

Question :-

- We have a collection of N images $[1000 \times 1000]$ pixels.
- We are given a "query" image $[1000 \times 1000]$.
- Goal: quickly find out the "closest" image to the query image.

How?

High level idea

- Transform each image into a vector in $1000^2 = 10^6$ dimensions.

$$d \begin{array}{|c|c|c|} \hline 0 & 2 & 1 \\ \hline .15 & .7 & .8 \\ \hline 0 & .5 & .4 \\ \hline \end{array} \rightarrow \leftarrow D = d^2 \rightarrow (0, .2, .1, .15, .7, .8, 0, .5, .4)$$

- Same for query image

- Output the nearest vector to the query vector.
(say l_2 distance)

(is) Given N vectors $v_1, v_2, \dots, v_N \in \mathbb{R}^d$

1.3

and given query $q \in \mathbb{R}^d$

find

$$\arg \min_{i=1}^N \|v_i - q\|_2^2$$

(Approx.)

NEAREST NEIGHBOR SEARCH

Given N "base" vectors in D dimensions,
pre-process and build a data structure.

So that when queries are presented,
we can quickly retrieve the
(approximate) closest vector to the query

Naive Solution

when query is presented, compare
distances to all base points
and output the closest.

operations per query : $3 \cdot N \cdot D + N$
↙ ↓
D subtractions, find MIN.

Multiplications
Additions per point

In our example
 $D = 10^6$.

Today's lecture

generic "dimension reduction" techniques
which can bring down $10^6 \rightarrow 10^3$
with very little "error".

JOHNSON-LINDENSTRAUSS LEMMA

Given n vectors in \mathbb{R}^d , v_1, v_2, \dots, v_n
there is a mapping $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$
such that

$$\forall u_i, u_j$$
$$(1-\epsilon) \leq \frac{\|f(u_i) - f(u_j)\|_2}{\|u_i - u_j\|_2} \leq (1+\epsilon)$$

with $k = \Theta\left(\frac{\log n}{\epsilon^2}\right)$.

Moreover we can find this efficiently
with good probability.

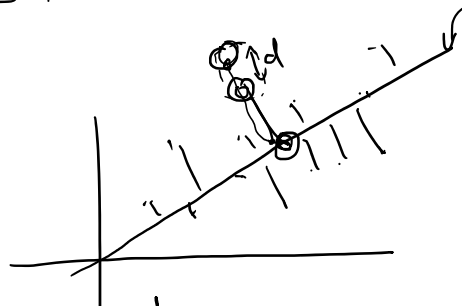
Goal for today

Idea: PCA?

- Unclear

- It's more of a

- global error reduction



- we want pairwise preservation

Idea 2% what about just picking
 k random coordinates?

$v = (x_1, x_2, \dots, x_d)$ then $f(v) = (x_{i_1}, x_{i_2}, \dots, x_{i_k})$
where i_1, i_2, \dots, i_k are
randomly chosen from $[d]$.

Same i_1, i_2, \dots, i_k are used for all vectors.

[sample with replacement, to make
analysis easy]

Good News

for any u and $v \in \mathbb{R}^d$

$$E [\|f(u) - f(v)\|^2] = \frac{k}{d} \|u - v\|^2$$

↓
good, b/c then we can think of

$$\hat{f}(u) = \sqrt{\frac{d}{k}} (u_{i_1}, u_{i_2}, \dots, u_{i_k})$$

and then we'll have

$$E [\| \hat{f}(u) - \hat{f}(v) \|^2] = \|u - v\|^2$$

Proof of Good News

Fix u, v and let $z = u - v$

$$\|z\|^2 = \sum_{i=1}^d (u_i - v_i)^2 = \sum_{i=1}^d z_i^2$$

Now, consider i_1 and focus on

$$\begin{aligned} E_{i_1} [(u_{i_1} - v_{i_1})^2] &= \sum_{l=1}^d \underbrace{\text{Pr}(i_1 = l)}_d \cdot (u_l - v_l)^2 \\ &= \frac{1}{d} \cdot \sum_{l=1}^d (u_l - v_l)^2 \\ &= \frac{1}{d} \|u - v\|_2^2 \end{aligned}$$

$$\begin{aligned} \Rightarrow E [\|f(u) - f(v)\|_2^2] &= \sum_{l=1}^k E [(u_{i_l} - v_{i_l})^2] \\ &= \frac{k}{d} \|u - v\|_2^2 \end{aligned}$$

If we can show some "concentration inequality" that $\|f(u) - f(v)\|$ is close to its average whp, then this scheme works!

and then union bound for all $1 \leq i < j \leq n$.

Q: Is this likely to have concentration?

A: Sadly, no 😞

Here's an example :-

$$v_1 = (\overbrace{v, 0, 0, \dots, 0}^{\text{---}})$$

$$v_2 = (0, \overbrace{v, 0, \dots, 0}^{\text{---}})$$

$$\vdots$$
$$v_n = (0, 0, \dots, v)$$

Unless we pick one of 1 or 2
coordinate in f ,

$$\|f(u_1) - f(u_2)\| = 0$$

\Rightarrow even if $k < n-1$, we will
always make huge
error on some pair.

$$\left. \begin{array}{l} \text{real distance} = \sqrt{2}v \\ \text{"our imagined" distance} = 0 \end{array} \right\}$$

badness because most coordinates are
same, few coords give all
the distance -

How to "spread" this over all coords?

TRICK:

- Nothing sacred about standard basis.
 Choose a random rotation $R \in \mathbb{R}^{d \times d}$

first rotate the data

$$V_i \rightarrow R V_i \quad (\text{still in } d\text{-dim})$$

then sample k coords from $R V_i$

- intuition, can be made rigorous

- can also be simplified greatly using Gaussians.

$$f(u) = \frac{1}{\sqrt{\lambda}} \begin{bmatrix} g_{11} & g_{12} & \dots & g_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ g_{k1} & g_{k2} & \dots & g_{kd} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_d \end{bmatrix}$$

\uparrow
G

09/04/2021

$$f(u) = \frac{1}{\sqrt{\lambda}} \begin{bmatrix} g_{11} & g_{12} & \dots & g_{1d} \\ g_{21} & \dots & & g_{2d} \\ \vdots & & & \vdots \\ g_{k1} & \dots & & g_{kd} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_d \end{bmatrix}$$

$\leftarrow d \rightarrow$

$G =$ random $k \times d$ gaussian matrix
 where each g_{ij} is an
 independent gaussian r.v
 with mean 0 & variance 1.

(1e) $g_{ij} \sim N(0, 1)$.
 λ will be chosen suitably later.

THEOREM

For suitable λ and k , this mapping
 preserves all pairs of distances
 upto $(1 \pm \epsilon)$ factor.

Q How do we prove this?

We'll fix a pair of vectors u_i, u_j
 and let $z = u_i - u_j$

Next we show that

$$\text{Pr}_{\text{random choice of } G} \left[\frac{\|f(u_i) - f(u_j)\|_2^2}{\|u_i - u_j\|_2^2} \notin (1 - \epsilon, 1 + \epsilon) \right]$$

$$\leq \frac{1}{n^2}$$

Then we just do union bound over all
 pairs u_i, u_j .

$$\Pr_{G, \lambda} \left[\exists i, j \text{ which is not preserved} \right] \leq \frac{\binom{n}{2}}{n^2} \leq \frac{1}{2}$$

With prob $\frac{1}{2}$, experiment was successful
 (ie) all pairs preserved!
 If not, just repeat 😊

Need to show that a single pair is preserved w.h.p.

Fix u & v and let $Z = u - v$

$$f(u) = \lambda \cdot G \cdot u$$

$$f(v) = \lambda \cdot G \cdot v$$

$$f(u) - f(v) = \lambda \cdot G \cdot (u - v) = \lambda G Z = f(Z)$$

Need to show that

$$1 - \epsilon \leq \frac{\|f(Z)\|_2^2}{\|Z\|_2^2} \leq 1 + \epsilon \quad \text{with prob} \geq 1 - \frac{1}{n^2}$$

let $Z = (z_1, z_2, \dots, z_d)$,

and consider $\hat{Z} = \frac{Z}{\|Z\|}$

Then $Z = \|Z\| \cdot \hat{Z}$

Moreover,

$$f(z) = \lambda \cdot G \cdot z = \lambda (\|z\| \cdot G \hat{z}) \\ = \|z\| \cdot f(\hat{z})$$

{ So, it suffices to show that } $w.p. \geq 1 - \frac{\epsilon}{n^2}$
 $(1-\epsilon) \leq \|f(\hat{z})\|_2^2 \leq (1+\epsilon)$ (*)

$$\Rightarrow \|z\|^2 (1-\epsilon) \leq \|f(z)\|_2^2 \leq (1+\epsilon) \|z\|^2$$

Consider $\hat{z} = (\hat{z}_1, \hat{z}_2, \dots, \hat{z}_d)$

We know $\|\hat{z}\|_2^2 = 1 \Rightarrow \sum \hat{z}_i^2 = 1$.

Let's look at $f(\hat{z}) = \lambda \cdot G \cdot \hat{z}$

$$= \lambda \begin{bmatrix} \sum_{j=1}^d g_{1j} \hat{z}_j \\ \sum_{j=1}^d g_{2j} \hat{z}_j \\ \vdots \\ \sum_{j=1}^d g_{dj} \hat{z}_j \end{bmatrix}$$

Each entry in $f(\hat{z})$ looks like

$$\lambda \cdot N\left(0, \sum_{j=1}^d \hat{z}_j^2\right) = \lambda N(0, 1)$$

Moreover because g_{ij} and $g_{i'j'}$ are independent for all i, j, i', j'

independent for all i, j, i', j' ,

these entries are themselves independent!

$$f(\hat{z}) = \lambda \begin{pmatrix} F_1 \\ F_2 \\ \vdots \\ F_k \end{pmatrix} \quad \text{where } F_i \sim N(0,1).$$

$$\text{What is } E[\|f(\hat{z})\|^2] = \lambda \sum_{i=1}^k E[F_i]^2 = \lambda \cdot k.$$

We want $f(\hat{z})$ to be close to 1.

$$\Rightarrow \text{set } \boxed{\lambda = \frac{1}{k}}$$

$$\Rightarrow E[\|f(\hat{z})\|^2] = 1.$$

Let $\gamma = k \cdot \|f(\hat{z})\|^2$ be the random variable.

$$\gamma = \underbrace{\sum_{i=1}^k F_i^2}$$

γ is the sum of many independent-random variables.

Moreover each F_i is a gaussian $N(0,1)$

χ^2 - distribution

1.1 - concentration

$\left\{ \begin{array}{l} Y \text{ actually sharply concentrates around its} \\ \text{mean,} \end{array} \right\}$

Try to prove this using ideas from Chernoff bounds proof

More or less,

Y behaves like $N(\mu(Y), \sigma^2(Y))$

$$E[Y] = \sum E[F_i^2] = k.$$

$$\text{Var}[Y] = \sum_{i=1}^k \text{Var}[F_i^2]$$

$$\text{Var}[F_i^2] = E[F_i^4] - E[F_i^2]^2$$

$$= E[F_i^4] - 1.$$

$$\left\{ \begin{array}{l} E[F_i^4] = O(1) \text{ [constant]} \\ \text{for gaussian distrib.} \end{array} \right.$$

$$\Rightarrow \text{Var}[Y] = Ck \text{ for constant } C.$$

$\left\{ \begin{array}{l} Y \text{ behaves like } N(k, Ck) \end{array} \right\}$

Very big cheat, crude approx.

Basically, we'll use Chernoff-like bounds on Y .

tail bounds on Y .

$$\Pr \left[|Y - E[Y]| \geq t \right] \leq \exp\left(\frac{-t^2}{\sigma^2}\right)$$

Can prove this type of inequality for Y .

$$E[Y] = k$$

$$\sigma^2(Y) = c \cdot k$$

$$\text{Set } t = \varepsilon \cdot k = \varepsilon E[Y]$$

$$\Pr \left[Y \notin (1-\varepsilon, 1+\varepsilon) E[Y] \right] \leq \exp\left(\frac{-\varepsilon^2 k^2}{c \cdot k}\right)$$

$$= \exp\left(-\frac{\varepsilon^2 k}{c}\right)$$

So just set $k = \frac{2c \cdot \ln n}{\varepsilon^2}$

$$\Pr(\text{bad event}) \leq \exp\left(\frac{-\varepsilon^2 \cdot 2c \ln n}{c \varepsilon^2}\right)$$

$$= \frac{1}{n^2} \quad \square$$

Whenever Y lies $\in [1-\varepsilon, 1+\varepsilon] E[Y]$

$$\Leftrightarrow f(\hat{z}) = \frac{1}{k} Y \text{ lies in } [1-\varepsilon, 1+\varepsilon]$$

$$\Leftrightarrow f(\tilde{z}) = \frac{1}{k} \gamma \text{ lies in } [1-\epsilon, 1+\epsilon].$$

$$\Rightarrow \Pr \left(f(\tilde{z}) \notin [1-\epsilon, 1+\epsilon] \right) \leq \frac{1}{n^2} \text{ for}$$
$$k = \frac{2c \ln n}{\epsilon^2}$$

and $\lambda = \frac{1}{k}$.

(C, r) - Approximate Near Neighbor search.

Motivation

- Approx vectors in high dimensional space
- look to retrieve "closest" vector to given query vector.

Concrete Example :-

N documents, N is very large
(all in english)

represented as a $\{0,1\}^D$ vector for
 D very large.

$D =$ all distinct english words Ordered as $\{w_1, w_2, \dots, w_D\}$

$v = (v_1, v_2, \dots, v_D)$

where $v_i = 0$ if i^{th} word is not in document
 $= 1$ if i^{th} word is in doc.

Goal: build a data structure / Algorithm

st when query vector q arrives, we quickly retrieve any nearby base vector.

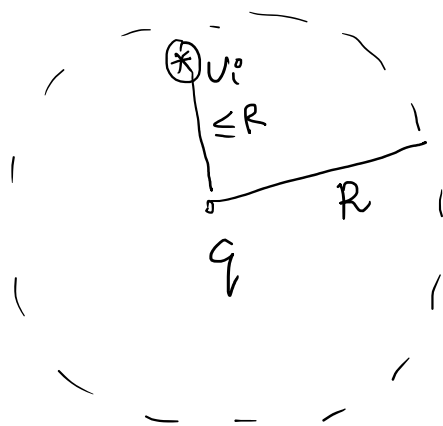
$$q = (q_1, q_2, \dots, q_D)$$

Given N base vectors u_1, u_2, \dots, u_N and $q = (q_1, q_2, \dots, q_D)$ arrives, and given parameter R , find any base vector which are at distance $\leq R$ from query q , if they exist.

Distance = Hamming distance

$$= \|u - v\|_1 = \sum_{j=1}^D |u_j - v_j|$$

= # indices where they differ



Naive approach

- Compare distances to all base vectors
Output the closest.

Query Time = $O(ND)$

Space of data structure = ND bits
(store all base vectors).

QN:

Can we get sub-linear dependence on N ?
eg, even if we bring query time
from ND to $\sqrt{N} \cdot D$, that is
really significant in real life
 N can be 10^9 , 10^{12} even!!!

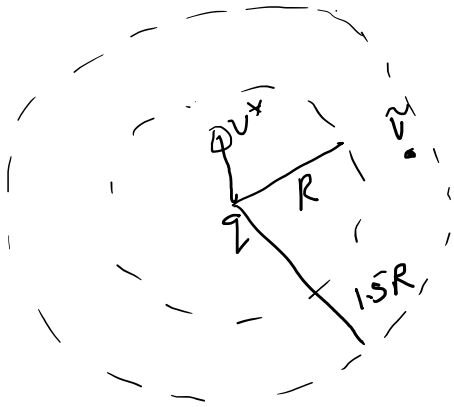
Today + Tomorrow:

Can do such a thing if we're OK
with Approximate Near Neighbor Search

Given q , target R ,

- Algo can output No if no base
point is distance $\leq R$
from q .

- Output any base point at distance $\leq CR$ from q - for some $C > 1$.



[Algo can output \tilde{v} instead of v^*]

In this case $C = 1.5$

C is fixed up-front, is like "approximation factor of algo".

THM [Indyk - Motwani '90s]

For above problem, can design C -Approx NNS algorithm with space $O(N^{1+1/c} + ND)$ but query time $= O(N^{1/c} D)$.

Idea: Design very generic technique called "LOCALITY SENSITIVE HASHING".

Imagine the following type of search Algo:-
Construct a hash function

Compute a hash function

$$h: \{0,1\}^D \rightarrow [k] \quad \text{eg. where } k \ll N$$

think of $k = \sqrt{N}$.

It maps each document vector to one of k buckets.

Similarly it can map query to a bucket.

Algo: ONLY look in the query bucket and output the closest-base vector.

Query time: $\approx \frac{N}{k} \cdot D \approx \sqrt{N} \cdot D$ if $k = \sqrt{N}$.

h is a random hash \Rightarrow each bucket has $\approx N/k$ pts.

Issue

No reason why query's nearest neighbor is hashed to same bucket as q .

Idea:

What if hash fn h is locality sensitive?
(i.e.) nearby points are in same bucket more often than far away points!

than far away points!
Then maybe such an algo would work!

LSH

A hash family $H = \{h_1, \dots, h_m\}$ is

a (C, R, P_1, P_2) LSH for $P_1 > P_2$ if

for any 2 vectors x and y ,

$$\textcircled{1} \quad \Pr_{h \in H} [h(x) = h(y)] \geq P_1 \quad \text{if } d(x, y) \leq R$$

$$\textcircled{2} \quad \Pr_{h \in H} [h(x) = h(y)] \leq P_2 \quad \text{if } d(x, y) > CR$$

Property $\textcircled{2}$ ensures that buckets are small on avg, and

$\textcircled{1}$ ensures that near neighbors of query are in same bucket on q.

{ From (C, R, P_1, P_2) LSH to (C, R) Algo. }

Idea :

"boost" the gap between P_1 and P_2

$\textcircled{1}$ h. taking

① by taking

$$g(x) = [h_1(x), h_2(x), \dots, h_k(x)]$$

where h_1, h_2, \dots, h_k are k independent samples from \mathcal{H} .

What does g satisfy:

$$\Pr [g(x) = g(y)] \geq p_1^k \quad \text{if } d(x, y) \leq R$$

$$\Pr [g(x) = g(y)] \leq p_2^k \quad \text{if } d(x, y) > CR$$

- This has driven down p_2^k \Downarrow
- Now, this also has driven down good case collisions (p_1^k)
- Fix this by using L different hash functions

$$g_1, g_2, \dots, g_L \quad \left\{ \begin{array}{l} \text{independently} \\ \text{for each } g_i \end{array} \right.$$

Overall Alg :-

- ① Compute $g_i = (h_{i1}, h_{i2}, \dots, h_{ik})$
for $i = 1, 2, \dots, L$.
- ② \dots

② Compute $g_i(v)$ for all $1 \leq i \leq L$ and all base vectors.

When query arrives

Compute $g_1(q), \dots, g_L(q)$.

For $i = 1, 2, \dots, L$

- Look at all base vectors with $g_i(v) = g_i(q)$.

- If we find any at distance $\leq cr$, output and terminate

Need to analyze time and space complexity.

Space complexity: $N \cdot D + O(N \cdot k \cdot L)$
↑ hash values.

Query Time :-

$Lk + L * \left(E[\# \text{ far away pts which fall into bucket of } q] \cdot D \right)$

↑ # hash fns g ↑ dimension

Sort of "wasteful computation"

sort of
"wasteful computation"



$$\leq L \cdot N \cdot P_2^k \cdot D$$

$$E[\text{query time}] \leq L N P_2^k \cdot D$$

Now its all about parameter choosing

$$\text{Choose } k \text{ st } P_2^k = \frac{1}{N}$$

(1e)

$$k \log \frac{1}{P_2} = \log N$$

$$k = \frac{\log N}{\log \frac{1}{P_2}}$$

Next choose L so that q 's near-neighbor
is in the bucket with good
probability

Let v^* be q 's neighbor within
distance $\leq r$ (if any).

What is the probability that
it falls into one of the
 L buckets of q ?

$P_0 [v^* \text{ falls in } q^S \text{ bucket in one of the } L \text{ hash fns}]$

$$\approx 1 - \prod_{i=1}^L (1 - p_i^k)$$

$$= 1 - (1 - p_i^k)^L$$

Need to formalize $\rightarrow \approx$

$$L \cdot p_i^k$$

$$(1 - \epsilon)^L \approx 1 - \epsilon L \text{ when } \epsilon \text{ is small}$$

Set L st

$$L \cdot p_i^k = 1$$

to ensure good success probability

$$L = \frac{1}{p_i^k} = \left(\frac{1}{p_i}\right)^{\frac{\log N}{\log(1/2)}}$$

$$L = N^{\frac{\log(1/p_i)}{\log(1/2)}} \leftarrow p$$

Parameter $p = \frac{\log(1/p_i)}{\log(1/2)} < 1$ is

$\log(1/p_L)$

the important parameter of LSH.

SUMMARY

$$k = \frac{\log N}{\log 1/p_L}; \quad L = N^p$$

Expected Query Time

$$= kL + O(\underbrace{N \cdot L \cdot p_L^k}_{\text{great if } p < 1} \cdot D)$$

$$= \frac{\log N}{\log 1/p_L} \cdot N^p + O(N^p \cdot D)$$

{great if $p < 1$ } 😊

Success Prob

if \exists a near nbr, we find some approx NN with Prob

$$\geq \left(1 - (1 - p_i^k)^L\right)$$

$- p_i^k \geq L$

$$\begin{aligned} (e^x \geq 1+x) &\Rightarrow 1 - e^{-x} \\ &\geq (1 - \gamma) \end{aligned} \quad \square$$

$$(e^x \geq 1+x) \quad \geq (1-\epsilon)$$

THM :

(C, R, P_1, P_2) LSH with $p = \frac{\log \frac{1}{P_1}}{\log \frac{1}{P_2}}$

implies a randomized

(C, R) -ANNS algorithm

with success prob $\geq (1-\epsilon)$ of finding a near nbr @

distance $\leq CR$ (if there exists

some nbr at dist $\leq R$)

and expected query processing time

$$= O\left(N^p \left(\frac{\log N}{\log \frac{1}{P_2}} + D\right)\right)$$

\leftarrow constant, can ignore

\Rightarrow LSH type hash fns are good as long as p is small !!

Qn

for vectors in $\{0,1\}^D$ with $d(x,y)$

$$= \sum_{i=1}^D |x_i - y_i|_1 \text{ as distance}$$

What is a good LSH?

- Want
- If x, y are close by ($\leq R$)
 $\Pr[h(x) = h(y)]$ is large.
 - If x, y are far apart ($\geq CR$)
 $\Pr[h(x) = h(y)]$ is small.

Dirt Simple Idea

Pick a random coordinate
 d from $\{1, 2, \dots, D\}$.

$$\boxed{h_d(x) = x_d}$$

$$H = \{h_1, h_2, \dots, h_D\}$$

P_1, P_2 ANALYSIS

if $d(x, y) \leq R$, (i.e.) they differ
 in at most
 R places,

$$\Pr_d [h_d(x) = h_d(y)] \geq \underbrace{\left(1 - \frac{R}{D}\right)}_{P_1}$$

$$\Pr_d [h_d(x) = h_d(y)] \leq \underbrace{\left(1 - \frac{CR}{D}\right)}_{P_2}$$

if $d(x, y) \geq cR$

What is the f of this hash fun.

$$f = \frac{\log\left(\frac{1}{1-R/D}\right)}{\log\left(\frac{1}{1-R/D}\right)}$$

$$= \frac{\log\left(\frac{1}{1-\frac{R}{D}}\right)}{\log\left(\frac{1}{1-\frac{R}{D}}\right)}$$

Need to be formalized

Assuming $R \ll D$,
we'll make some simplifications

a) $\frac{1}{1-\epsilon} \approx 1+\epsilon$

b) $\log(1+\epsilon) \approx \epsilon$

} if $\epsilon \ll 1$.

Hence $f \approx \frac{R/D}{D \cdot cR} = \frac{1}{c}$.

Using this hash family gives us
the (C, R) - Near Neighbor search
problem with
space $O(N^{1+1/c} + ND)$ and
query time $O(N^{1/c} \cdot (D + \log N))$.

⇓
· Represents a significant milestone
for "nearest vector search" and
is used in variety of applications
today !!