

# Similarity Estimation

# Similarity Between Data Points

# Similarity Between Data Points

- Similarity estimation arises in many big data applications

# Similarity Between Data Points

- Similarity estimation arises in many big data applications
  - Find near-duplicate documents/webpages to remove redundancy

# Similarity Between Data Points

- Similarity estimation arises in many big data applications
  - Find near-duplicate documents/webpages to remove redundancy
  - Detect pirated files

# Similarity Between Data Points

- Similarity estimation arises in many big data applications
  - Find near-duplicate documents/webpages to remove redundancy
  - Detect pirated files
- How do we define and compute/estimate similarity?

# Jaccard Similarity

# Jaccard Similarity

- An article may be represented by the set of words it contains



# Jaccard Similarity

- An article may be represented by the set of words it contains
  - Such a representation is high-dimensional but sparse

# Jaccard Similarity

- An article may be represented by the set of words it contains
  - Such a representation is high-dimensional but sparse

• **Def.** The **Jaccard similarity** between two sets  $S$  and  $T$  is  $\frac{|S \cap T|}{|S \cup T|}$ , denoted as  $\text{SIM}(S, T)$ .

# Jaccard Similarity

- An article may be represented by the set of words it contains
  - Such a representation is high-dimensional but sparse
- **Def.** The **Jaccard similarity** between two sets  $S$  and  $T$  is  $\frac{|S \cap T|}{|S \cup T|}$ , denoted as  $\text{SIM}(S, T)$ .
- We consider  $S, T$  very similar if  $\text{SIM}(S, T) \geq \alpha$  for some threshold  $\alpha$ .

# Jaccard Similarity

- An article may be represented by the set of words it contains
  - Such a representation is high-dimensional but sparse
- **Def.** The **Jaccard similarity** between two sets  $S$  and  $T$  is  $\frac{|S \cap T|}{|S \cup T|}$ , denoted as  $\text{SIM}(S, T)$ .
- We consider  $S, T$  very similar if  $\text{SIM}(S, T) \geq \alpha$  for some threshold  $\alpha$ .
- Is there a quick way to estimate  $\text{SIM}(S, T)$ ?

# Min Hashing

# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .

# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .
- For a permutation  $\sigma$  on  $[n]$  and  $S \subseteq [n]$ , define  $\sigma_{\min}(S) := \min_{i \in S} \sigma(i)$

# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .
- For a permutation  $\sigma$  on  $[n]$  and  $S \subseteq [n]$ , define  $\sigma_{\min}(S) := \min_{i \in S} \sigma(i)$

**Lemma.** For  $S, T \subseteq [n]$ , if  $\sigma$  is a random permutation of  $[n]$ , then  $\mathbb{P}[\sigma_{\min}(S) = \sigma_{\min}(T)] = \frac{|S \cap T|}{|S \cup T|} = \text{SIM}(S, T)$



# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .
- For a permutation  $\sigma$  on  $[n]$  and  $S \subseteq [n]$ , define  $\sigma_{\min}(S) := \min_{i \in S} \sigma(i)$

**Lemma.** For  $S, T \subseteq [n]$ , if  $\sigma$  is a random permutation of  $[n]$ , then  $\mathbb{P}[\sigma_{\min}(S) = \sigma_{\min}(T)] = \frac{|S \cap T|}{|S \cup T|} = \text{SIM}(S, T)$

- A data structure that estimates Jaccard similarity:

# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .
- For a permutation  $\sigma$  on  $[n]$  and  $S \subseteq [n]$ , define  $\sigma_{\min}(S) := \min_{i \in S} \sigma(i)$

**Lemma.** For  $S, T \subseteq [n]$ , if  $\sigma$  is a random permutation of  $[n]$ , then  $\mathbb{P}[\sigma_{\min}(S) = \sigma_{\min}(T)] = \frac{|S \cap T|}{|S \cup T|} = \text{SIM}(S, T)$

- A data structure that estimates Jaccard similarity:
  - Take  $\ell$  random permutations  $\sigma^1, \dots, \sigma^\ell$  of  $[n]$

# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .
- For a permutation  $\sigma$  on  $[n]$  and  $S \subseteq [n]$ , define  $\sigma_{\min}(S) := \min_{i \in S} \sigma(i)$

**Lemma.** For  $S, T \subseteq [n]$ , if  $\sigma$  is a random permutation of  $[n]$ , then  $\mathbb{P}[\sigma_{\min}(S) = \sigma_{\min}(T)] = \frac{|S \cap T|}{|S \cup T|} = \text{SIM}(S, T)$

- A data structure that estimates Jaccard similarity:
  - Take  $\ell$  random permutations  $\sigma^1, \dots, \sigma^\ell$  of  $[n]$
  - For each record  $S \subseteq [n]$ , store  $(\sigma_{\min}^1(S), \dots, \sigma_{\min}^\ell(S))$

# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .
- For a permutation  $\sigma$  on  $[n]$  and  $S \subseteq [n]$ , define  $\sigma_{\min}(S) := \min_{i \in S} \sigma(i)$

**Lemma.** For  $S, T \subseteq [n]$ , if  $\sigma$  is a random permutation of  $[n]$ , then  $\mathbb{P}[\sigma_{\min}(S) = \sigma_{\min}(T)] = \frac{|S \cap T|}{|S \cup T|} = \text{SIM}(S, T)$

- A data structure that estimates Jaccard similarity:
  - Take  $\ell$  random permutations  $\sigma^1, \dots, \sigma^\ell$  of  $[n]$
  - For each record  $S \subseteq [n]$ , store  $(\sigma_{\min}^1(S), \dots, \sigma_{\min}^\ell(S))$

To estimate  $\text{SIM}(S, T)$ , output  $\frac{|\{i : \sigma_{\min}^i(S) = \sigma_{\min}^i(T)\}|}{\ell}$

# Min Hashing

- Let the vocabulary be  $[n] := \{1, \dots, n\}$ .
- For a permutation  $\sigma$  on  $[n]$  and  $S \subseteq [n]$ , define  $\sigma_{\min}(S) := \min_{i \in S} \sigma(i)$

**Lemma.** For  $S, T \subseteq [n]$ , if  $\sigma$  is a random permutation of  $[n]$ , then  $\mathbb{P}[\sigma_{\min}(S) = \sigma_{\min}(T)] = \frac{|S \cap T|}{|S \cup T|} = \text{SIM}(S, T)$

- A data structure that estimates Jaccard similarity:
  - Take  $\ell$  random permutations  $\sigma^1, \dots, \sigma^\ell$  of  $[n]$
  - For each record  $S \subseteq [n]$ , store  $(\sigma_{\min}^1(S), \dots, \sigma_{\min}^\ell(S))$

To estimate  $\text{SIM}(S, T)$ , output  $\frac{|\{i : \sigma_{\min}^i(S) = \sigma_{\min}^i(T)\}|}{\ell}$

How large should  $\ell$  be for  $(1 \pm \epsilon)$ -approximation w.p.  $1 - \delta$ ?

# Minwise Independent Permutations

# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive

# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive
- Recurring idea: use pseudo-randomness (e.g. from hash families)



# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive
- Recurring idea: use pseudo-randomness (e.g. from hash families)
  - Let  $S_n$  be the set of permutations of  $[n]$ .

# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive
- Recurring idea: use pseudo-randomness (e.g. from hash families)
  - Let  $S_n$  be the set of permutations of  $[n]$ .
  - Think of the definition of  $k$ -wise independent hash families

# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive
- Recurring idea: use pseudo-randomness (e.g. from hash families)
  - Let  $S_n$  be the set of permutations of  $[n]$ .
  - Think of the definition of  $k$ -wise independent hash families

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a *minwise independent family* of permutations if for every  $S \subseteq [n]$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\mathbb{P}[\sigma_{\min}(S) = \sigma(a)] = \frac{1}{|S|}$ .

# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive
- Recurring idea: use pseudo-randomness (e.g. from hash families)
  - Let  $S_n$  be the set of permutations of  $[n]$ .
  - Think of the definition of  $k$ -wise independent hash families

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a *minwise independent family* of permutations if for every  $S \subseteq [n]$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\mathbb{P}[\sigma_{\min}(S) = \sigma(a)] = \frac{1}{|S|}$ .

Minwise independent families suffice for estimation of Jaccard similarity.

# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive
- Recurring idea: use pseudo-randomness (e.g. from hash families)
  - Let  $S_n$  be the set of permutations of  $[n]$ .
  - Think of the definition of  $k$ -wise independent hash families

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a *minwise independent family* of permutations if for every  $S \subseteq [n]$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\mathbb{P}[\sigma_{\min}(S) = \sigma(a)] = \frac{1}{|S|}$ .

Minwise independent families suffice for estimation of Jaccard similarity.

Minwise independent families of size  $4^n$  exist ( $\ll |S_n| = n!$ )

# Minwise Independent Permutations

- Storing random permutations and computing  $\sigma_{\min}(S)$  is expensive
- Recurring idea: use pseudo-randomness (e.g. from hash families)
  - Let  $S_n$  be the set of permutations of  $[n]$ .
  - Think of the definition of  $k$ -wise independent hash families

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a *minwise independent family* of permutations if for every  $S \subseteq [n]$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\mathbb{P}[\sigma_{\min}(S) = \sigma(a)] = \frac{1}{|S|}$ .

Minwise independent families suffice for estimation of Jaccard similarity.

Minwise independent families of size  $4^n$  exist ( $\ll |S_n| = n!$ )

Any minwise independent family has size  $e^{(1-o(1))n}$

# Relaxation

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a *minwise independent family* of permutations if for every  $S \subseteq [n]$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\mathbb{P}[\sigma_{\min}(S) = \sigma(a)] = \frac{1}{|S|}$ .

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a  $(\epsilon, k)$  *minwise independent family* of permutations if for every  $S \subseteq [n]$  with  $|S| \leq k$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\frac{1 - \epsilon}{|S|} \leq \mathbb{P}[\sigma_{\min}(S) = \sigma(a)] \leq \frac{1 + \epsilon}{|S|}$ .

# Relaxation

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a *minwise independent family* of permutations if for every  $S \subseteq [n]$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\mathbb{P}[\sigma_{\min}(S) = \sigma(a)] = \frac{1}{|S|}$ .

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a  $(\epsilon, k)$  *minwise independent family* of permutations if for every  $S \subseteq [n]$  with  $|S| \leq k$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\frac{1 - \epsilon}{|S|} \leq \mathbb{P}[\sigma_{\min}(S) = \sigma(a)] \leq \frac{1 + \epsilon}{|S|}$ .

**Thm.** [Indyk] Let  $\mathcal{H}$  be a  $t$ -wise independent hash family from  $[n]$  to  $[n]$ , with  $t = \Omega(\log \frac{1}{\epsilon})$ , then  $\mathcal{H}$  is a  $(\epsilon, k)$  minwise independent family of permutations for  $k = O(\epsilon n)$ .



# Relaxation

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a *minwise independent family* of permutations if for every  $S \subseteq [n]$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\mathbb{P}[\sigma_{\min}(S) = \sigma(a)] = \frac{1}{|S|}$ .

**Def.** [Broder, Charikar, Frieze, Mitzenmacher] A family  $\mathcal{F} \subseteq S_n$  is a  $(\epsilon, k)$  *minwise independent family* of permutations if for every  $S \subseteq [n]$  with  $|S| \leq k$  and any  $a \in S$ , for  $\sigma$  sampled uniformly from  $\mathcal{F}$ ,  $\frac{1 - \epsilon}{|S|} \leq \mathbb{P}[\sigma_{\min}(S) = \sigma(a)] \leq \frac{1 + \epsilon}{|S|}$ .

**Thm.** [Indyk] Let  $\mathcal{H}$  be a  $t$ -wise independent hash family from  $[n]$  to  $[n]$ , with  $t = \Omega(\log \frac{1}{\epsilon})$ , then  $\mathcal{H}$  is a  $(\epsilon, k)$  minwise independent family of permutations for  $k = O(\epsilon n)$ .

How do you use minwise independent family to sample near-uniformly from the distinct elements in a streaming input?

# Angular Distance

# Angular Distance

- Two vectors in  $\mathbb{R}^d$  may be considered similar if they point roughly the same direction

# Angular Distance

- Two vectors in  $\mathbb{R}^d$  may be considered similar if they point roughly the same direction
- For  $u, v \in \mathbb{R}^d$ , let  $\text{dist}(u, v) := \frac{\theta(u, v)}{\pi}$ , where  $\theta(u, v)$  is the angle between  $u$  and  $v$

# Angular Distance

- Two vectors in  $\mathbb{R}^d$  may be considered similar if they point roughly the same direction
- For  $u, v \in \mathbb{R}^d$ , let  $\text{dist}(u, v) := \frac{\theta(u, v)}{\pi}$ , where  $\theta(u, v)$  is the angle between  $u$  and  $v$
- If  $u, v$  are unit vectors, then  $\cos(\theta(u, v)) = u \cdot v$

# Angular Distance

- Two vectors in  $\mathbb{R}^d$  may be considered similar if they point roughly the same direction
- For  $u, v \in \mathbb{R}^d$ , let  $\text{dist}(u, v) := \frac{\theta(u, v)}{\pi}$ , where  $\theta(u, v)$  is the angle between  $u$  and  $v$ 
  - If  $u, v$  are unit vectors, then  $\cos(\theta(u, v)) = u \cdot v$
- Similarity between  $u, v$  is  $1 - \text{dist}(u, v)$

# Sim Hash

# Sim Hash

- Can we estimate the similarity without computing the full inner product?



# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$
  - For each  $v_i$ , let  $h_r(v_i) := \text{sign}(r \cdot v_i)$

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$
  - For each  $v_i$ , let  $h_r(v_i) := \text{sign}(r \cdot v_i)$

Recall: how do we sample a random direction in  $\mathbb{R}^d$ ?

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$
  - For each  $v_i$ , let  $h_r(v_i) := \text{sign}(r \cdot v_i)$

Recall: how do we sample a random direction in  $\mathbb{R}^d$ ?

**Lemma.** For  $u, v \in \mathbb{R}^d$ ,  $\mathbb{P}[h_r(u) = h_r(v)] = \frac{\theta(u, v)}{\pi}$

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$
  - For each  $v_i$ , let  $h_r(v_i) := \text{sign}(r \cdot v_i)$

Recall: how do we sample a random direction in  $\mathbb{R}^d$ ?

**Lemma.** For  $u, v \in \mathbb{R}^d$ ,  $\mathbb{P}[h_r(u) = h_r(v)] = \frac{\theta(u, v)}{\pi}$

Proof: The projection of  $r$  onto the plane spanned by  $u, v$  is again a random direction.

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$
  - For each  $v_i$ , let  $h_r(v_i) := \text{sign}(r \cdot v_i)$

Recall: how do we sample a random direction in  $\mathbb{R}^d$ ?

**Lemma.** For  $u, v \in \mathbb{R}^d$ ,  $\mathbb{P}[h_r(u) = h_r(v)] = \frac{\theta(u, v)}{\pi}$

Proof: The projection of  $r$  onto the plane spanned by  $u, v$  is again a random direction.

- SimHash: a data structure that estimates angular similarity:



# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$
  - For each  $v_i$ , let  $h_r(v_i) := \text{sign}(r \cdot v_i)$

Recall: how do we sample a random direction in  $\mathbb{R}^d$ ?

**Lemma.** For  $u, v \in \mathbb{R}^d$ ,  $\mathbb{P}[h_r(u) = h_r(v)] = \frac{\theta(u, v)}{\pi}$

Proof: The projection of  $r$  onto the plane spanned by  $u, v$  is again a random direction.

- SimHash: a data structure that estimates angular similarity:
  - Take  $\ell$  random directions  $r^1, \dots, r^\ell \in \mathbb{R}^d$

# Sim Hash

- Can we estimate the similarity without computing the full inner product?
  - Say we have a data set  $v_1, \dots, v_n \in \mathbb{R}^d$
- Idea [Charikar]: use random projection!
  - Pick random direction (unit vector)  $r \in \mathbb{R}^d$
  - For each  $v_i$ , let  $h_r(v_i) := \text{sign}(r \cdot v_i)$

Recall: how do we sample a random direction in  $\mathbb{R}^d$ ?

**Lemma.** For  $u, v \in \mathbb{R}^d$ ,  $\mathbb{P}[h_r(u) = h_r(v)] = \frac{\theta(u, v)}{\pi}$

Proof: The projection of  $r$  onto the plane spanned by  $u, v$  is again a random direction.

- SimHash: a data structure that estimates angular similarity:
  - Take  $\ell$  random directions  $r^1, \dots, r^\ell \in \mathbb{R}^d$
  - For each record  $v_i$ , store the  $\ell$ -tuple  $(h_{r^1}(v_i), \dots, h_{r^\ell}(v_i))$

SimHash was showcased in this popular book by Jun Wu

