

Fast Similarity Graph Construction via Data Sketching Techniques

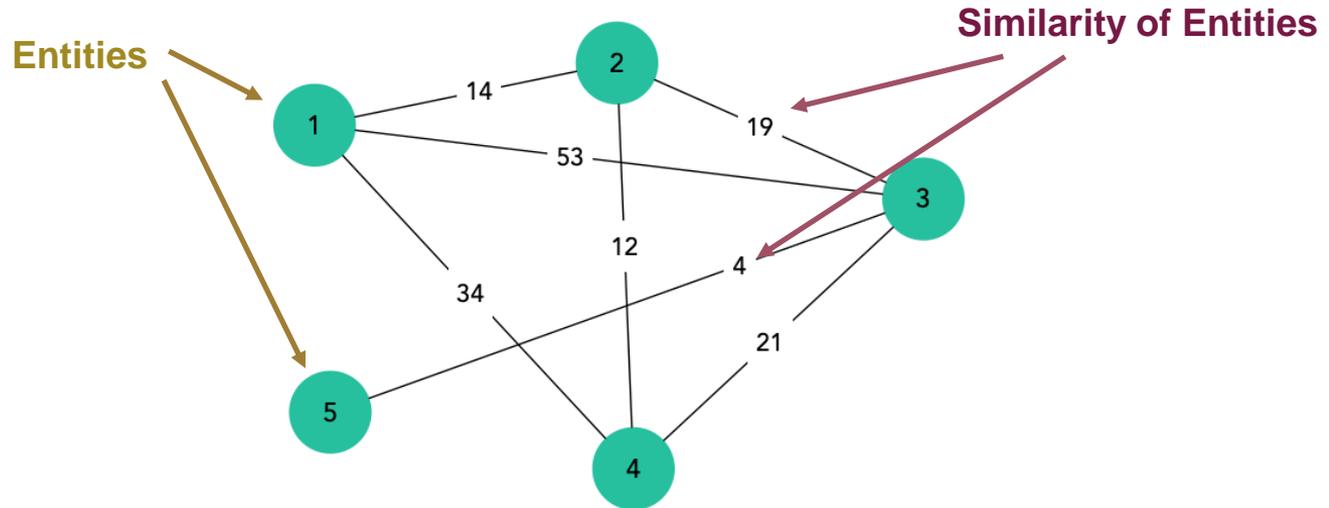
By: Hoorieh Marefat

Supervisors: Aijun An, Manos Papagelis

Introduction

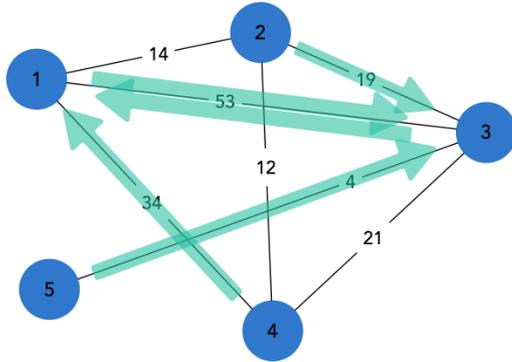
Motivation

We explore how to **build similarity graphs** in an **efficient** way

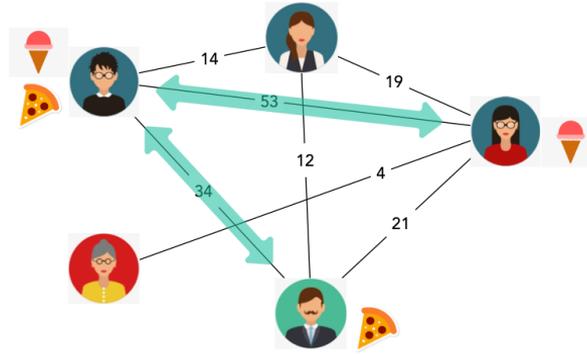


Applications

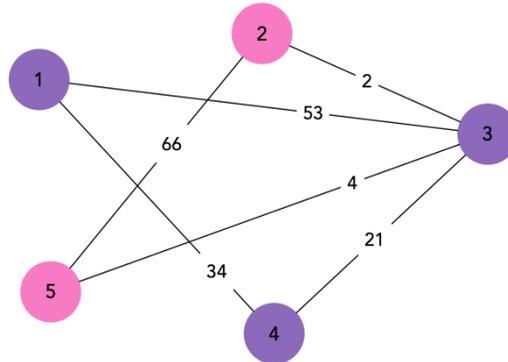
NN Search



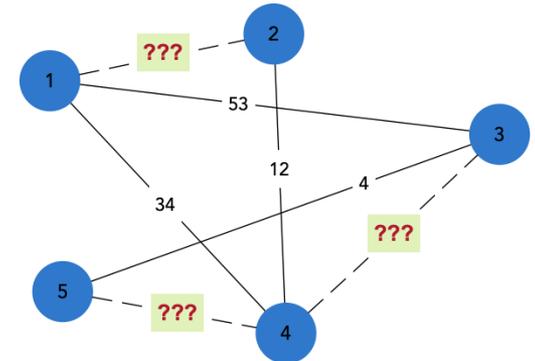
Collaborative Filtering



Clustering



Link Prediction



Similarity-based Graphs

- Similarity Graph
- ϵ -Graph
- Nearest Neighbour (NN) Graph

Similarity graph

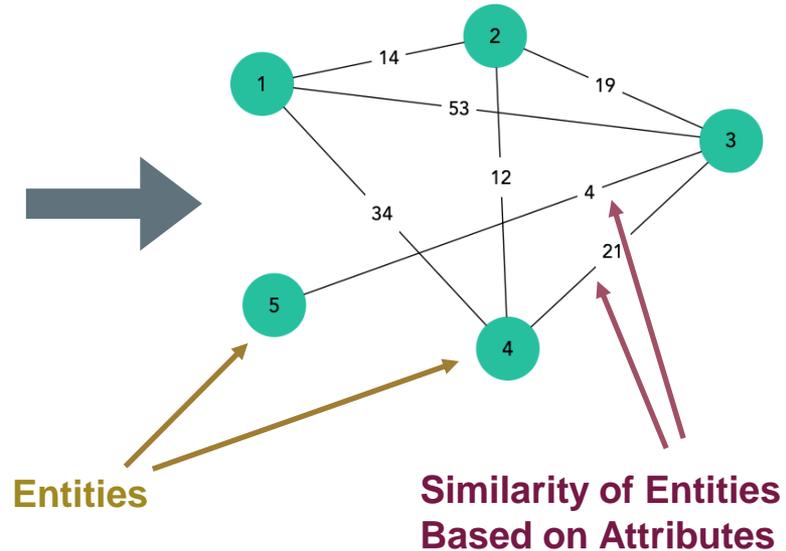
The Dataset

Attributes

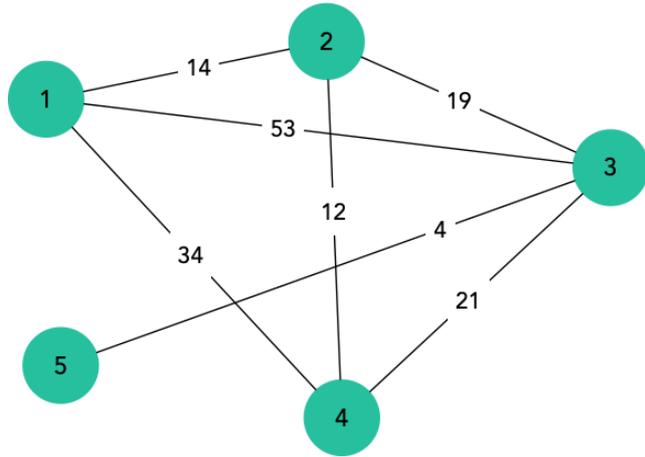
Entities

| UserID | Item1 | Item2 | Item3 | Item4 | Item5 |
|--------|-------|-------|-------|-------|-------|
| 1 | 3 | 0 | 2 | 0 | 10 |
| 2 | 4 | 5 | 1 | 0 | 0 |
| 3 | 1 | 3 | 0 | 1 | 5 |
| 4 | 0 | 2 | 2 | 0 | 3 |
| 5 | 0 | 0 | 0 | 4 | 0 |

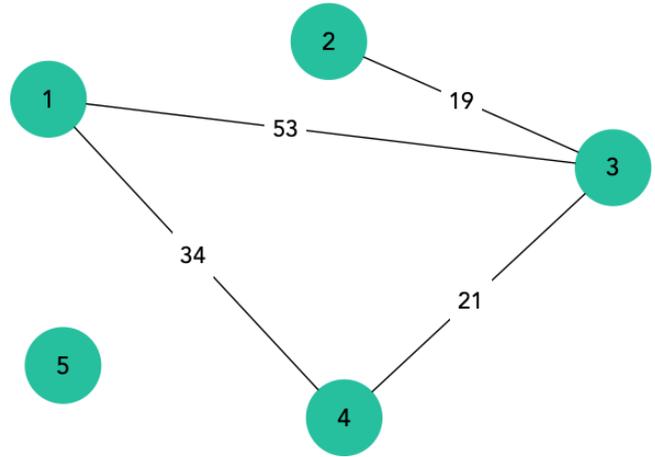
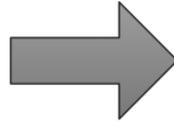
Dataset's Similarity Graph



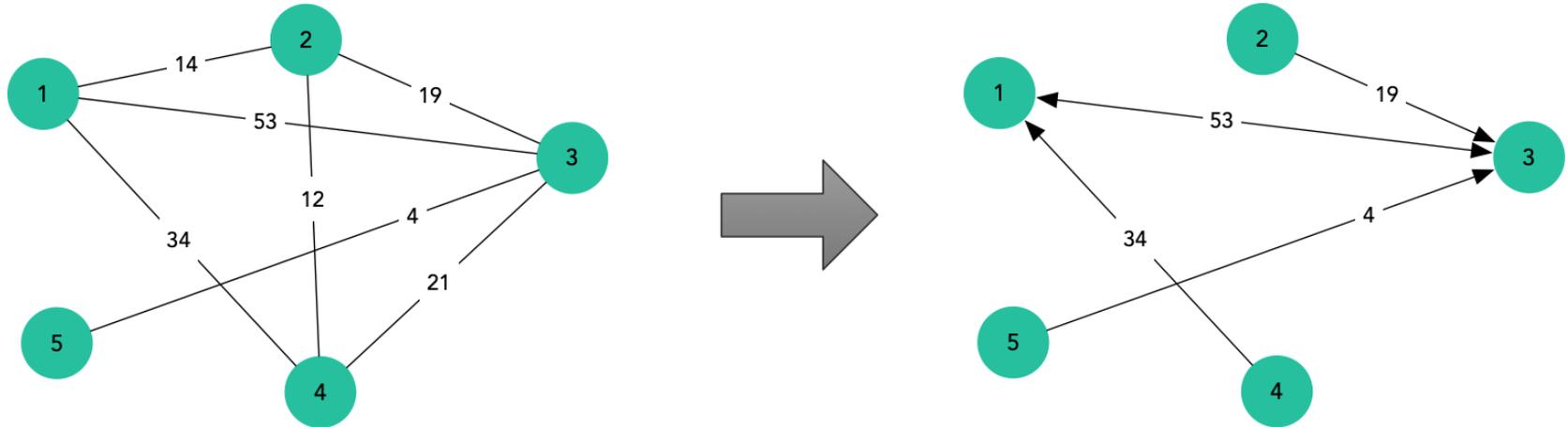
ϵ -graph



$\epsilon = 15$



Nearest Neighbour (NN) Graph



If we have k nearest neighbours for each node,
graph would be a k NN graph

Similarity-based graphs

- Each of these graphs lead to **different problems**
- Each of them have **different solutions**
- In our work, we focus on the **ϵ -similarity graph**, a similarity graph whose edges are above the ϵ threshold

Similarity Graph Construction Challenges

When we have to **build similarity graphs many times**
like in data streams for different snapshots or windows

$$O(n^2 d)$$

Scalability



Similarities should be computed
for **all pairs of entities**
based on **all attributes**

N
Entities

D Attributes

| UserID | Item1 | Item2 | Item3 | Item4 | Item5 |
|--------|-------|-------|-------|-------|-------|
| 1 | 3 | 0 | 2 | 0 | 10 |
| 2 | 4 | 5 | 1 | 0 | 0 |
| 3 | 1 | 3 | 0 | 1 | 5 |
| 4 | 0 | 2 | 2 | 0 | 3 |
| 5 | 0 | 0 | 0 | 4 | 0 |

Main Objective

Proposing an **efficient** and **effective** method
for **similarity graph construction** from **high-dimensional** data

Approaches

- Distributed solutions on MapReduce
- GPU-based solutions

- Efficient algorithmic optimizations
 - Using inverted index
 - Sampling/sketching based methods

Our work

Inverted index

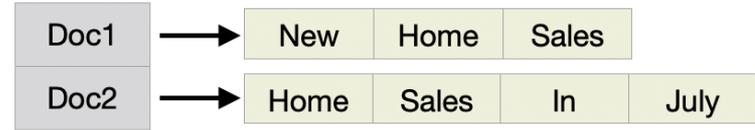


Doc1: New Home Sales

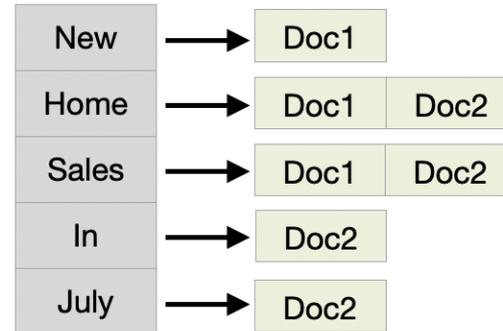


Doc2: Home Sales In July

The Forward Index



The Inverted Index



Data Sketching

Summarizing data that might be thought of as a **high dimensional** vector, or matrix

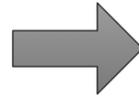
Data sketches have mathematically **proven error bounds**

Problem Statement

Sparse Vector Representation of Matrix

The Matrix

| | A1 | A2 | A3 | A4 | A5 |
|----|----|----|----|----|----|
| R1 | 5 | 1 | 0 | 0 | 4 |
| R2 | 0 | 0 | 3 | 0 | 0 |
| R3 | 0 | 0 | 2 | 0 | 1 |
| R4 | 3 | 2 | 0 | 0 | 4 |
| R5 | 0 | 3 | 0 | 0 | 2 |



Its Sparse Vector Representation

R1 = {(A1, 5), (A2, 1), (A5, 4)}

R2 = {(A3, 3)}

R3 = {(A3, 2), (A5, 1)}

R4 = {(A1, 3), (A2, 2), (A5, 4)}

R5 = {(A2, 3), (A5, 2)}

Approximate Similarity Graph Construction

- Given:
 - a **similarity threshold** ϵ
 - a **data matrix** or its **sparse vector representation**
- the problem is:
 - to build a **similarity graph** $G(V, E)$ where
 - V is the set of entities in the data matrix and
 - E is the set of edges representing the similarity between two nodes and
 - the **similarity is above the ϵ threshold**

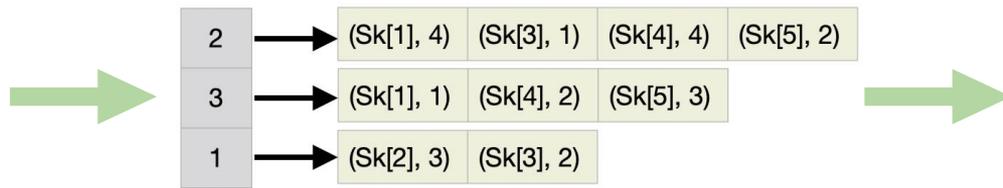
Proposed Methodology

Overview of the Algorithm

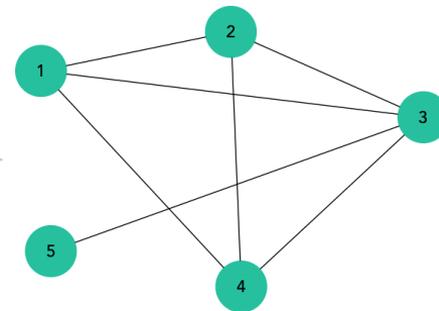
Step 1: Data Sketching

$Sk[1] = \{(2, 4), (3, 1)\}$
 $Sk[2] = \{(1, 3)\}$
 $Sk[3] = \{(1, 2), (2, 1)\}$
 $Sk[4] = \{(2, 4), (3, 2)\}$
 $Sk[5] = \{(2, 2), (3, 3)\}$

Step 2: Pairwise Similarity Computations



Step 3: Similarity Graph Construction



Step 1: Data Sketching

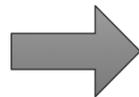
Main goal:

Start from a large dataset

Make it smaller

The dataset

| | 1 | 2 | 3 | 4 | 5 | ... | ... | ... | D |
|---|---|---|---|---|---|-----|-----|-----|---|
| 1 | █ | □ | █ | □ | □ | □ | □ | █ | □ |
| 2 | █ | □ | □ | █ | □ | □ | █ | □ | █ |
| 3 | □ | □ | □ | □ | □ | █ | □ | █ | □ |
| 4 | □ | □ | █ | □ | □ | □ | □ | █ | □ |
| 5 | □ | █ | □ | █ | █ | █ | █ | □ | □ |



Dataset's sketches

| 1 | █ | █ | █ | □ |
|---|---|---|---|---|
| 2 | █ | █ | █ | █ |
| 3 | █ | █ | █ | □ |
| 4 | █ | █ | █ | █ |
| 5 | █ | █ | █ | █ |

Sketch size: $k \ll d$

Step 1: Data Sketching

- Input
 - Two different kinds of dataset
 - A data matrix
 - Sparse vector representation of a data matrix
 - The sketch size (k)
- Output
 - Sketches of data

Data Sketching From a Data Matrix

Dataset: D

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | 5 | 1 | 0 | 0 | 4 |
| 2 | 0 | 0 | 3 | 0 | 0 |
| 3 | 0 | 0 | 2 | 0 | 1 |
| 4 | 3 | 2 | 0 | 0 | 4 |
| 5 | 0 | 3 | 0 | 0 | 2 |

Random
Column Permutation



[1, 2, 3, 4, 5]
To
[4, 3, 1, 5, 2]

Permuted Dataset

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | 0 | 4 | 1 | 5 | 0 |
| 2 | 3 | 0 | 0 | 0 | 0 |
| 3 | 2 | 1 | 0 | 0 | 0 |
| 4 | 0 | 4 | 2 | 3 | 0 |
| 5 | 0 | 2 | 3 | 0 | 0 |

Taking the first k
column ID/value pairs
with nonzero values



Recording the highest
column ID in each
sketch



Sketches of D
Sketch size (K) = 2

Sk[1] = {(2, 4), (3, 1)}
Sk[2] = {(1, 3)}
Sk[3] = {(1, 2), (2, 1)}
Sk[4] = {(2, 4), (3, 2)}
Sk[5] = {(2, 2), (3, 3)}

Sk_max_id[1] = 3
Sk_max_id[2] = 1
Sk_max_id[3] = 2
Sk_max_id[4] = 3
Sk_max_id[5] = 3

**Sketch Max
Column IDs**

Data Sketching From Sparse Vector Representation

Original Dataset

R1 = {(1, 5), (2, 1), (5, 4)}
R2 = {(3, 3)}
R3 = {(3, 2), (5, 1)}
R4 = {(1, 3), (2, 2), (5, 4)}
R5 = {(2, 3), (5, 2)}

Column Permutation



[1, 2, 3, 4, 5]
To
[4, 3, 1, 5, 2]

Permuted Dataset

R1 = {(4, 5), (3, 1), (2, 4)}
R2 = {(1, 3)}
R3 = {(1, 2), (2, 1)}
R4 = {(4, 3), (3, 2), (2, 4)}
R5 = {(3, 3), (2, 2)}

Sorting



R1 = {(2, 4), (3, 1), (4, 5)}
R2 = {(1, 3)}
R3 = {(1, 2), (2, 1)}
R4 = {(2, 4), (3, 2), (4, 3)}
R5 = {(2, 2), (3, 3)}

Sketch size (K) = 2



Sketch Max IDs

Sk_max_id[1] = 3
Sk_max_id[2] = 1
Sk_max_id[3] = 2
Sk_max_id[4] = 3
Sk_max_id[5] = 3

Sketches

Sk[1] = {(2, 4), (3, 1)}
Sk[2] = {(1, 3)}
Sk[3] = {(1, 2), (2, 1)}
Sk[4] = {(2, 4), (3, 2)}
Sk[5] = {(2, 2), (3, 3)}

Data Sketching - Time Complexity

- From Data Matrix: $O(d) + O(nm)$
 - dimensionality (points to d)
 - # entities (points to m)
 - average position of k^{th} nonzero entry (points to m)
- From Sparse Vector Rep.: $O(d) + O(n \times (f + f \cdot \log(f) + k))$
 - average number of nonzero entries (points to f)
 - sketch size (points to k)
- $O(n \times (f + f \cdot \log(f) + k)) < O(nm)$

Step 2: Pairwise similarity computations

Input: Data sketches

Output: Similarities of all pairs of sketches

Our similarity measure is the **inner-product**

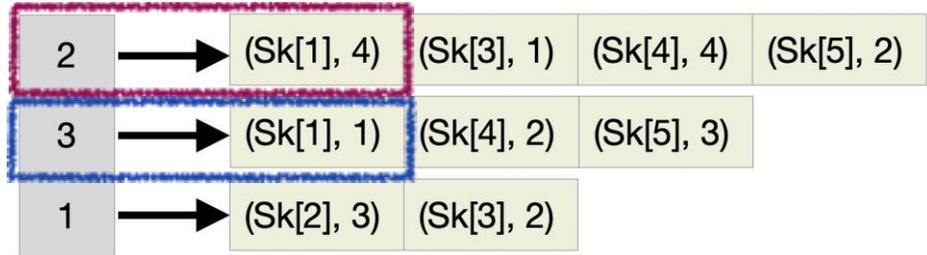
Structure of An Inverted Index

Sketches of D ($k = 2$)

$Sk[1] = \{(2, 4), (3, 1)\}$
 $Sk[2] = \{(1, 3)\}$
 $Sk[3] = \{(1, 2), (2, 1)\}$
 $Sk[4] = \{(2, 4), (3, 2)\}$
 $Sk[5] = \{(2, 2), (3, 3)\}$



The Inverted Index



Pairwise similarity computations

Sketches

$$Sk[1] = \{(2, 4), (3, 1)\}$$

$$Sk[2] = \{(1, 3)\}$$

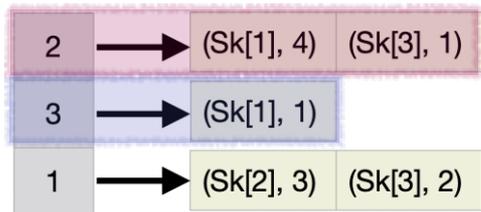
$$Sk[3] = \{(1, 2), (2, 1)\}$$

$$Sk[4] = \{(2, 4), (3, 2)\}$$

$$Sk[5] = \{(2, 2), (3, 3)\}$$

Current Sketch
Iterating Over

Inverted Index



$$(Sk[1], 16), (Sk[3], 4)$$

$$(Sk[1], 16+2)$$

Similar Sketches of Sk[4]

$$(Sk[1], 18), (Sk[3], 4)$$

Effective sample size (d_s) Computation

Our similarity measure is the **inner-product**

Original data

Dimensionality

Rows of original dataset

$$a = \sum_{i=1}^d u_{1,i} u_{2,i}$$

Sketches

Sketches

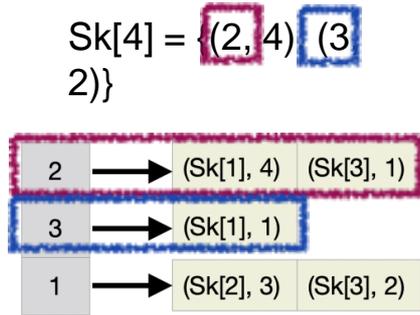
$$\hat{a} = \frac{d}{d_s} \sum_{i=1}^{d_s} \tilde{u}_{1,i} \tilde{u}_{2,i}$$

Computed pairwise

$$d_s = \min(\max(ID(\tilde{u}_1)), \max(ID(\tilde{u}_2)))$$

Three approaches for d_S computation

Online



Offline via Sorting

Sorting sketches based on their max ID

$Sk[2] = \{(1, 3)\}$
 $Sk[3] = \{(1, 2), (2, 1)\}$
 $Sk[1] = \{(2, 4), (3, 1)\}$
 $Sk[4] = \{(2, 4), (3, 2)\}$
 $Sk[5] = \{(2, 2), (3, 3)\}$

Ascending

Offline via Matrix Precomputation

d_S pairwise matrix

| | sk[1] | sk[2] | sk[3] | sk[4] | sk[5] |
|-------|-------|-------|-------|-------|-------|
| sk[1] | | | | | |
| sk[2] | | | | | |
| sk[3] | | | | | |
| sk[4] | | | | | |
| sk[5] | | | | | |

Preprocessing

Extra Space Needed for d_S computations

- Online: $O(1)$
- Offline via sorting: $O(n)$
- Offline via matrix precomputations: $O(n^2)$

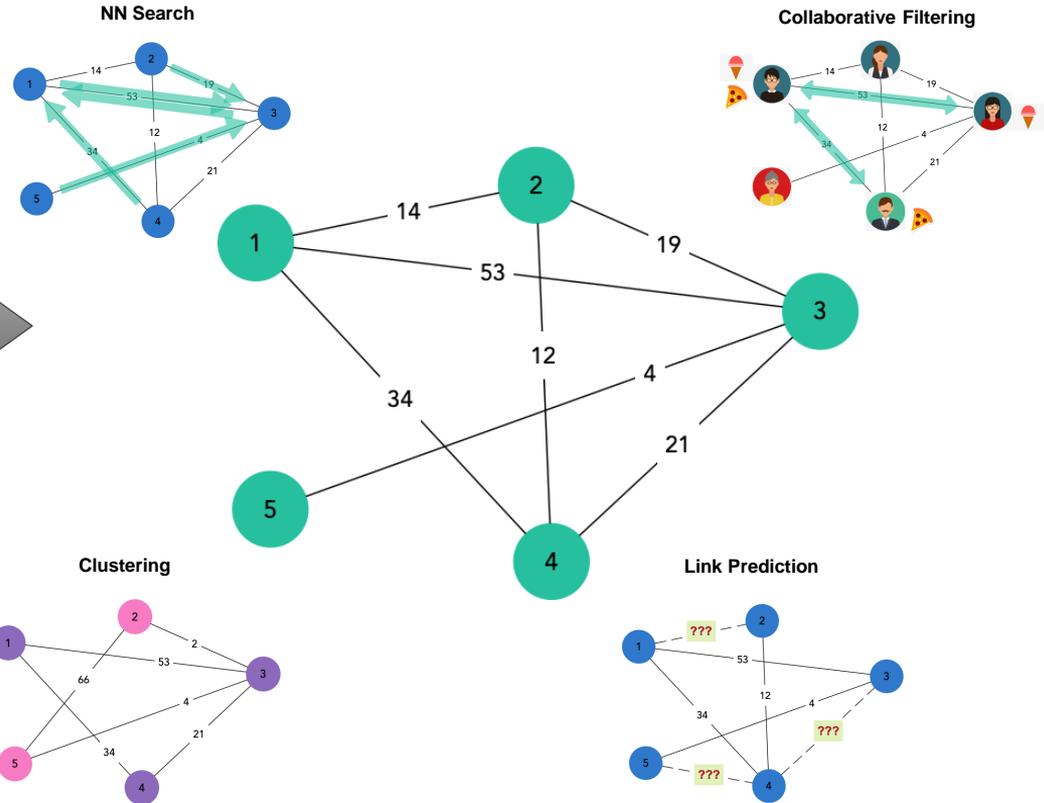
Pairwise Similarity Computation - Time Complexity

- Online: $O(\underbrace{c}_{\substack{\text{\# operations} \\ \text{in the inner-most loop}}} \times \underbrace{n}_{\substack{\text{\# entities}} \times \underbrace{k}_{\substack{\text{sketch size}}} \times \underbrace{l}_{\substack{\text{average size of lists} \\ \text{in the inverted index}}})$
- Offline via sorting: $O(c' \times n \times k \times l) + O(n \cdot \log(n))$
- Offline via matrix precomputation: $O(c'' \times n \times k \times l) + O(n^2)$

$$c', c'' < c$$

Step 3: Similarity Graph Construction

Similarities of all pairs of entities



Evaluation

Evaluation Scenarios

- Runtime Cost
- Accuracy
- Effectiveness

Datasets

| Name | Distribution | Density | Size |
|-----------|---------------------------------|---------|------------------|
| Normal2 | Normal($\mu=200, \sigma=50$) | 2% | 10k \times 10k |
| Normal4 | Normal($\mu=400, \sigma=100$) | 4% | |
| Normal6 | Normal($\mu=600, \sigma=100$) | 6% | |
| Normal8 | Normal($\mu=800, \sigma=100$) | 8% | |
| Binomial2 | Binomial($n=10k, p=0.02$) | 2% | 10k \times 10k |
| Binomial4 | Binomial($n=10k, p=0.04$) | 4% | |
| Binomial6 | Binomial($n=10k, p=0.06$) | 6% | |
| Binomial8 | Binomial($n=10k, p=0.08$) | 8% | |

Evaluated methods

- Original



| UserID | Item1 | Item2 | Item3 | Item4 | Item5 |
|--------|-------|-------|-------|-------|-------|
| 1 | 3 | 0 | 2 | 0 | 10 |
| 2 | 4 | 5 | 1 | 0 | 0 |
| 3 | 1 | 3 | 0 | 1 | 5 |
| 4 | 0 | 2 | 2 | 0 | 3 |
| 5 | 0 | 0 | 0 | 4 | 0 |



sim. operations on the **vectors**
of the original dataset

- Normal Random Projection (NRP)



$$A \in \mathbb{R}^{n \times D} \quad R \in \mathbb{R}^{D \times k}$$
$$AR \in \mathbb{R}^{n \times k}$$

- Sk_naive



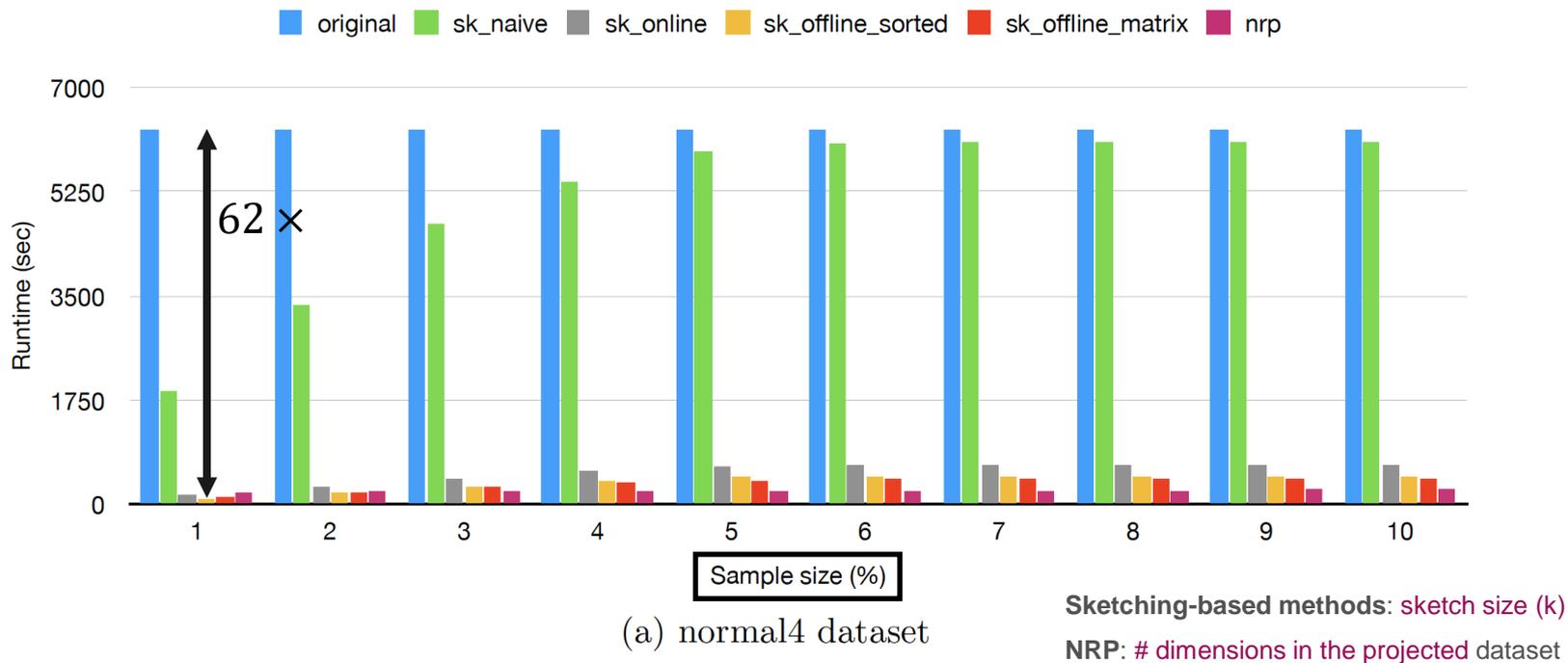
Data sketching
without using an inverted index

- Sk_online
- Sk_offline_sorted
- Sk_offline_matrix



Our proposed methods

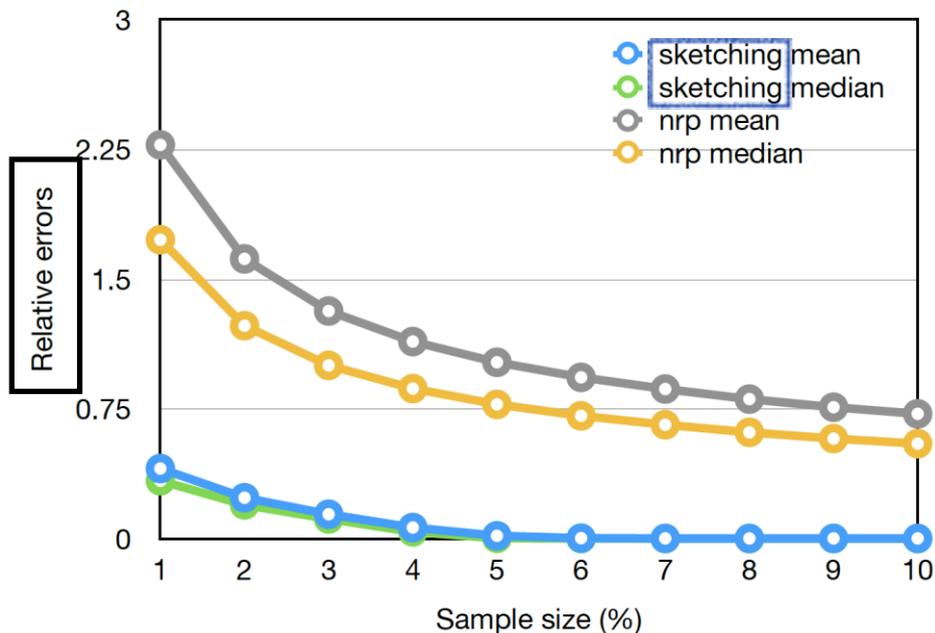
Runtime vs. sample size



Accuracy vs. sample size

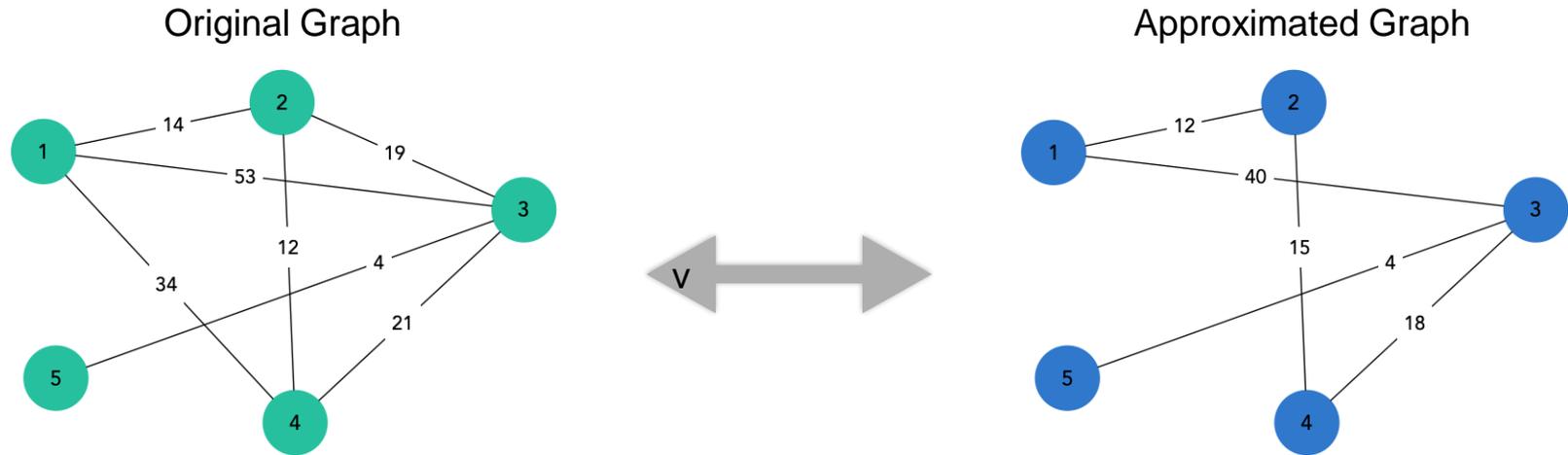
Relative error for each pair (i, j) :

$$e_{ij} = \frac{|s_{ij}^{sample} - s_{ij}^{original}|}{s_{ij}^{original}}$$



(a) normal4 dataset

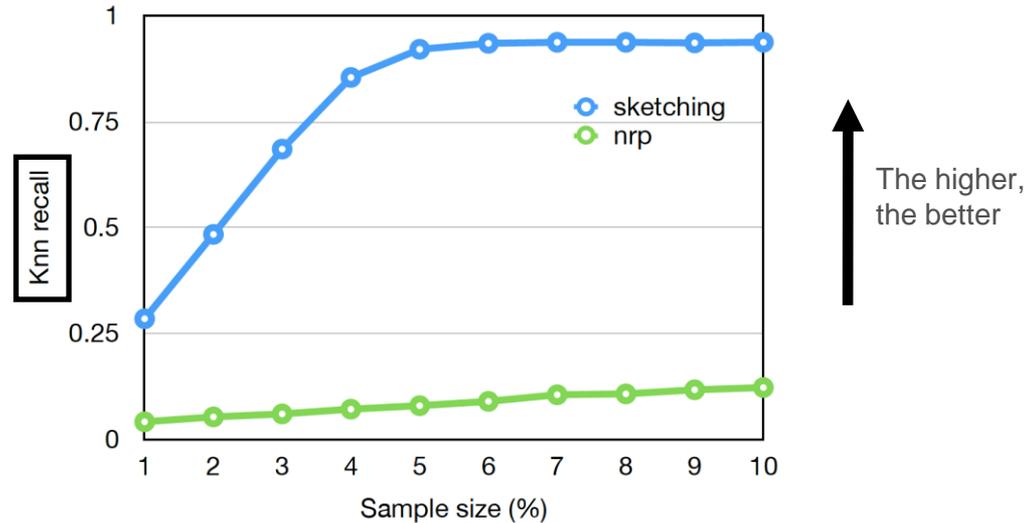
Effectiveness



Effectiveness

- K -nearest neighbours
- Node centrality values
 - show **importance** of the nodes in the graph
 - we work with **eigenvector centrality**
- Node rankings
 - Based on the centrality of the nodes in the graph, we have a ranking for them
 - We use **Spearman's ranking correlation coefficient (ρ)**

Effectiveness - kNN



(a) normal4 dataset

kNN recall:

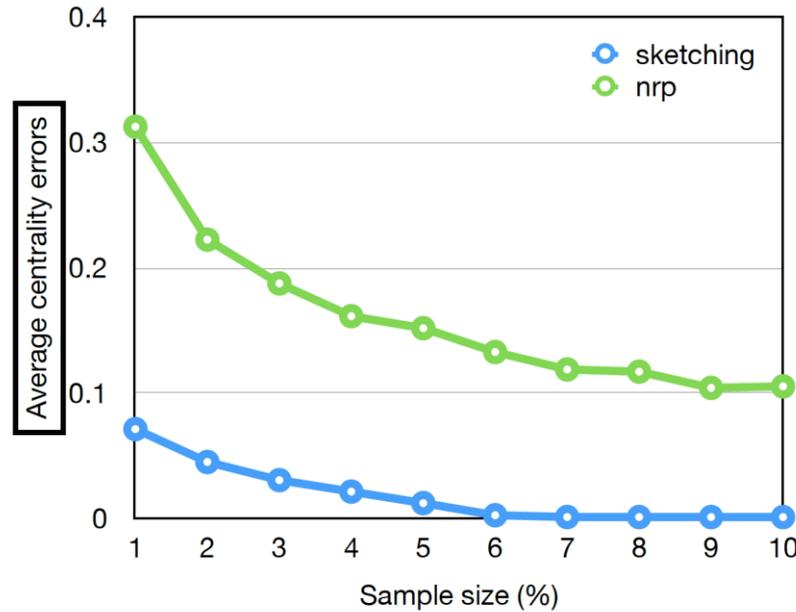
How many of the k real nearest neighbours of each node we are returning

Precision and recall are the same in this case

Effectiveness - Centrality Errors

Centrality error for node i :

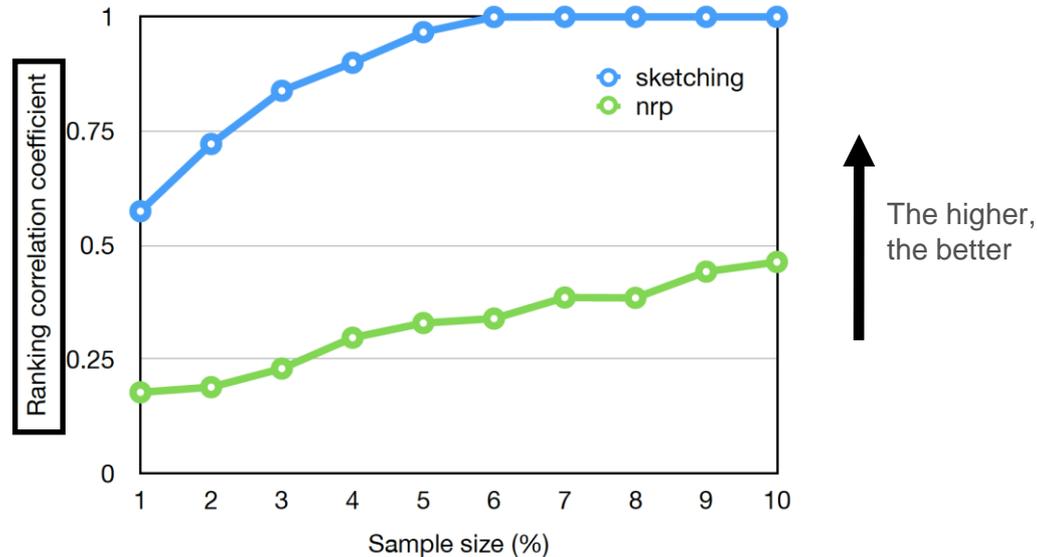
$$e_i = \frac{|c_i^{sample} - c_i^{original}|}{c_i^{original}}$$



The lower,
the better

(a) normal4 dataset

Effectiveness - Node Ranking Correlations



(a) normal4 dataset

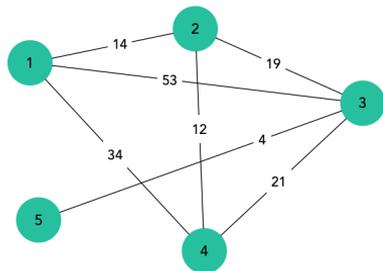
We use Spearman's ranking correlation coefficient:

A measure to see how well node rankings of the approximated graphs are compared to the original.

Conclusion & Future Work

Summary of contributions

building similarity graphs
from high-dimensional data



$$O(n^2 d)$$

Data Sketching



Inverted Index

| | | | |
|---|---|------------|------------|
| 2 | → | (Sk[1], 4) | (Sk[3], 1) |
| 3 | → | (Sk[1], 1) | |
| 1 | → | (Sk[2], 3) | (Sk[3], 2) |

efficient, accurate and effective way of

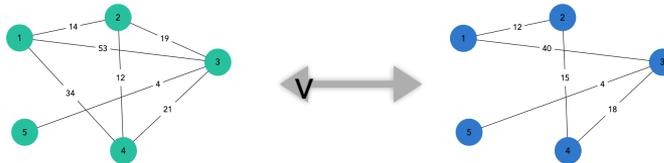
construction of similarity graphs

Summary of contributions

three algorithms each of which has a trade-off for speed and space

Time efficient
Accurate

Effective on **different graph analysis tasks**



Limitations

Scalability for very large datasets

Storage overhead for the inverted index + **maintenance cost**

Future Work

Providing **theoretical bounds**
for the quality of graph downstream task results

Making the methods **distributed**
to increase scalability

Working with **multi-dimensional arrays** instead of matrix

Thank You!

Appendix

Online Pairwise Similarities

Algorithm 2: Online Pairwise Similarities

Input: sketches of D : \mathbf{Sk} , sketch max ids: $\mathbf{Sk_Max_Id}$, similarity threshold: ϵ ,
dimensionality of original data: d

Output: pairwise similarities: \mathbf{S}

```
1  $S \leftarrow \emptyset$ 
2  $I_1, I_2, \dots, I_d \leftarrow \emptyset$  ( $I_i$  is the entry for the  $i$ th attribute in the inverted index. It will
   contain a list of  $(x, v)$ 's, where  $x$  is a row/sketch id and  $v$  is the value of  $x$  for the
    $i$ th attribute.)
3 for  $x \in \mathbf{Sk}$  do
4    $M \leftarrow \emptyset$  ( $M$  holds the similarity values between  $x$  and each of the sketches before
    $x$  in  $\mathbf{Sk}$ )
5   for  $(a, v) \in x$  do
6     for  $(y, y_v) \in I_a$  do
7        $d_s = \min(\mathbf{Sk\_Max\_Id}[x], \mathbf{Sk\_Max\_Id}[y])$ 
8       if  $a \leq d_s$  then
9          $M[y] \leftarrow M[y] + (d/d_s) \cdot v \cdot y_v$ 
10       $I_a \leftarrow I_a \cup \{(x, v)\}$ 
11    $S_x \leftarrow \mathbf{Filter\_Similarities}(M, \epsilon)$  (Remove similarities in  $M$  whose value is less
   than  $\epsilon$ )
12    $S \leftarrow S \cup (x, S_x)$ 
13 return  $\mathbf{S}$ 
```

Offline Pairwise Similarities via Sorting

Algorithm 3: Offline Pairwise Similarities via Sorting

Input: sketches of D : \mathbf{Sk} , sketch max ids: $\mathbf{Sk_Max_Id}$, similarity threshold: ϵ ,
dimensionality of original data: d

Output: pairwise similarities: \mathcal{S}

```
1  $S \leftarrow \emptyset$ 
2  $I_1, I_2, \dots, I_d \leftarrow \emptyset$ 
3  $Sorted\_Sk\_Indices \leftarrow \mathbf{arg\_sort}(Sk\_Max\_Id, ascending)$ 
4 for  $i \in range(|Sk|)$  do
5    $x_{id} = Sorted\_Sk\_Indices[i]$ 
6    $x \leftarrow Sk[x_{id}]$ 
7    $M \leftarrow \emptyset$ 
8   for  $(a, v) \in x$  do
9     for  $(y, y_v) \in I_a$  do
10       $d_s = Sk\_Max\_Id[y]$ 
11       $M[y] \leftarrow M[y] + (d/d_s) \cdot v \cdot y_v$ 
12       $I_a \leftarrow I_a \cup \{(x, v)\}$ 
13    $S_x \leftarrow \mathbf{Filter\_Similarities}(M, \epsilon)$ 
14    $S \leftarrow S \cup (x, S_x)$ 
15 return  $S$ 
```

Offline Pairwise Similarities via Matrix Precomputations

Algorithm 4: Offline Pairwise Similarities via Matrix Precomputations

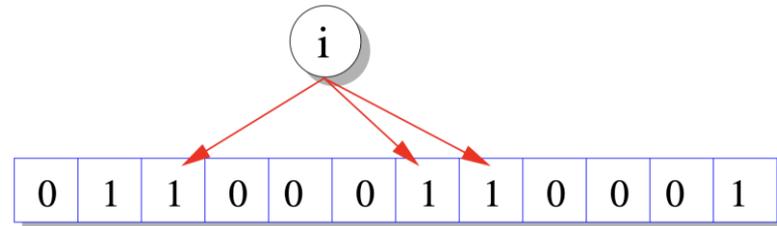
Input: sketches of D : \mathbf{Sk} , sketch max ids: $\mathbf{Sk_Max_Id}$, similarity threshold: ϵ ,
dimensionality of original data: d

Output: pairwise similarities: \mathbf{S}

```
1  $S \leftarrow \emptyset$ 
2  $I_1, I_2, \dots, I_d \leftarrow \emptyset$ 
3  $n \leftarrow |\mathbf{Sk}|$ 
4 for  $i \in \text{range}(n)$  do
5   for  $j \in \text{range}(n)$  do
6      $d_s\text{-pairwise}[i][j] \leftarrow \min(\mathbf{Sk\_Max\_Id}[i], \mathbf{Sk\_Max\_Id}[j])$ 
7 for  $x \in \mathbf{Sk}$  do
8    $M \leftarrow \emptyset$ 
9   for  $(a, v) \in x$  do
10    for  $(y, y_v) \in I_a$  do
11       $d_s = d_s\text{-pairwise}[x][y]$ 
12      if  $a < d_s$  then
13         $M[y] \leftarrow M[y] + (d/ds) \cdot v \cdot y_v$ 
14     $I_a \leftarrow I_a \cup \{(x, v)\}$ 
15   $S_x \leftarrow \mathbf{Filter\_Similarities}(M, \epsilon)$ 
16   $S \leftarrow S \cup (x, S_x)$ 
17 return  $\mathbf{S}$ 
```

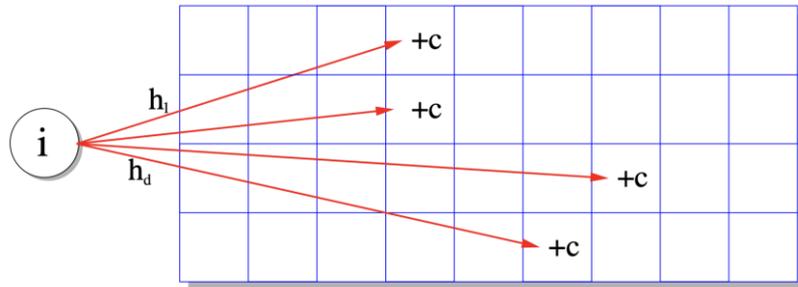
Bloom filters for set summarization

- Set membership
- The item has *definitely not been stored*, or the item has *probably been stored*
- Having k hash functions and map each item with each of them
- Set all the corresponding bits to 1. If all were one for an item, say it is a member, if any of them were 0, say it is not a member



Counting with count-min sketch

- Counts the number of items of a certain type
- Sketch: an array of counters, and a set of hash functions which map items into the array
- Count of the desired item is to take the smallest of counters in each row as our estimate.



PCA vs. RP

- PCA
 - Extracting a small number of directions from the data which captures most of variation of dataset
 - Finding the direction requires finding **eigenvectors of the covariance matrix** 
Substantial amount of work
- Random projection
 - Rather than finding “the best” directions, it suffices to use random vectors.
 - Picking a moderate number of random directions captures a *comparable amount of variation*, while requiring *much less computation*.