Similarity Estimation

Hu Fu @SHUFE, 2022.

Similarity Between Data Points

- Recall the Nearest Neighbour Search problem we discussed
- Similarity estimation arises in other applications as well
 - Finding near-duplicate documents/webpages and removing redundancy
 - Detecting pirated files
- How do we define and compute/estimate similarity?

Jaccard Similarity

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

Min Hashing

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

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

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

$i \in S$

$$S), \cdots, \sigma_{\min}^{\ell}(S))$$

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

 ℓ

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

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

<u>Def.</u> [Broder, Charikar, Frieze, Mitzenmacher] A family $\mathscr{F} \subseteq S_n$ is a minwise independent family of permutations if for every $S \subseteq [n]$ and any $a \in S$, for σ sampled uniformly from \mathscr{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 ($\ll |S_{\nu}|$

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

$$|n| = n!$$
)

Relaxation

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

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

Thm. [Indyk] Let \mathscr{H} be a *t*-wise independent hash family from [n] to [n], with $t = \Omega(\log \frac{1}{2})$, then \mathscr{H} is a (ϵ, 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

- Two vectors in \mathbb{R}^d may be considered similar if they point roughly to the same direction
- π and v
 - If u, v are unit vectors, then $\cos(\theta)$
 - Fact:
- Similarity between u, v is 1 dist(u, v)

• For $u, v \in \mathbb{R}^d$, let dist $(u, v) := \frac{\theta(u, v)}{\dots}$, where $\theta(u, v)$ is the angle between u

$$u(u,v)) = u \cdot v$$

Sim Hash

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

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

- SimHash: a data structure that estimates angular similarity:
 - Take ℓ random directions $r^1, \ldots, r^\ell \in \mathbb{R}^d$
 - For each record v_i , store the ℓ -tuple $(h_{r_1}(v_i), \dots, h_{r_{\ell}}(v_i))$

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

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



SimHash was showcased in this popular book by Jun Wu

