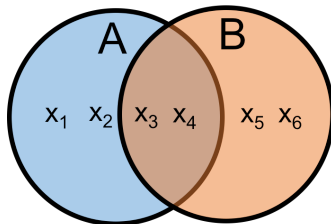## MinHash Analysis

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

$$\Pr\left(\min_{x \in A} h(x) = \min_{y \in B} h(y)\right) = ?$$

## MinHash Analysis

For two sets *A* and *B*, what is $\Pr(MinHash(A) = MinHash(B))$?

$$\Pr\left(\min_{x \in A} \mathsf{h}(x) = \min_{y \in B} \mathsf{h}(y)\right) =?$$

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathsf{h}(x) = \mathsf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)

## MinHash Analysis

For two sets *A* and *B*, what is $\Pr(MinHash(A) = MinHash(B))$?

$$\Pr\left(\min_{x \in A} h(x) = \min_{y \in B} h(y)\right) = ?$$

- Since we are hashing into the continuous range $[0, 1]$, we will never have $h(x) = h(y)$ for $x \neq y$ (i.e., no spurious collisions)
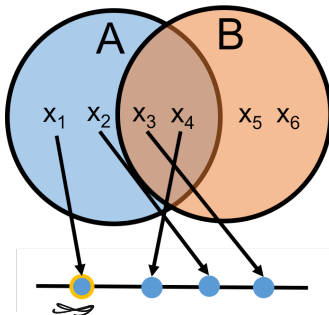
$J(A,B) = \dfrac{2}{6} = \dfrac{1}{3}$

## MinHash Analysis

For two sets *A* and *B*, what is $\Pr(MinHash(A) = MinHash(B))$?

$$\Pr\left(\min_{x \in A} h(x) = \min_{y \in B} h(y)\right) = ?$$

- Since we are hashing into the continuous range $[0, 1]$, we will never have $h(x) = h(y)$ for $x \neq y$ (i.e., no spurious collisions)
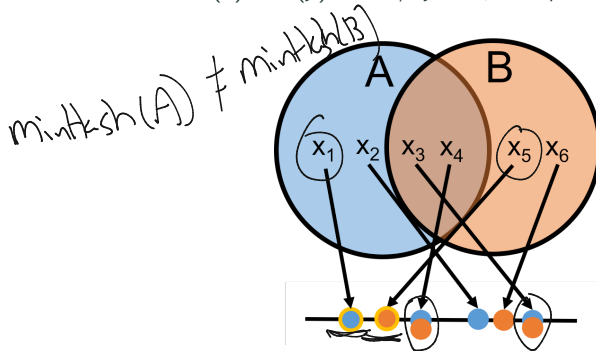
# MinHash Analysis

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

$$\Pr\left(\min_{x \in A} \mathsf{h}(x) = \min_{y \in B} \mathsf{h}(y)\right) = ?$$

- Since we are hashing into the continuous range $[0, 1]$, we will never have $\mathsf{h}(x) = \mathsf{h}(y)$ for $x \neq y$ (i.e., no spurious collisions)

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

$$\Pr\left(\min_{x \in A} h(x) = \min_{y \in B} h(y)\right) = ?$$

- Since we are hashing into the continuous range $[0, 1]$, we will never have $h(x) = h(y)$ for $x \neq y$ (i.e., no spurious collisions)
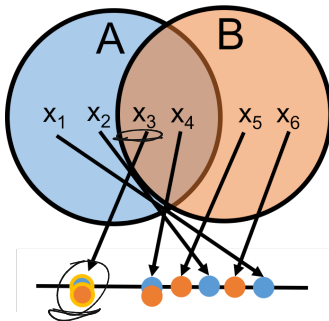
For two sets $A$ and $B$, what is $Pr(MinHash(A) = MinHash(B))$?

$$Pr\left(\min_{x \in A} h(x) = \min_{y \in B} h(y)\right) =?$$

- Since we are hashing into the continuous range $[0, 1]$, we will never have $h(x) = h(y)$ for $x \neq y$ (i.e., no spurious collisions)
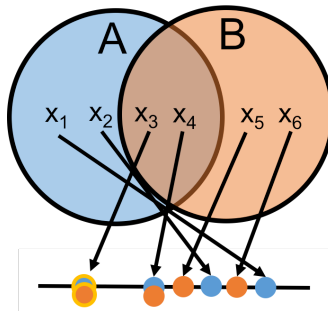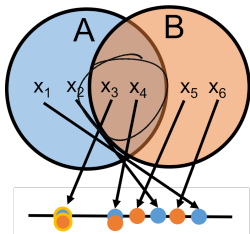
## MinHash Analysis

For two sets *A* and *B*, what is $\Pr(MinHash(A) = MinHash(B))$? $\sim J(A,B)$

**Claim:** *MinHash(A) = MinHash(B)* only if an item in $A \cap B$ has the minimum hash value in both sets.



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

## MinHash Analysis

For two sets *A* and *B*, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$\Pr(MinHash(A) = MinHash(B)) = ?$

## MinHash Analysis

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$$\Pr(MinHash(A) = MinHash(B)) = \frac{|A \cap B|}{\text{total \# items hashed}}$$

## MinHash Analysis

For two sets *A* and *B*, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$$\Pr(MinHash(A) = MinHash(B)) = \frac{|A \cap B|}{\text{total \# items hashed}}$$
$$= \frac{|A \cap B|}{|A \cup B|}$$

## MinHash Analysis

For two sets *A* and *B*, what is $\Pr(MinHash(A) = MinHash(B))$?

**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



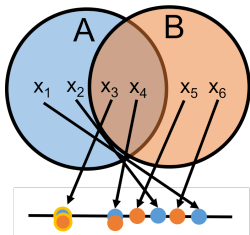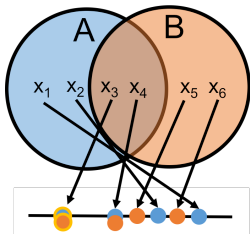$$\Pr(MinHash(A) = MinHash(B)) = \frac{|A \cap B|}{\text{total \# items hashed}}$$

$$= \frac{|A \cap B|}{|A \cup B|} = J(A, B).$$

## MinHash Analysis

For two sets $A$ and $B$, what is $\Pr(MinHash(A) = MinHash(B))$?

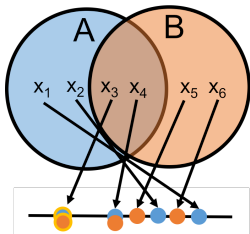**Claim:** $MinHash(A) = MinHash(B)$ only if an item in $A \cap B$ has the minimum hash value in both sets.



$$\Pr(MinHash(A) = MinHash(B)) = \frac{|A \cap B|}{\text{total \# items hashed}}$$
$$= \frac{|A \cap B|}{|A \cup B|} = J(A, B).$$

Locality sensitive: the higher $J(A, B)$ is, the more likely $MinHash(A), MinHash(B)$ are to collide.

## Similarity Search with MinHash

**Goal:** Given a document $y$, identify all documents $x$ in a database with Jaccard similarity (of their shingle sets) $J(x, y) \geq 1/2$.

## Similarity Search with MinHash

**Goal:** Given a document $y$, identify all documents $x$ in a database with Jaccard similarity (of their shingle sets) $J(x, y) \geq 1/2$.

### Our Approach:

- Create a hash table of size $m$, choose a random hash function $\mathbf{g} : [0, 1] \rightarrow [m]$, and insert every item $x$ into bucket $\mathbf{g}(MinHash(x))$. Search for items similar to $y$ in bucket $\mathbf{g}(MinHash(y))$.



**Locality Sensitive Hash Function**

**Goal:** Given a document $y$, identify all documents $x$ in a database with Jaccard similarity (of their shingle sets) $J(x, y) \geq 1/2$.

**Our Approach:**

- Create a hash table of size $m$, choose a random hash function $\mathbf{g} : [0, 1] \to [m]$, and insert every item $x$ into bucket $\mathbf{g}(MinHash(x))$. Search for items similar to $y$ in bucket $\mathbf{g}(MinHash(y))$.

- What is $\Pr\left[\mathbf{g}(MinHash(x)) = \mathbf{g}(MinHash(y))\right]$ assuming $J(x, y) = 1/2$ and $\mathbf{g}$ is collision free?

$$\frac{1}{m}$$

$$\gtrsim \frac{1}{m}$$

$$\frac{1}{2} + \underset{\sim \frac{1}{m}}{\partial}$$

**Goal:** Given a document $y$, identify all documents $x$ in a database with Jaccard similarity (of their shingle sets) $J(x, y) \geq 1/2$.
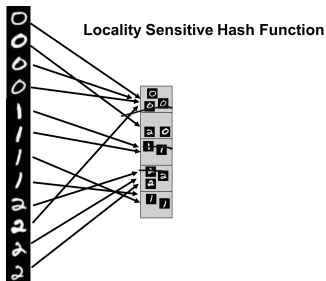
**Our Approach:**

- Create a hash table of size $m$, choose a random hash function $\mathbf{g} : [0, 1] \to [m]$, and insert every item $x$ into bucket $\mathbf{g}(MinHash(x))$. Search for items similar to $y$ in bucket $\mathbf{g}(MinHash(y))$.

- What is $\Pr\left[\mathbf{g}(MinHash(x)) = \mathbf{g}(MinHash(y))\right]$ assuming $J(x, y) = 1/2$ and $\mathbf{g}$ is collision free? $1/2$

- For every document $x$ in your database with $J(x, y) \geq 1/2$ what is the probability you will find $x$ in bucket $\mathbf{g}(MinHash(y))$? $\geq 1/2$

10

## Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

## Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability $1/2$. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

# Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items <u>similar to $y$</u>, look at all items in bucket **g**$(MH_1(y))$ of the $1^{st}$ table, bucket **g**$(MH_2(y))$ of the $2^{nd}$ table, etc.

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.
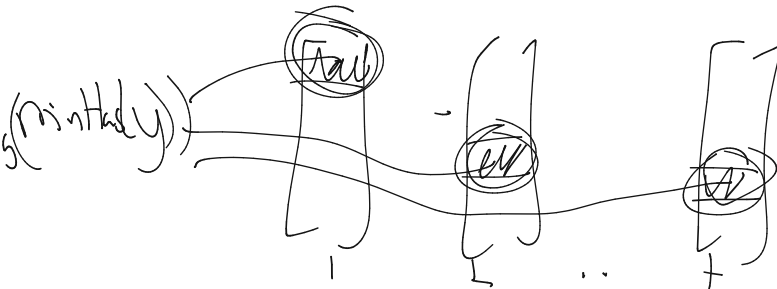
- To search for items similar to $y$, look at all items in bucket **g**$(MH_1(y))$ of the $1^{st}$ table, bucket **g**$(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions?

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items similar to $y$, look at all items in bucket $\mathbf{g}(MH_1(y))$ of the $1^{st}$ table, bucket $\mathbf{g}(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions? $1-$ (probability in *no* buckets)

## Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items similar to $y$, look at all items in bucket $\mathbf{g}(MH_1(y))$ of the $1^{st}$ table, bucket $\mathbf{g}(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions?
  $1-$ (probability in *no* buckets) $= 1 - \left(\frac{1}{2}\right)^t$

## Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items similar to $y$, look at all items in bucket $\mathbf{g}(MH_1(y))$ of the $1^{st}$ table, bucket $\mathbf{g}(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions? $1-$ (probability in *no* buckets) $= 1 - \left(\frac{1}{2}\right)^t \approx .99$ for $t = 7$.

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items similar to $y$, look at all items in bucket $\mathbf{g}(MH_1(y))$ of the $1^{st}$ table, bucket $\mathbf{g}(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions? $1-$ (probability in *no* buckets) $= \underline{1 - \left(\frac{1}{2}\right)^t} \approx .99$ for $t = 7$.
- What is the probability that $x$ with $J(x, y) = 1/4$ is in at least one of these buckets, assuming for simplicity **g** has no collisions?

11

# Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability 1/2. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items similar to $y$, look at all items in bucket **g**$(MH_1(y))$ of the $1^{st}$ table, bucket **g**$(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions? $1-$ (probability in *no* buckets) $= 1 - \left(\frac{1}{2}\right)^t \approx .99$ for $t = 7$.
- What is the probability that $x$ with $J(x, y) = 1/4$ is in at least one of these buckets, assuming for simplicity **g** has no collisions? $1-$ (probability in *no* buckets) $= 1 - \left(\frac{3}{4}\right)^t$

# Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability $1/2$. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items similar to $y$, look at all items in bucket $\mathbf{g}(MH_1(y))$ of the $1^{st}$ table, bucket $\mathbf{g}(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions?
  $1-$ (probability in *no* buckets) $= 1 - \left(\frac{1}{2}\right)^t \approx .99$ for $t = 7$.
- What is the probability that $x$ with $J(x, y) = 1/4$ is in at least one of these buckets, assuming for simplicity **g** has no collisions?
  $1-$ (probability in *no* buckets) $= 1 - \left(\frac{3}{4}\right)^t \approx .87$ for $t = 7$.

# Reducing False Negatives

With a simple use of MinHash, we miss a match $x$ with $J(x, y) = 1/2$ with probability $1/2$. How can we reduce this false negative rate?

**Repetition:** Run MinHash $t$ times independently, to produce hash values $MH_1(x), \ldots, MH_t(x)$. Apply random hash function **g** to map all these values to locations in $t$ hash tables.

- To search for items similar to $y$, look at all items in bucket $\mathbf{g}(MH_1(y))$ of the $1^{st}$ table, bucket $\mathbf{g}(MH_2(y))$ of the $2^{nd}$ table, etc.
- What is the probability that $x$ with $J(x, y) = 1/2$ is in at least one of these buckets, assuming for simplicity **g** has no collisions?
  $1-$ (probability in *no* buckets) $= 1 - \left(\frac{1}{2}\right)^t \approx .99$ for $t = 7$.
- What is the probability that $x$ with $J(x, y) = 1/4$ is in at least one of these buckets, assuming for simplicity **g** has no collisions?
  $1-$ (probability in *no* buckets) $= 1 - \left(\frac{3}{4}\right)^t \approx .87$ for $t = 7$.
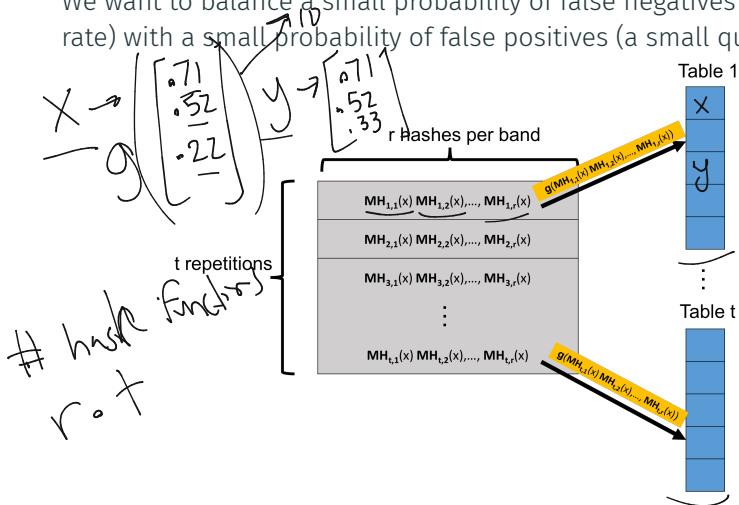
Potential for a lot of false positives! Slows down search time.

## Balancing Hit Rate and Query Time

We want to balance a small probability of false negatives (a high hit rate) with a small probability of false positives (a small query time.)

We want to balance a small probability of false negatives (a high hit rate) with a small probability of false positives (a small query time.)



Create $t$ hash tables. Each is indexed into not with a single MinHash value, but with $r$ values, appended together. A length $r$ signature.

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x, y) = s$:

## Balancing Hit Rate and Query Time

Consider searching for matches in *t* hash tables, using MinHash $1/t$ signatures of length *r*. For *x* and *y* with Jaccard similarity $J(x, y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x, y) = s$.

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x, y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x, y) = s$.

$$x \rightarrow \begin{bmatrix} .71 \\ .63 \\ .22 \end{bmatrix} \quad y \rightarrow \begin{bmatrix} .71 \\ .63 \\ .22 \end{bmatrix}$$

- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right]$

$$s^r$$

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x, y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x, y) = s$.

- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right] = s^r$.

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x, y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x, y) = s$.

- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right] = s^r$.

- Probability that $x$ and $y$ don't match in repetition $i$:

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x, y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x, y) = s$.

- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right] = s^r$.

- Probability that $x$ and $y$ don't match in repetition $i$: $1 - s^r$.

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x, y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x, y) = s$.

- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right] = s^r$.

- Probability that $x$ and $y$ don't match in repetition $i$: $1 - s^r$.

- Probability that $x$ and $y$ don't match in *all repetitions*:

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x, y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x, y) = s$.
- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right] = s^r$.
- Probability that $x$ and $y$ don't match in repetition $i$: $1 - s^r$.
- Probability that $x$ and $y$ don't match in *all repetitions*: $(1 - s^r)^t$.

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x,y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x,y) = s$.

- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right] = s^r$.

- Probability that $x$ and $y$ don't match in repetition $i$: $1 - s^r$.

- Probability that $x$ and $y$ don't match in *all repetitions*: $(1 - s^r)^t$.

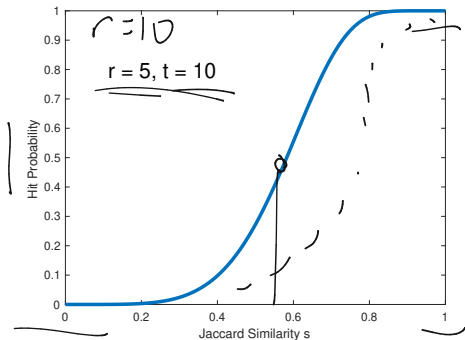- Probability that $x$ and $y$ match in at least one repetition:

## Balancing Hit Rate and Query Time

Consider searching for matches in $t$ hash tables, using MinHash signatures of length $r$. For $x$ and $y$ with Jaccard similarity $J(x,y) = s$:

- Probability that a single hash matches.
  $\Pr\left[MH_{i,j}(x) = MH_{i,j}(y)\right] = J(x,y) = s$.

- Probability that $x$ and $y$ having matching signatures in repetition $i$. $\Pr\left[MH_{i,1}(x), \ldots, MH_{i,r}(x) = MH_{i,1}(y), \ldots, MH_{i,r}(y)\right] = s^r$.

- Probability that $x$ and $y$ don't match in repetition $i$: $1 - s^r$.

- Probability that $x$ and $y$ don't match in *all repetitions*: $(1 - s^r)^t$.

- Probability that $x$ and $y$ match in at least one repetition:

$$\text{Hit Probability: } 1 - (1 - s^r)^t.$$

## The $s$-curve

Using $t$ repetitions each with a signature of $r$ MinHash values, the probability that $x$ and $y$ with Jaccard similarity $J(x, y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.
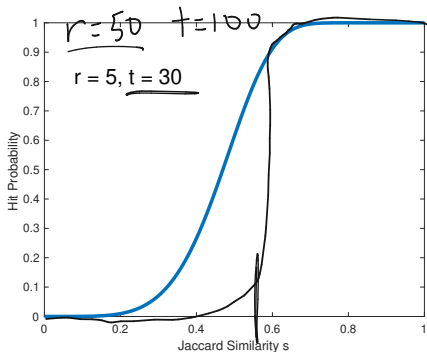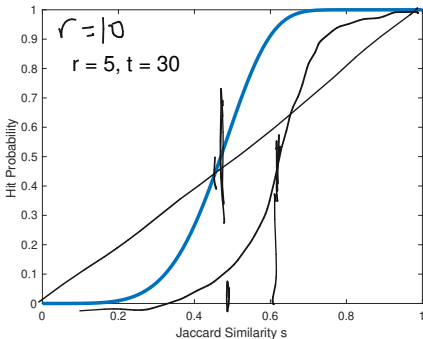
Using *t* repetitions each with a signature of *r* MinHash values, the
probability that *x* and *y* with Jaccard similarity $J(x, y) = s$ match in at
least one repetition is: $1 - (1 - s^r)^t$.

## The $s$-curve

Using $t$ repetitions each with a signature of $r$ MinHash values, the probability that $x$ and $y$ with Jaccard similarity $J(x, y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.

# The s-curve

Using $t$ repetitions each with a signature of $r$ MinHash values, the probability that $x$ and $y$ with Jaccard similarity $J(x, y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.

# The *s*-curve

Using *t* repetitions each with a signature of *r* MinHash values, the probability that *x* and *y* with Jaccard similarity $J(x, y) = s$ match in at least one repetition is: $1 - (1 - s^r)^t$.



$r = 10$
$r = 5, t = 30$

Hit Probability

Jaccard Similarity s

*r* and *t* are tuned depending on application. 'Threshold' when hit probability is 1/2 is $\approx (1/t)^{1/r}$. E.g., $\approx (1/30)^{1/5} = .51$ in this case.

14

## s-curve Example

For example: Consider a database with 10, 000, 000 audio clips. You are given a clip $x$ and want to find any $y$ in the database with $J(x, y) \geq .9$.

## s-curve Example

For example: Consider a database with 10, 000, 000 audio clips. You are given a clip *x* and want to find any *y* in the database with $J(x, y) \geq .9$.

- There are 10 true matches in the database with $J(x, y) \geq .9$.
- There are 10, 000 near matches with $J(x, y) \in [.7, .9]$.

## *s*-curve Example

For example: Consider a database with 10, 000, 000 audio clips. You are given a clip *x* and want to find any *y* in the database with $J(x, y) \geq .9$.

- There are 10 true matches in the database with $J(x, y) \geq .9$.
- There are 10, 000 near matches with $J(x, y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x, y) = s$ is $1 - (1 - s^{25})^{50}$.

## s-curve Example

For example: Consider a database with $10,000,000$ audio clips. You are given a clip $x$ and want to find any $y$ in the database with $J(x, y) \geq .9$.

- There are 10 true matches in the database with $J(x, y) \geq .9$.
- There are $10,000$ near matches with $J(x, y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x, y) = s$ is $1 - (1 - s^{25})^{50}$.

- Hit probability for $J(x, y) \geq .9$ is $\geq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \leq .7$ is $\leq 1 - (1 - .7^{25})^{50} \approx .007$

15

## *s*-curve Example

For example: Consider a database with 10, 000, 000 audio clips. You are given a clip *x* and want to find any *y* in the database with $J(x, y) \geq .9$.

$O(n)$ ~.66   $O(n^c)$ ~.5

- There are 10 true matches in the database with $J(x, y) \geq .9$
- There are 10, 000 near matches with $J(x, y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x, y) = s$ is $1 - (1 - s^{25})^{50}$.

- Hit probability for $J(x, y) \geq .9$ is $\geq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \leq .7$ is $\leq 1 - (1 - .7^{25})^{50} \approx .007$

Expected Number of Items Scanned: (proportional to query time)

$$\leq 10 + .98 * 10, 000 + .007 * 9, 989, 990 \approx 80, 000$$

15

## *s*-curve Example

For example: Consider a database with 10,000,000 audio clips. You are given a clip *x* and want to find any *y* in the database with $J(x, y) \geq .9$.

- There are 10 true matches in the database with $J(x, y) \geq .9$.
- There are 10,000 near matches with $J(x, y) \in [.7, .9]$.

With signature length $r = 25$ and repetitions $t = 50$, hit probability for $J(x, y) = s$ is $1 - (1 - s^{25})^{50}$.

- Hit probability for $J(x, y) \geq .9$ is $\geq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{25})^{50} \approx .98$
- Hit probability for $J(x, y) \leq .7$ is $\leq 1 - (1 - .7^{25})^{50} \approx .007$

Expected Number of Items Scanned: (proportional to query time)

$\leq 10 + .98 * 10,000 + .007 * 9,989,990 \approx 80,000 \ll 10,000,000.$

# Hashing for Duplicate Detection

|  | Hash Table | Bloom Filters | MinHash Similarity Search | Distinct Elements |
|---|---|---|---|---|
| **Goal** | Check if x is a duplicate of any y in database and return y. | Check if x is a duplicate of y in database. | Check if x is a duplicate of any y in database and return y. | Count # of items, excluding duplicates. |
| **Space** | $O(n)$ items | $O(n)$ bits | $O(n \cdot t)$ items (when t tables used) | $O\left(\frac{\log \log n}{\epsilon^2}\right)$ |
| **Query Time** | $O(1)$ | $O(1)$ | Potentially o($n$) | NA |
| **Approximate Duplicates?** | ✘ | ✘ | ✔ | ✘ |

All different variants of detecting duplicates/finding matches in large datasets. An important problem in many contexts!

# Generalizing Locality Sensitive Hashing

Repetition and $s$-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.
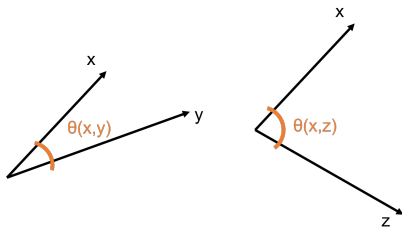
## Generalizing Locality Sensitive Hashing

Repetition and s-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.

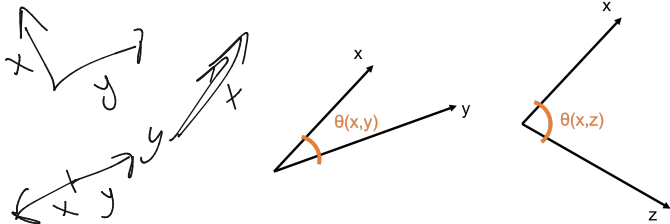- LSH schemes exist for many similarity/distance measures: hamming distance, cosine similarity, etc.

## Generalizing Locality Sensitive Hashing

Repetition and *s*-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.

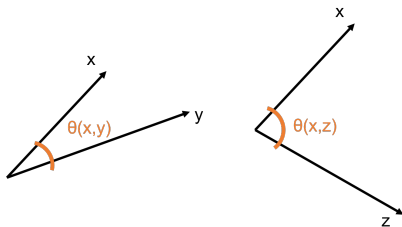- LSH schemes exist for many similarity/distance measures: hamming distance, cosine similarity, etc.

## Generalizing Locality Sensitive Hashing

Repetition and *s*-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.

- LSH schemes exist for many similarity/distance measures: hamming distance, cosine similarity, etc.

## Generalizing Locality Sensitive Hashing

Repetition and *s*-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.

- LSH schemes exist for many similarity/distance measures: hamming distance, cosine similarity, etc.



Cosine Similarity: $\cos(\theta(x, y))$

# Generalizing Locality Sensitive Hashing

Repetition and *s*-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.

- LSH schemes exist for many similarity/distance measures: hamming distance, cosine similarity, etc.
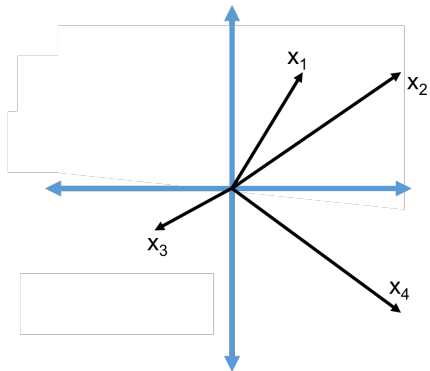


Cosine Similarity: $\cos(\theta(x, y))$

- $\cos(\theta(x, y)) = 1$ when $\theta(x, y) = 0°$ and $\cos(\theta(x, y)) = 0$ when $\theta(x, y) = 90°$, and $\cos(\theta(x, y)) = -1$ when $\theta(x, y) = 180°$

# Generalizing Locality Sensitive Hashing

Repetition and *s*-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.

- LSH schemes exist for many similarity/distance measures: hamming distance, cosine similarity, etc.



**Cosine Similarity:** $\cos(\theta(x, y)) = \frac{\langle x, y \rangle}{\|x\|_2 \cdot \|y\|_2}$.

- $\cos(\theta(x, y)) = 1$ when $\theta(x, y) = 0°$ and $\cos(\theta(x, y)) = 0$ when $\theta(x, y) = 90°$, and $\cos(\theta(x, y)) = -1$ when $\theta(x, y) = 180°$

## SimHash for Cosine Similarity

**SimHash Algorithm:** LSH for cosine similarity.

# SimHash for Cosine Similarity

**SimHash Algorithm:** LSH for cosine similarity.

# SimHash for Cosine Similarity

SimHash Algorithm: LSH for cosine similarity.



random plane

# SimHash for Cosine Similarity

**SimHash Algorithm:** LSH for cosine similarity.

# SimHash for Cosine Similarity

**SimHash Algorithm:** LSH for cosine similarity.



$SimHash(x) = \mathrm{sign}(\langle x, t \rangle)$ for a random vector $t$.

## SimHash for Cosine Similarity

What is $\Pr[SimHash(x) = SimHash(y)]$?

What is $\Pr[SimHash(x) = SimHash(y)]$?

$SimHash(x) \neq SimHash(y)$ when the plane separates $x$ from $y$.

What is $\Pr[SimHash(x) = SimHash(y)]$?

$SimHash(x) \neq SimHash(y)$ when the plane separates $x$ from $y$.

What is $\Pr[SimHash(x) = SimHash(y)]$?

$SimHash(x) \neq SimHash(y)$ when the plane separates $x$ from $y$.



- $\Pr[SimHash(x) \neq SimHash(y)] = \frac{\theta(x,y)}{\pi}$

What is $\Pr[SimHash(x) = SimHash(y)]$?

$SimHash(x) \neq SimHash(y)$ when the plane separates $x$ from $y$.



$J(x,y)$

- $\Pr[SimHash(x) \neq SimHash(y)] = \frac{\theta(x,y)}{\pi}$
- $\Pr[SimHash(x) = SimHash(y)] = 1 - \frac{\theta(x,y)}{\pi} \approx \frac{\cos(\theta(x,y))+1}{2}$

Questions on MinHash and Locality Sensitive Hashing?