

# Randomized and Big Data Algorithms, Approximate Nearest Neighbor Search

**Anton Belyy**

# Roadmap

- Nearest neighbor search
  - Task definition
  - Naïve approaches
  - LSH
- Minwise hashing
  - Motivation
  - Connection to LSH
  - Efficient algorithm

# NN and c-ANN search: recap

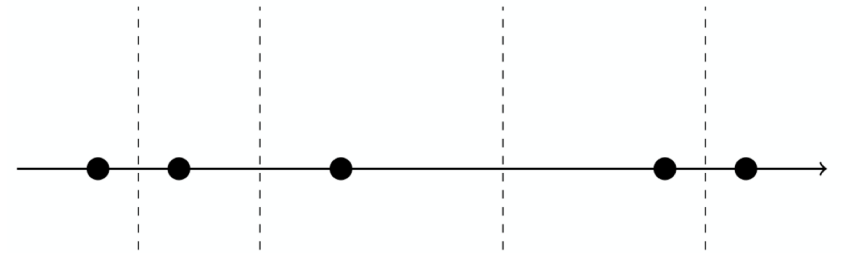
Given a set  $P$  of  $n$  points in  $\mathbb{R}^d$  (assuming  $d \gg 1$ ) and a metric  $D$ , construct a data structure that, given any point  $q$  in  $\mathbb{R}^d$  ...

**Nearest Neighbor Search (NN):** return the **closest** point  $p \in P$ , such that  $p = \operatorname{argmin}_{p \in P} D(p, q)$

**c-Approximate Nearest Neighbor Search (c-ANN):** return a **close** point  $p \in P$ , such that  $D(p, q) \leq c \min_{p'} D(p', q)$

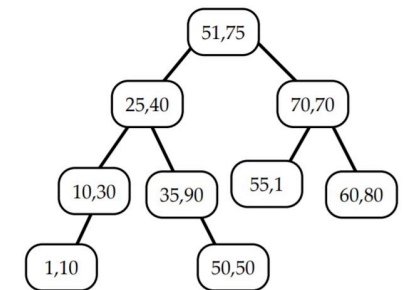
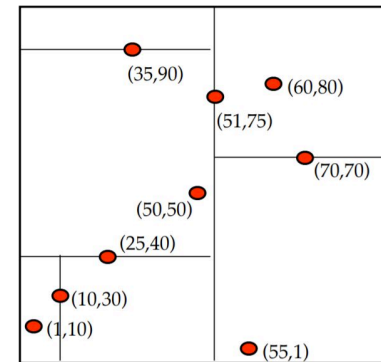
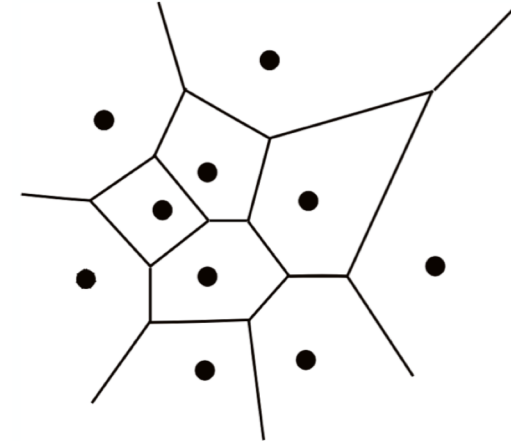
# Naïve Approaches to NN search

- Approach #1: Brute Force
  - preprocessing time:  $O(1)$
  - query time:  $O(nd)$
  - space:  $O(nd)$
- Approach #2: Binary Search Tree ( $d = 1$ )
  - preprocessing time:  $O(n \log n)$
  - query time:  $O(\log n)$
  - space:  $O(n)$



# Less Naïve Approaches to NN search

- Voronoi Diagram ( $d = 2$ )  
preprocessing time:  $O(n \log n)$   
query time:  $O(\log n)$   
space:  $O(n)$
- k-d Tree (small  $d$ )  
preprocessing time:  $O(nd \log n)$   
query time:  **$O(\log n + 2^d)$**   
space:  $O(nd)$  “Curse of dimensionality”



# Solve $c$ -ANN with dimensionality reduction?


**JL Lemma:** Given a set of  $P$  points in  $\mathbb{R}^d$ , there exists a mapping  $f: \mathbb{R}^d \mapsto \mathbb{R}^k$ , such that for any two points  $p_i$  and  $p_j$ , their distance  $\|p_i - p_j\|_2$  is preserved with  $(1 \pm \varepsilon)$  multiplicative error, and  $k = O\left(\frac{1}{\varepsilon^2} \log n\right)$

**Corollary:** k-d trees with JL dimensionality reduction has:

preprocessing time:  $O\left(n \frac{\log^2 n}{\varepsilon^2}\right)$

query time:  $n^{O(1/\varepsilon^2)}$

space:  $O\left(n \frac{\log n}{\varepsilon^2}\right)$



Still super-linear  
dependence on  $n$   
for small  $\varepsilon$

# Sublinear time algorithms for $c$ -ANN

We will show an alternative way of distance-preserving dimensionality reduction from  $d$  to  $k$ , based on **locality sensitive** hash functions.

Using a simple lookup over multiple hash tables, we will solve  $c$ -ANN problem in provably **sublinear** time.

All we need is to construct a locality-sensitive **hash family** for a given metric  $D$  we care about.

# Locality-Sensitive Hashing (LSH) Framework

(Indyk and Motwani, 1998)

**Definition:** A hash family  $\mathcal{H} = \{h: U \mapsto S\}$  is called  $(r_1, r_2, p_1, p_2)$ -**LSH** for a metric function  $D$ , if for all points  $x, y \in U$ ,

1. If  $D(x, y) \leq r_1$ , then  $\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \geq p_1$
2. If  $D(x, y) > r_2$ , then  $\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq p_2$

**Compare:**

2-universal hashes:  $\Pr[h(x) = h(y)] \leq \frac{1}{|S|}$

Locality sensitive hashes:  $\Pr[h(x) = h(y)] \propto \frac{1}{D(x,y)}$



# Locality-Sensitive Hashing (LSH) Framework

(Indyk and Motwani, 1998)

**Exercise (LSH for Hamming distance):**

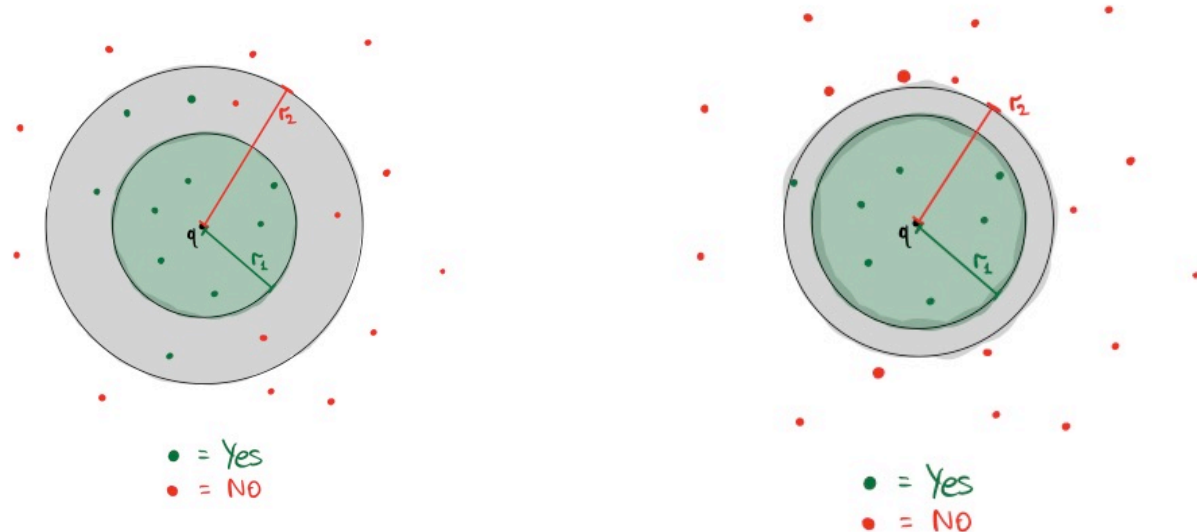
$$U = \{0, 1\}^d, \mathcal{H} = \{h_i: h_i(x) = x_i, i \in [d]\},$$

$$D(x, y) = \sum_{i=1}^d [x_i \neq y_i]. \text{ Show that } \mathcal{H} \text{ is } \left(r, cr, 1 - \frac{r}{d}, 1 - \frac{cr}{d}\right)\text{-LSH.}$$

# $(c,r)$ -ANN: a bridge between LSH and $c$ -ANN

**$(c,r)$ -ANN**: Given a set  $P$  of points in  $\mathbb{R}^d$ , construct a data structure for the set  $P$ , such that given any point  $q$  in  $\mathbb{R}^d$ ,

1. If there exists  $p \in P$  such that  $D(p, q) \leq r$ , return **YES** and *any* point  $p' \in P$  with  $D(p', q) \leq cr$
2. If there is no point  $p \in P$  s.t.  $D(p, q) \leq cr$ , return **NO**



# Using $(c,r)$ -ANN to solve $c$ -ANN

**Lemma 1:** Given an algorithm for  $(1 + \varepsilon, r)$ -ANN (for all  $r$ ), there exists an algorithm for  $(1 + \varepsilon)$ -ANN

**Idea:** construct  $(c,r)$ -ANN instances over increasing radii  $r$ :

$$D_{min}, \quad (1 + \varepsilon)D_{min}, \quad (1 + \varepsilon)^2 D_{min}, \quad \dots, \quad D_{max}$$

Use **binary search** to find the smallest ball that contains points from  $P$ .

# Using LSH to solve $(c,r)$ -ANN

**Theorem 1:** Given an  $(r, cr, p_1, p_2)$ -LSH family for a metric  $D$ , there exists an algorithm for  $(c, r)$ -ANN, which succeeds with probability  $> 0$  and uses:

preprocessing time:  $O(n^{1+\rho} \log n)^*$

query time:  $O(n^\rho \log n)^*$

space:  $O(n^{1+\rho} \log n)$ , where  $\rho = \frac{\ln 1/p_1}{\ln 1/p_2} < 1$

We will prove this constructively, i.e. by showing an algorithm.

**Corollary:** Having LSH for a metric  $D$  allows for a sublinear  $(c,r)$ -ANN search, since  $O(n^\rho \log n) = o(n)$ .

\* measured in the number of hash evaluations

# Using LSH to solve $(c,r)$ -ANN

**Index:** for each data point  $p \in P$ , compute  $L$  hash signatures of size  $K$  each, using LSH (denoted as  $g_i(p)$ ,  $i = 1 \dots L$ ), and store  $p$  in  $L$  hash tables ( $HT_i$ ,  $i = 1 \dots L$ ), where each  $HT_i$  is indexed by  $g_i(p)$ :

```
for  $i \in \{1,2, \dots L\}$  do // Initialize L empty hash tables
     $HT_i = \text{HashTable}()$ 
forall  $p \in P$  do // Insert p's into L hash tables
    for  $i \in \{1,2, \dots L\}$  do
         $HT_i[g_i(p)].\text{add}(p)$ 
```

## Using LSH to solve $(c,r)$ -ANN

**Query:** given a query  $q$ , look for potential c-ANN candidates of  $q$  in each of  $HT_i$  hash tables. Stop after checking exactly  $2L$  candidates.

```
num_checked = 0
for  $i \in \{1, 2, \dots, L\}$  do
    for  $p' \in HT_i[g_i(q)]$ :
        num_checked += 1
        if  $D(p', q) \leq cr$ : return (YES,  $p'$ )
        if num_checked ==  $2L$ : return NO
return NO
```

# Using LSH to solve $(c,r)$ -ANN

## Correctness proof:

We need to show that, given a query  $q$ ,

- 1) there are at most  $2L - 1$  “false positives” ( $D(p', q) > cr$ , but  $\exists j': g_{j'}(p') = g_{j'}(q)$ )
- 2) if there is a “true positive”  $c$ -ANN  $p \in P$ , then  $\exists j: g_j(p) = g_j(q)$

With probability  $> 0$

# Using LSH to solve $(c,r)$ -ANN

**Claim #1:** there are at most  $2L - 1$  “false positives” ( $D(p', q) > cr$ , but  $\exists j' : g_{j'}(p') = g_{j'}(q)$ ) with probability  $> 1/2$

Recall from LSH that, if  $D(p', q) > cr$ , then  $\Pr_{h \in \mathcal{H}}[h(p') = h(q)] \leq p_2$

If we compute  $K$  independent hashes, then the probability is  $\leq p_2^K$

Since we have  $n$  points and  $L$  hash tables, we can upper bound:

$$\mathbb{E}[\text{false positives}] \leq p_2^K \cdot n \cdot L = \frac{1}{n} \cdot n \cdot L = L \quad (\text{choose } K = \log_{1/p_2} n)$$

Then, by Markov's inequality:

$$\Pr[\text{false positives} > 2L] \leq \frac{\mathbb{E}[\text{false positives}]}{2L} \leq \frac{L}{2L} \leq \frac{1}{2}.$$



# Using LSH to solve $(c,r)$ -ANN

**Claim #2:** if there is a “true positive”  $c$ -ANN  $p \in P$ , then  $\exists j: g_j(p) = g_j(q)$

Recall from LSH that if  $D(p, q) \leq r$ , then  $\Pr_{h \in \mathcal{H}}[h(p) = h(q)] \geq p_1$

If we compute  $K$  independent hashes, then the probability is at least

$$p_1^K = p_1^{\log_{1/p_2} n} = n^{-\frac{\ln 1/p_1}{\ln 1/p_2}} = n^{-\rho}$$

Considering all  $L$  hash tables (**and choose  $L = n^\rho$** ), we get:

$$1 - (1 - n^{-\rho})^L = 1 - (1 - n^{-\rho})^{n^\rho} \geq 1 - \frac{1}{e} \geq \frac{1}{2}$$

$$\left(1 - \frac{1}{x}\right)^x < \frac{1}{e}$$

## Using LSH to solve $(c,r)$ -ANN

Since both claims are true w.p. at least 0.5, then by the union bound,

$$\begin{aligned} \Pr[(1 \text{ is true}) \text{ AND } (2 \text{ are true})] &= 1 - \Pr[(1 \text{ is false}) \text{ OR } (2 \text{ is false})] \\ &\geq 1 - \Pr[1 \text{ is false}] - \Pr[2 \text{ is false}] > 1 - \frac{1}{2} - \frac{1}{2} > 0. \end{aligned}$$

\end{proof}

Thus, with a constant probability, we can find  $c$ -ANNs by “mapping” to hash signatures of size  $K*L$ , and checking at most  $2*L$  candidates.

**How small is this compared to  $n$ ?  $n*d$ ?**

# Discussion

- Computing many hash signatures per data point is a bottleneck of LSH:
  - Theorem 1 requires a signature of size  $K * L = O(n^\rho \log_{1/p_2} n)$ :

$n \setminus (p_1, p_2)$	(0.6, 0.4)	(0.7, 0.3)	(0.8, 0.2)	(0.9, 0.1)
$n = 10^5$	7,702	290	35	8
$n = 10^6$	33,365	687	58	11
$n = 10^7$	140,518	1586	94	15

$$\rho = \frac{\ln 1/p_1}{\ln 1/p_2}$$

- In practice,  $K*L$  is around 100..1000

**$\Rightarrow$  Fast computation of signatures is key to practical success in LSH**

# Outline

- Nearest neighbor search
  - Task definition
  - Naïve approaches
  - LSH
- Minwise hashing
  - Motivation
  - Connection to LSH
  - Efficient algorithm

# Fingerprints for Duplicate Web Page Detection

Parallel to LSH framework, a method to filter duplicate web pages in search results was developed by Andrei Broder

**Idea:** compute short **fingerprints** of web pages, such that similar web pages get similar fingerprints, where similarity is measured by **Jaccard Index**

This could be seen as a special case LSH where  $D = 1$  - Jaccard Index



Andrei Broder

# Minwise Hashing

**Jaccard Index** is a similarity function between two sets  $A$  and  $B$ , defined as

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

**Lemma 2 (Minwise Hashing):** Let  $A$  and  $B$  be subsets of  $[d]$  and  $h_\pi(X) = \min_{x \in X} \pi(x)$ , where  $\pi$  is a permutation of  $[d]$ . Then  $\Pr_\pi[h_\pi(A) = h_\pi(B)] = J(A, B)$

**Corollary:**  $\mathcal{H} = \{h_\pi : \pi \text{ is a permutation of } [d]\}$  is  $(r_1, r_2, 1 - r_1, 1 - r_2)$ -LSH for Jaccard distance:  $D(A, B) = 1 - J(A, B)$

# Minwise Hashing

**Corollary:**  $\mathcal{H} = \{h_\pi : \pi \text{ is a permutation of } [d]\}$  is  $(r_1, r_2, 1 - r_1, 1 - r_2)$ -LSH for Jaccard distance:  $D(A, B) = 1 - J(A, B)$

# Practical Issues with Minwise Hashing

- Suppose we want to compute  $m$  hashes  $\Rightarrow$  need  $m$  permutations of  $[d]$
- Storing  $m$  permutations naively requires  $m \log d! = O(md \log d)$  bits
  - Typically  $m \sim 10^2 \dots 10^3$  and  $d = 2^{32} \dots 2^{64}$ , so this is not acceptable
- Later, Broder et al.\* showed this to be sufficient for  $\pi$  (for any  $X \subseteq [d]$  and any  $x \in X$ ):

$$\Pr_{\pi \sim S_d} [\min\{\pi(X)\} = \pi(x)] = \frac{1}{|X|}$$

	$\pi_1$	$\pi_2$	$\pi_3$	$\pi_4$
2	2	3	10	3
3	6	4	4	6
4	15	14	16	15
6	11	10	7	5
7	8	7	5	1

$X$   $d = 2^4$

\* Broder, Charikar, and Mitzenmacher. Min-Wise Independent Permutations. 1998



# Practical Issues with Minwise Hashing

- In practice,  $\pi(x)$  is replaced with a 2-universal hash  $g(x)$  for  $x \in X$
- There might be collisions, i.e.  $g(x_1) = g(x_2)$  when  $x_1 \neq x_2$ 
  - When  $d$  is large and  $n \ll d$ , this can be neglected
- Still, time complexity  $O(nm)$ , where  $m$  is the number of hashes
- In LSH, we need 1000s of hashes  $\Rightarrow$  this is not acceptable

	$g_1$	$g_2$	$g_3$	$g_4$
2	15	14	12	5
3	1	8	10	15
4	6	3	4	9
6	5	8	9	8
7	3	1	15	7

$X$   $d = 2^4$

# Minwise Hashing with Exponential Sampling

- Suppose we replace 2-universal hash  $g(x)$  with a random variable  $\tilde{g}_x \sim \text{Exp}(1)$  ( $x \in X \subset [d]$ )
- **Lemma 3:** Let  $A$  and  $B$  be subsets of  $[d]$  and  $h_{\tilde{g}}(X) = \min_{x \in X} \tilde{g}_x$ , where  $\tilde{g}_x$ 's are i.i.d. r.v.s drawn from  $\text{Exp}(1)$ . Then  $\Pr_{\tilde{g}}[h_{\tilde{g}}(A) = h_{\tilde{g}}(B)] = J(A, B)$

	$\tilde{g}_1$	$\tilde{g}_2$	$\tilde{g}_3$	$\tilde{g}_4$
2	0.80	1.26	0.92	0.78
3	0.55	1.04	0.58	2.22
4	3.31	<b>0.48</b>	1.57	0.75
6	0.84	2.60	<b>0.07</b>	<b>0.09</b>
7	<b>0.02</b>	1.79	1.51	2.04

**X**  $d = 2^4$

# Minwise Hashing with Exponential Sampling

**Lemma 3:** Let  $A$  and  $B$  be subsets of  $[d]$  and  $h_{\tilde{g}}(X) = \min_{x \in X} \tilde{g}_x$ , where  $\tilde{g}_x$ 's are i.i.d. r.v.s drawn from  $Exp(1)$ . Then  $\Pr_{\tilde{g}}[h_{\tilde{g}}(A) = h_{\tilde{g}}(B)] = J(A, B)$

# Minwise Hashing with Poisson Processes

- Why sample  $n$  values and throw away  $(n - 1)$  of them, when we can **sample the minimum** of exponentials **directly**?
- This can be done with **Poisson processes**\*

# Intro to Poisson Processes



- Definition:
  - **Continuous time** process with parameter  $\lambda$
  - Number of arrivals in disjoint time intervals are **independent**
  - Number of arrivals in interval of duration  $\tau$ , or  $N_\tau$ , is **Poisson**( $\lambda\tau$ )-distributed and is the same for any interval  $[t_0; t_0 + \tau]$
  - For very small intervals  $\delta$ : 
$$P(k, \delta) \approx \begin{cases} 1 - \lambda\delta, & \text{if } k = 0 \\ \lambda\delta, & \text{if } k = 1 \\ 0, & \text{if } k > 1 \end{cases}$$
- Useful property:  $X_1$  – time of the 1<sup>st</sup> arrival – is **Exp**( $\lambda$ )-distributed, i.e.  $P(X_1 > t) = e^{-\lambda t}$

# Minwise Hashing with Poisson Processes

- Why sample  $n$  values and throw away  $(n - 1)$  of them, when we can **sample the minimum** of exponentials **directly**?
- This can be done with **Poisson processes**\*
- **Idea**: per each  $x \in X$ , start a Poisson process  $P_x$  that generates hash values for all  $m$  signatures and stops, when cannot improve current min

$$X = \{2,3\}, \quad m = 4$$

time = hash value

$P_2$ : (3, 0.8), (1, 2.0), (2, 2.4), (4, 3.1)

$P_3$ : (4, 0.1), (2, 1.6), (2, 3.5)

3.5 > 3.1, stop early

# Minwise Hashing with Poisson Processes

- Why sample  $n$  values and throw away  $(n - 1)$  of them, when we can **sample the minimum** of exponentials **directly**?
- This can be done with **Poisson processes**\*
- **Idea**: per each  $x \in X$ , start a Poisson process  $P_x$  that generates hash values for all  $m$  signatures and stops, when cannot improve current min

$$X = \{2, 3\}, \quad m = 4$$

time = hash value

$P_2$ : (3, 0.8), (1, 2.0), (2, 2.4), (4, 3.1)

$P_3$ : (4, 0.1), (2, 1.6), (2, 3.5)

3.5 > 3.1, stop early

	$h_1$	$h_2$	$h_3$	$h_4$
{2}	2.0	2.4	0.8	3.1
{2, 3}	2.0	<u>1.6</u>	0.8	<u>0.1</u>

# ProbMinHash algorithm (TKDE 2020)

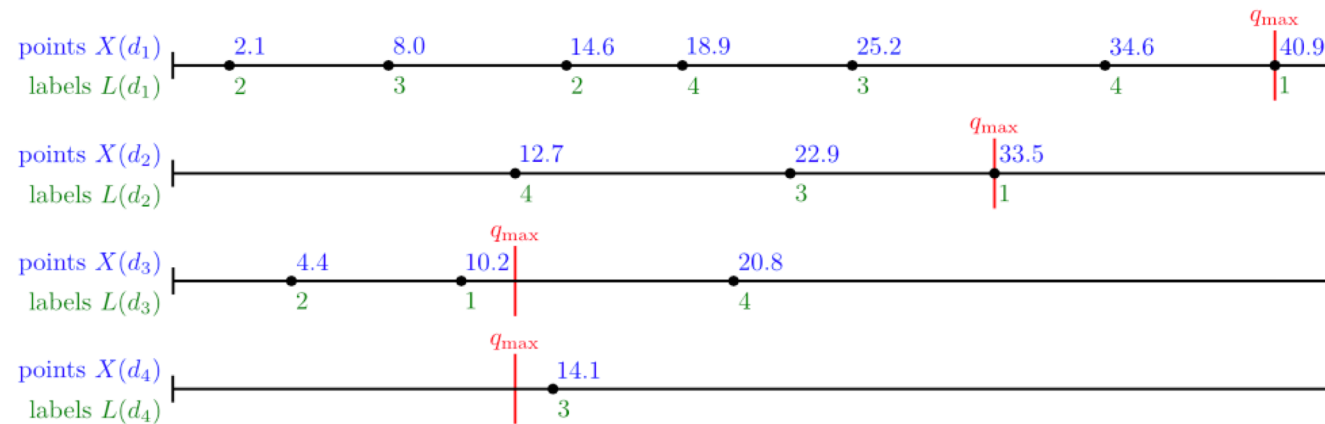
initialization

↓  
add  $d_1$ ,  $w(d_1) = 0.17$

↓  
add  $d_2$ ,  $w(d_2) = 0.09$

↓  
add  $d_3$ ,  $w(d_3) = 0.14$

↓  
add  $d_4$ ,  $w(d_4) = 0.07$



$q =$	1	2	3	4	
	$\infty$	$\infty$	$\infty$	$\infty$	$q_{\max} = \infty$
$z =$	—	—	—	—	

$q =$	40.9	2.1	8.0	18.9	$q_{\max} = 40.9$
$z =$	$d_1$	$d_1$	$d_1$	$d_1$	

$q =$	33.5	2.1	8.0	12.7	$q_{\max} = 33.5$
$z =$	$d_2$	$d_1$	$d_1$	$d_2$	

$q =$	10.2	2.1	8.0	12.7	$q_{\max} = 12.7$
$z =$	$d_3$	$d_1$	$d_1$	$d_2$	

$q =$	10.2	2.1	8.0	12.7	$q_{\max} = 12.7$
$z =$	$d_3$	$d_1$	$d_1$	$d_2$	

- Each arrival is labeled with a hash signature label  $L \sim Uniform(\{1, 2, \dots, m\})$
- By the **split** property, a *sub-process* where all labels are  $L = l_i$  is also Poisson with rate  $\frac{\lambda}{m}$
- If  $\lambda = m$ , then the **arrival time** of the **first event** in each sub-process is **Exp(1)-distributed**



# ProbMinHash algorithm (TKDE 2020)

- Processing of the first element takes  $O(m \log m)$  time (Coupon Collector argument)
- Processing of the  $i$ -th element takes  $O(\frac{1}{i} m \log m)$  time (due to early stop)
- Overall, the algorithm runs in  $O\left(n + \sum_{i=1}^n \frac{1}{i} m \log m\right) = O(n + m \log^2 m)$  time
- Substantial improvement over  $O(nm)$  in Broder's method!

**Input:**  $X$

**Output:**  $h_1, h_2, \dots, h_m$

$h_1, h_2, \dots, h_m \leftarrow (\infty, \infty, \dots, \infty)$

```
forall  $x \in X$  do  
     $R \leftarrow$  new PRNG with seed  $x$   
     $h \leftarrow R[\text{Exp}(1)]$   
    while  $h < h_{\max}$  do  
         $k \leftarrow R[\text{Uniform}(\{1, 2, \dots, m\})]$   
        if  $h < h_k$  then  
             $h_k \leftarrow h$   
             $h_{\max} \leftarrow \max(\{h_1, h_2, \dots, h_m\})$   
            if  $h \geq h_{\max}$  then break  
         $h \leftarrow h + R[\text{Exp}(1)]$ 
```



early stop

# Summary

- NN search is challenging in high-dimensional settings
- Exact algorithms scale at least linearly with either  $n$  or  $d$
- LSH provides sublinear algorithms for  $c$ -ANN search
- Minwise hashing is a conceptually simple way to do  $c$ -ANN search with Jaccard distance as a target metric
- Poisson process-based sampling provides a much faster and statistically equivalent way to perform Minwise Hashing