

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Andrew McGregor

Lecture 9

- MinHash reduces estimating the Jaccard similarity of two sets to checking equality of a *single number*.

$$\Pr(\text{MinHash}(A) = \text{MinHash}(B)) = J(A, B) = |A \cap B| / |A \cup B|$$

LOCALITY SENSITIVE HASHING

- MinHash reduces estimating the Jaccard similarity of two sets to checking equality of a *single number*.

$$\Pr(\text{MinHash}(A) = \text{MinHash}(B)) = J(A, B) = |A \cap B| / |A \cup B|$$

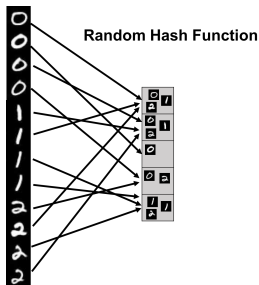
- An instance of **locality sensitive hashing** (LSH), i.e., where the collision probability is higher when two inputs are more similar.

LOCALITY SENSITIVE HASHING

- MinHash reduces estimating the Jaccard similarity of two sets to checking equality of a *single number*.

$$\Pr(\text{MinHash}(A) = \text{MinHash}(B)) = J(A, B) = |A \cap B| / |A \cup B|$$

- An instance of **locality sensitive hashing** (LSH), i.e., where the collision probability is higher when two inputs are more similar.

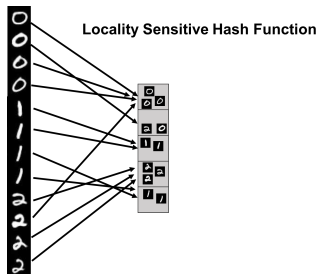
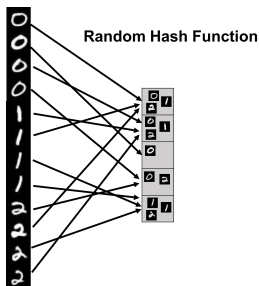


LOCALITY SENSITIVE HASHING

- MinHash reduces estimating the Jaccard similarity of two sets to checking equality of a *single number*.

$$\Pr(\text{MinHash}(A) = \text{MinHash}(B)) = J(A, B) = |A \cap B| / |A \cup B|$$

- An instance of **locality sensitive hashing** (LSH), i.e., where the collision probability is higher when two inputs are more similar.



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

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:

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 x into bucket $\mathbf{g}(\text{MinHash}(x))$.
Search for items similar to y in bucket $\mathbf{g}(\text{MinHash}(y))$.

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 $g : [0, 1] \rightarrow [m]$, and insert x into bucket $g(\text{MinHash}(x))$. Search for items similar to y in bucket $g(\text{MinHash}(y))$.
- What is $\Pr [g(\text{MinHash}(x)) = g(\text{MinHash}(y))]$ assuming $J(x, y) \leq 1/2$ and g is collision free?

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 $g : [0, 1] \rightarrow [m]$, and insert x into bucket $g(\text{MinHash}(x))$. Search for items similar to y in bucket $g(\text{MinHash}(y))$.
- What is $\Pr [g(\text{MinHash}(x)) = g(\text{MinHash}(y))]$ assuming $J(x, y) \leq 1/2$ and g is collision free? At most $1/2$

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 $g : [0, 1] \rightarrow [m]$, and insert x into bucket $g(\text{MinHash}(x))$. Search for items similar to y in bucket $g(\text{MinHash}(y))$.
- What is $\Pr[g(\text{MinHash}(x)) = g(\text{MinHash}(y))]$ assuming $J(x, y) \leq 1/2$ and g is collision free? At most $1/2$
- For each document x in your database with $J(x, y) \geq 1/2$ what is the probability you will find x in bucket $g(\text{MinHash}(y))$?

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 $g : [0, 1] \rightarrow [m]$, and insert x into bucket $g(\text{MinHash}(x))$. Search for items similar to y in bucket $g(\text{MinHash}(y))$.
- What is $\Pr[g(\text{MinHash}(x)) = g(\text{MinHash}(y))]$ assuming $J(x, y) \leq 1/2$ and g is collision free? At most $1/2$
- For each document x in your database with $J(x, y) \geq 1/2$ what is the probability you will find x in bucket $g(\text{MinHash}(y))$? At least $1/2$

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 $MH_1(x), \dots, MH_t(x)$. 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 $MH_1(x), \dots, MH_t(x)$. 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.

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 $MH_1(x), \dots, MH_t(x)$. 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?

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 $MH_1(x), \dots, MH_t(x)$. 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?**
 $1 - (\text{probability in no buckets})$

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 $MH_1(x), \dots, MH_t(x)$. 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 1st table, bucket $g(MH_2(y))$ of the 2nd 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{1}{2}\right)^t$

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 $MH_1(x), \dots, MH_t(x)$. 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{1}{2}\right)^t \approx .99$ 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 $MH_1(x), \dots, MH_t(x)$. 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{1}{2}\right)^t \approx .99$ 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?**

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 $MH_1(x), \dots, MH_t(x)$. 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{1}{2}\right)^t \approx .99$ 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{3}{4}\right)^t$

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 $MH_1(x), \dots, MH_t(x)$. 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{1}{2}\right)^t \approx .99$ 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{3}{4}\right)^t \approx .87$ 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 $MH_1(x), \dots, MH_t(x)$. 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{1}{2}\right)^t \approx .99$ 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?**
 $1 - (\text{probability in no buckets}) = 1 - \left(\frac{3}{4}\right)^t \approx .87$ for $t = 7$.

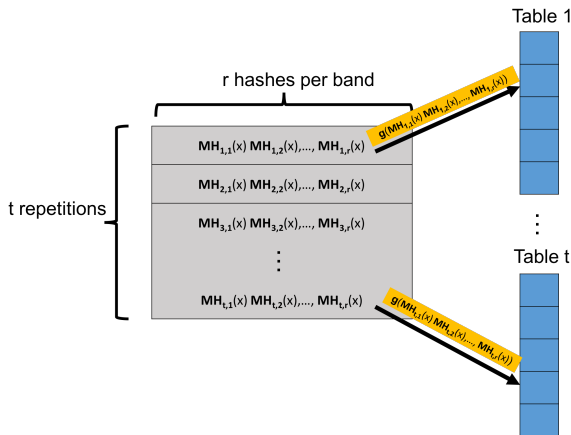
Potential for a lot of false positives! Slows down search time.

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.)

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.)



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 signatures of length r . For x and y with Jaccard similarity $J(x, y) = s$:

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

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

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

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

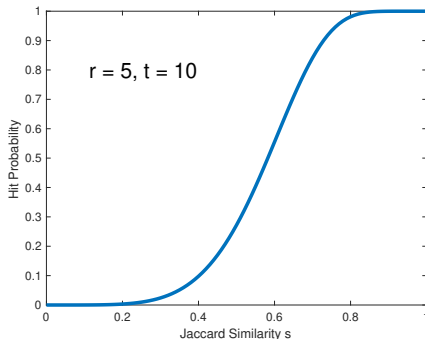
- Probability that a single hash matches.
 $\Pr[MH_{i,j}(x) = MH_{i,j}(y)] = J(x, y) = s$.
- Probability that x and y having matching signatures in repetition i .
 $\Pr[MH_{i,1}(x), \dots, MH_{i,r}(x) = MH_{i,1}(y), \dots, MH_{i,r}(y)] = s^r$.
- Probability that x and y don't match in repetition i :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

Hit Probability: $1 - (1 - s^r)^t$.

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: $1 - (1 - s^r)^t$.

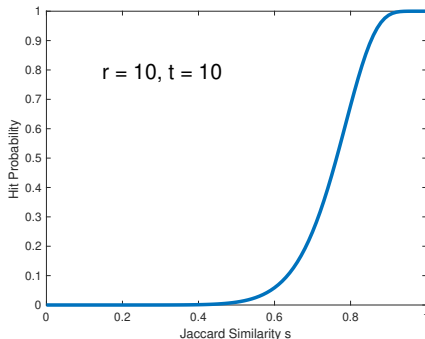
THE s -CURVE

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: $1 - (1 - s^r)^t$.



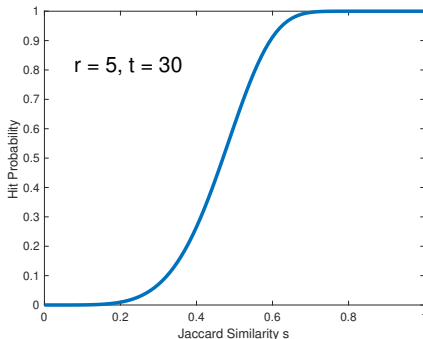
THE s -CURVE

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: $1 - (1 - s^r)^t$.



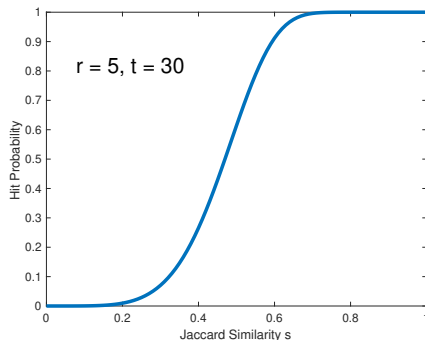
THE s -CURVE

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: $1 - (1 - s^r)^t$.



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 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.

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 $\geq 1 - (1 - .9^{20})^{40} \approx .98$
- Hit probability for $J(x, y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{20})^{40} \approx .98$
- Hit probability for $J(x, y) \leq .7$ is $\leq 1 - (1 - .7^{20})^{40} \approx .007$

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)

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

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 $\geq 1 - (1 - .9^{20})^{40} \approx .98$
- Hit probability for $J(x, y) \in [.7, .9]$ is $\leq 1 - (1 - .9^{20})^{40} \approx .98$
- Hit probability for $J(x, y) \leq .7$ is $\leq 1 - (1 - .7^{20})^{40} \approx .007$

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

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

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

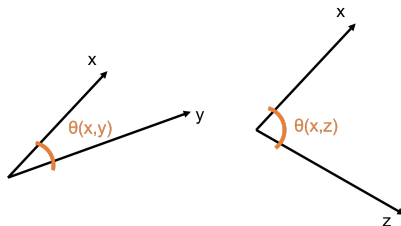
Repetition and *s*-curve tuning can be used for fast similarity search with any similarity metric, given a locality sensitive hash function for that metric.

- 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.

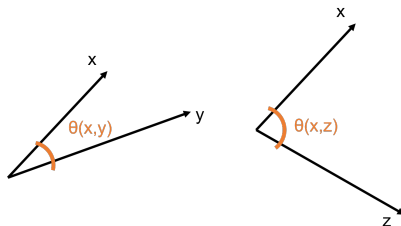
- LSH schemes exist for many similarity/distance measures: hamming distance, **cosine similarity**, etc.



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.

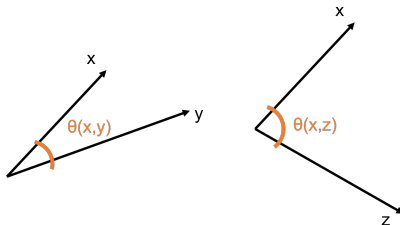


Cosine Similarity: $\cos(\theta(x, y))$

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.



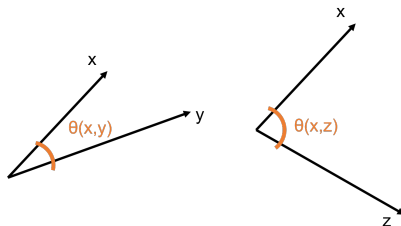
Cosine Similarity: $\cos(\theta(x, y))$

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

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.



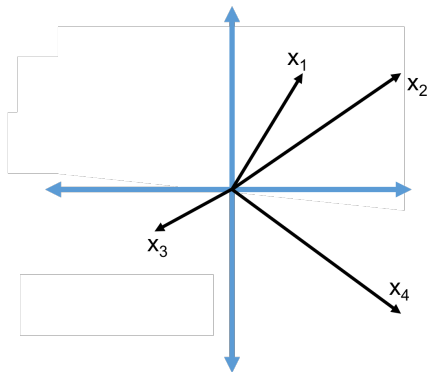
Cosine Similarity: $\cos(\theta(x, y)) = \frac{\langle x, y \rangle}{\|x\|_2 \cdot \|y\|_2}$.

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

SimHash Algorithm: LSH for cosine similarity.

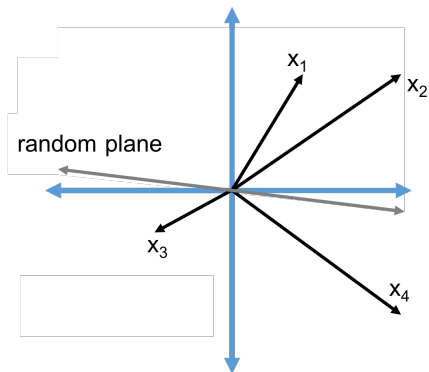
SIMHASH FOR COSINE SIMILARITY

SimHash Algorithm: LSH for cosine similarity.



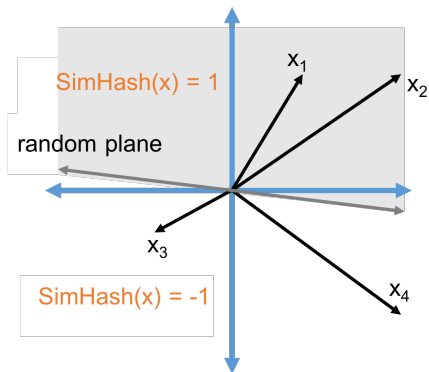
SIMHASH FOR COSINE SIMILARITY

SimHash Algorithm: LSH for cosine similarity.



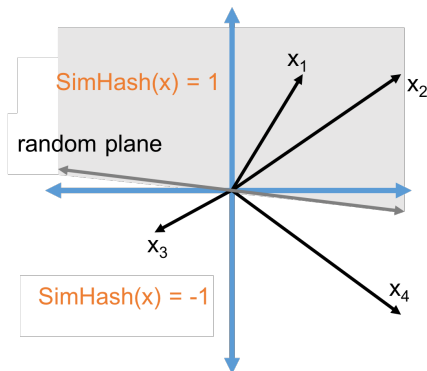
SIMHASH FOR COSINE SIMILARITY

SimHash Algorithm: LSH for cosine similarity.



SIMHASH FOR COSINE SIMILARITY

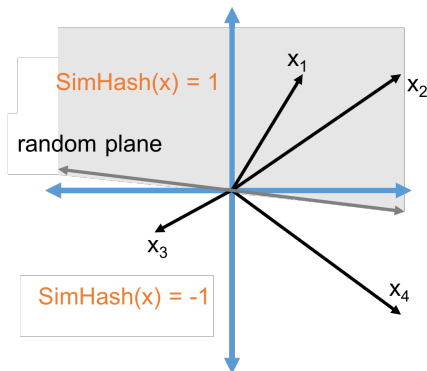
SimHash Algorithm: LSH for cosine similarity.



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

SIMHASH FOR COSINE SIMILARITY

SimHash Algorithm: LSH for cosine similarity.



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

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

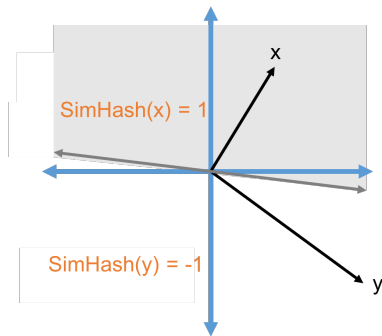
SIMHASH FOR COSINE SIMILARITY

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

SIMHASH FOR COSINE SIMILARITY

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

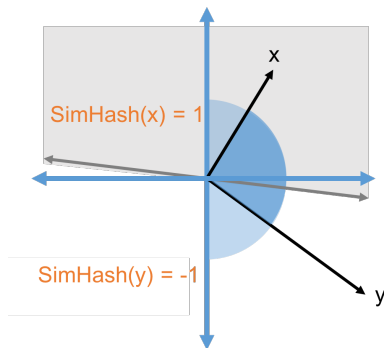
$\text{SimHash}(x) \neq \text{SimHash}(y)$ when the plane separates x from y .



SIMHASH FOR COSINE SIMILARITY

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

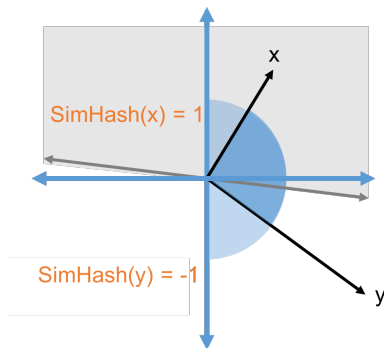
$\text{SimHash}(x) \neq \text{SimHash}(y)$ when the plane separates x from y .



SIMHASH FOR COSINE SIMILARITY

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

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

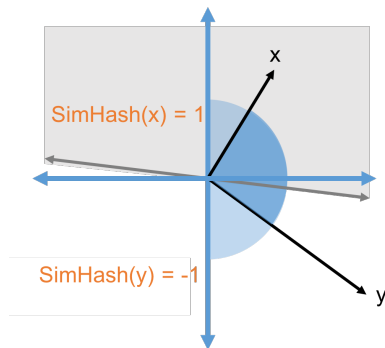


- $\Pr[\text{SimHash}(x) \neq \text{SimHash}(y)] = \frac{\theta(x,y)}{180}$

SIMHASH FOR COSINE SIMILARITY

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

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



- $\Pr[\text{SimHash}(x) \neq \text{SimHash}(y)] = \frac{\theta(x,y)}{180}$
- $\Pr[\text{SimHash}(x) = \text{SimHash}(y)] = 1 - \frac{\theta(x,y)}{180}$

Questions on MinHash and Locality Sensitive Hashing?