

NN-Search:

Given: set P of n points in \mathbb{R}^d (query point q from set Q in \mathbb{R}^d)

Goal: find the nearest-neighbor p of q in P

Background:

Tree Based Methods: (small dimensionality)

Stochastic Methods: (approximate) \rightarrow LSH

Filtering Based Methods: (exact)

Approximate Nearest Neighbors : input $\mathbb{P}, q, r, \epsilon \in \mathbb{R}$

output: if some points exist with $\text{dist}(p, q) \leq r$, output any such ANN (YES)

if no points within $\text{dist}(p, q) \leq (1 + \epsilon)r$, return NO

otherwise return any points within $\text{dist}(p, q) \leq (1 + \epsilon)r$

Local Sensitive Hashing: a set of functions is locally-sensitive if for a random hash function h in the set, for any pair of points p, q

$\Pr[h(p) = h(q)]$ is "high" if p close to q

$\Pr[h(p) = h(q)]$ is "low" if p far from q

preprocessing: hash data using several LSH functions so pr of collision for closer objects \uparrow

Querying: hash query point and retrieve elements in buckets containing query point

give $\#$ pull between $\#$ hashes $\&$ $\#$ tables

LSH: Hamming Distance

$$\sum_{i=1}^d 0, 1$$

$\text{ham}(p, q) = \#$ positions in which they differ

define hash function by choosing set S of k random dimensions

and setting $h(p) =$ projection of p on S

$$\Pr[h(p) = h(q)] = (1 - \text{ham}(p, q)/d)^k$$

ex: $d = 12, k = 3, p = \begin{array}{cccccccc} 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 1 \end{array}$

$S = \{1, 5, 6\} \rightarrow$ bucket $h(p) = 110$

L_p Norms: probability Distribution D is p -stable if

\rightarrow for any $v \in \mathbb{R}$ and random variables $X_1, \dots, X_n \sim D$, $\sum_i v_i X_i$ has same distribution as $[\sum_i v_i^p]^{1/p} X$ where $X \sim D$

Gaussian - normal distribution is 2-stable

Hash functions for L_p norm:

choose random $v_i \in [w]$ entries randomly drawn from p -stable distribution

compute $h_i(x) = \left\lfloor \frac{\langle v_i, x \rangle + b}{w} \right\rfloor$ where w is bucket width $\&$ b is uniformly random in $[0, w]$

LSH: Jaccard coefficient (consider binary vectors p, q)

$$\text{Jaccard}(p, q) = \frac{|P \cap Q|}{|P \cup Q|}$$

→ disjoint sets have $J(A, B) = 0 \neq J(A, A) = 1$

LSH: Cosine Similarity

→ choose random unit vector r

→ $h_r(p) = 1$ if $\langle r, p \rangle > 0$, 0 otherwise

→ build sketches/signatures by concatenating hash functions

IndexJoin (straight-forward exact solution)

$$\cos(d_1, d_2) = \frac{\langle d_1, d_2 \rangle}{\|d_1\|_2 \times \|d_2\|_2}$$

→ main idea: leverage data sparsity

→ $\text{sim}(d_1, d_2) = d_{1,1} \times d_{2,1} + d_{1,2} \times d_{2,2} \dots d_{1,n} \times d_{2,n}$ (n features)

method: normalize object vectors

construct inverted index from objects

for each obj:

→ compare only w/ objects containing features in common

→ select neighbors

Advantages:

→ skipping some comparisons

→ small memory footprint

Inverted Index:

$f_1 \ f_2 \ f_3 \dots$ features

documents
 d_1
 d_2
 d_3

=>

features
 f_1
 f_2
 f_3
:

Inverted index (so you know which obj to compare to each other)

$d_1 \ d_2 \ d_3 \dots$ documents

Accumulator(A): data structure tracking similarities while being computed

Sort the remaining results in Accumulator

select top k results and/or $\text{sim}(p, q) \geq \epsilon$

LZAP: (Fast Cosine Sim Search w/ Prefix L-2 Norm Bounds)

main idea: leverage similarity estimates

$$\langle d_a, d_c \rangle = \langle d_a^{\leq p}, d_c^{\leq p} \rangle + \langle d_a^{> p}, d_c^{> p} \rangle$$

	f_1	f_2	f_3	f_4	f_5	f_6
d_a	"	"	"	"	"	"
d_c	"	"	"	"	"	"

estimate (under f_1, f_2, f_3) *compute* (under f_4, f_5, f_6)

Filter objects not in final graph based on

similarity estimates

LZAP: Algorithm

1. Accumulate Similarity using inverted index

2. For each un-pruned obj, finish similarity computation using forward index

LZAP leads to greatly improved runtime

→ filters most objects without computing their similarity

→ outperforms all exact & most approximate baselines

key concepts:

→ more likely to be neighbors if we have common high weight feature

→ you are my neighbor's neighbor

Step 1: choose initial candidates

→ sort vectors & inverted lists in non-increasing weight order

→ traverse two lists at a time & gather $M \geq k$ candidates with higher dot product

$$C_M = S = [d_s, d_1, d_4]$$

→ compute similarities w/ candidates & keep k NN

Step 2:

→ Find potential better results (γ times)

→ visit neighbors in non-increasing similarity order

→ consider d_s , a neighbor of d_r , as a candidate if

• have collected less than M candidates

• $\text{sim}(d_s, d_r) \geq \text{sim}(d_r, d_q)$

→ compute similarities w/ candidates & update both neighborhoods of d_s & d_q

LZKnnng:

→ filtering Framework for k -NN construction

→ build initial approximate graph \hat{G}

• provides threshold for filtering

→ improve \hat{G} until exact

In practice:

higher # candidates = higher recall, but slower

Step 2 critical to high recall

higher # Step 2 iterations = higher recall