# Approximate Nearest Neighbor Search (ANN)

store $n$ pts $P \subseteq \mathbb{R}^d$  ($d$ const)
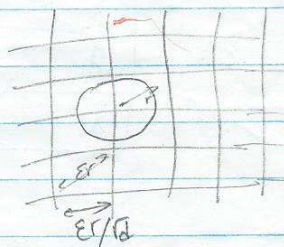
st. given query pt $q \in \mathbb{R}^d$, find $p \in P$ with

$$d(p,q) \leq (1+\varepsilon) \min_{p' \in P} d(p',q)$$

Exact problem ($\varepsilon = 0$): by halfspace range searching in $d+1$ dims

$O(n)$ space, $O^*(n^{1 - 1/\lceil d/2 \rceil})$ time

or $O^*(n^{\lceil d/2 \rceil})$ space, $O(\log n)$ time

Approx decision problem: given $r$,

return some pt of distance $\leq (1+\varepsilon) r$

or declare all pts have distance $> r$.

Method 0: for decision with fixed $r$

form grid of side length $\varepsilon r / \sqrt{d}$

store $S = $ all nonempty grid cells



query($q$):

check if any grid cell intersecting $ball(q, r)$

is in $S$ ← by hashing

space $O(n)$

time $O(\#$grid cells intersecting $ball(q,r))$

$$= O\left( \frac{volume(ball(q, (1+\varepsilon)r))}{(\varepsilon r / \sqrt{d})^d} \right)$$

naively: $O\left( \frac{r^d}{(\varepsilon r / \sqrt{d})^d} \right)$

$$= O\left( \frac{\frac{2\pi^{d/2}}{d\Gamma(d/2)} ((1+\varepsilon)r)^d}{(\varepsilon r / \sqrt{d})^d} \right)$$

$\equiv O\left( \left( \frac{\sqrt{d}}{\varepsilon} \right)^d \right)$

$$= O\left( \left( \frac{c}{\varepsilon} \right)^d \right)$$

Rmk - alternatively, $O\left( \left( \frac{c}{\varepsilon} \right)^d n \right)$ space, $O(1)$ time
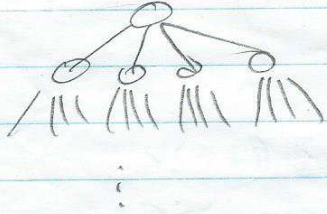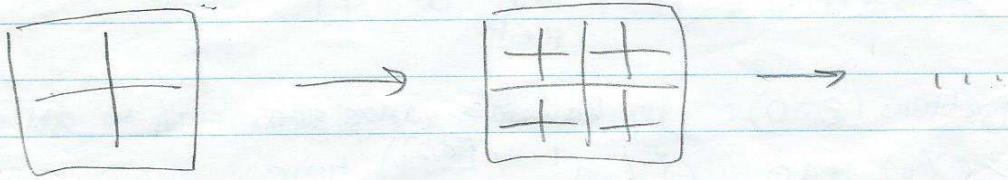
what if $r$ is not fixed? (store $S = $ all grid cells intersecting $ball(p, r)$ for some $p \in P$)
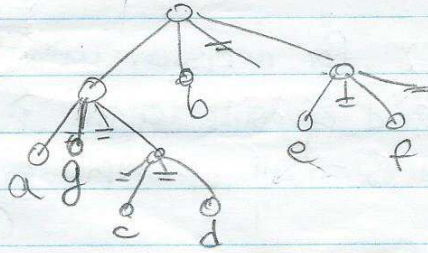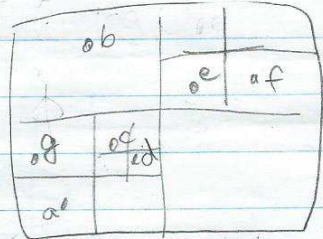
/A

# Method 1 : Quadtree
### idea - hierarchy of grids



**Def** A quadtree cell B is
a grid cell of side length $2^\ell$
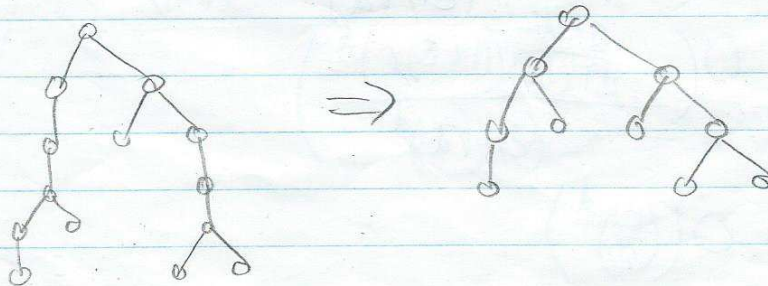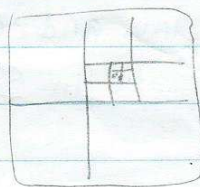for some $\ell \in \{0, 1, .., \log U\}$

e.g.



height
$O(\log U)$

space $O(n \log U)$

reducible to $\boxed{O(n)}$
by compressed quadtree
(shortcutting "deg-1 nodes)

decis-query $(B, q, r)$:
    if ball$(q,r)$ does not intersect $B$ return
    if $B$ has side length $< \varepsilon r/\sqrt{d}$ return any pt in $B$
    for each child $B_i$ of $B$
       decis-query $(B_i, q, r)$

$\Rightarrow$ time $= O\Big(\ \#$ quadtree cells of side length $> \frac{\varepsilon r}{\sqrt{d}}$
                        (intersecting $B$) $\Big)$

(ignore const factors depending on $d$)

$$= O\left( \sum_{\ell:\, 2^\ell > \frac{\varepsilon r}{\sqrt{d}}} \left\lceil \frac{r^d}{(2^\ell)^d} \right\rceil \right)$$

$$= O\left( \left\lceil \frac{1}{(\varepsilon/\sqrt{d})^d} \right\rceil + \left\lceil \frac{1}{(2\varepsilon/\sqrt{d})^d} \right\rceil + \cdots \right)$$

$$= \boxed{O\left( \log U + \frac{1}{\varepsilon^d} \right)}$$
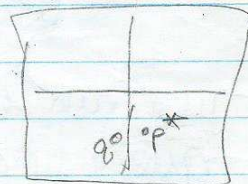
(how to adapt to ANN?

[maintain curr. min $r$, expand cells in increas order of dist]

Method 2    ... with Shifting   (Bern '93) $\leftarrow$ first aim for const factor first approx.

   naive-ANN-query $(B, q)$:
     find child $B_i$ containing $q$
     return naive-ANN-query $(B_i, q)$

Let $p^* =$ nearest neighbor of $q$
    $r^* = d(p^*, q)$
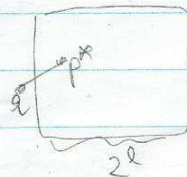
Def $q$ is good if $p^*$ and $q$ lie in a quadtree cell
              of side length $< 2(d+1)r^*$

Shifting Lemma   Shift all pts by random vector in $[U]^d$
      Then $q$ is good with const prob.
Pf: Say $(d+1)r^* < 2^\ell \leq 2(d+1)r^*$
     $\Pr[q \text{ bad}] = \Pr[\, p^* q \text{ crosses boundary}$
                of quadtree cell of
                   side length $2^\ell\,]$

$$\leq \Pr[\text{shift lies in a bad interval of length } r^* \text{ along some dim}]$$

$$\leq d \cdot \frac{r^*}{2^\ell} \leq \frac{d}{d+1} \quad \square$$

## Shifting Lemma (Deterministic version)

Say $d$ even, $U = 2^w$.

Shift all pts by vector $v_i = \left( \frac{i 2^w}{d+1}, \ldots, \frac{i 2^w}{d+1} \right)$, $i = 0, \ldots, d$.

Then $q$ is good for some $i$.

Pf: $q$ bad for $i \Rightarrow \frac{i 2^w}{d+1} \bmod 2^\ell$ lies in a bad interval of length $r^*$
   along one dim.

$\underset{\substack{\text{multiply by} \\ \frac{d+1}{2^\ell}}}{\Longrightarrow}$ $i 2^{w-\ell} \bmod (d+1)$ lies in a bad interval of length $\frac{(d+1) r^*}{2^\ell} < 1$

Fact: $ax \equiv b \bmod n$ has unique soln if $a, n$ are relatively prime

$\Longrightarrow$ $i 2^{w-\ell} \bmod (d+1)$ is equal to a specific bad integer.

$\underset{\substack{\text{since } 2^{w-\ell} \text{ and } d+1 \\ \text{are rel. prime}}}{\Longrightarrow}$ 'at most' one bad $i$ per dim
   so by pigeonhole, $\exists$ good $i$. $\square$

Just run naive-ANN-query for each of these $d+1$ shifts.

query time: $\boxed{O(\log U)}$

Analysis:

   let $p$ be returned pt at leaf.



Suppose $q$ is good,
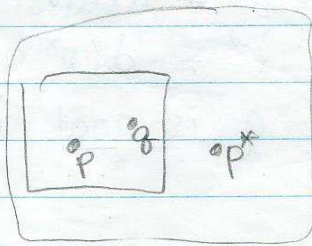   $p^*$ and $q$ are in quadtree cell of side length $< 2(d+1) r^*$

$\Rightarrow$ so are $p$ and $q$
$\Rightarrow$ $d(p,q) < 2(d+1)\sqrt{d} \, r^*$
$\Rightarrow$ approx factor $\boxed{O(1)}$.

reduce query time to $O(\log n)$? [Arya, Mount et al. '94]

## Method 3: Balanced Quadtree

tree centroid ←

**Lemma** $\exists$ quadtree cell $B$ s.t. $|P \cap B|, |P-B| \leq \dfrac{2^d}{2^d+1} n$

Pf: let $B$ be smallest quadtree cell with $|P \cap B| \geq \dfrac{1}{2^d+1} n$

Then $|P-B| \leq \dfrac{2^d}{2^d+1} n$

and for each child $B_i$ of $B$,

$$|P \cap B_i| \leq \frac{1}{2^d+1} n \implies |P \cap B| \leq \frac{2^d}{2^d+1} n \quad \square$$

recurse in $P \cap B$ & $P-B$

$\implies$ binary tree of height $\log_{\frac{2^d+1}{2^d}} n = O(\log n)$

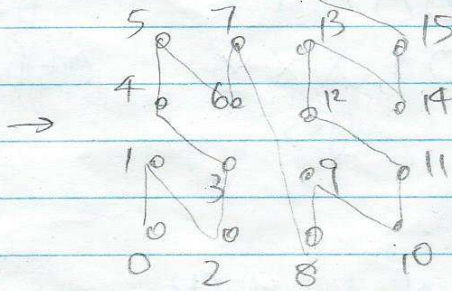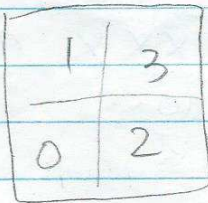$\implies$ query time $O(\log n)$
approx factor $O(1)$    by shifting

(other variants: Arya, Mount et al.'s BBD trees, Goodrich et al's BAR trees, ...)

## Method 4: Z-Ordering   (C.'02)

idea — don't use any trees!
directly reduce to 1D, by sorting?



space-filling curve

(other exs: Hilbert curve etc.)

Def  Given point $p = (x, y) \in [U]^2$,
   write  $x = a_{w-1} a_{w-2} \cdots a_0$  in binary
        $y = b_{w-1} b_{w-2} \cdots b_0$
  define shuffle$(p) = a_{w-1} b_{w-1} a_{w-2} b_{w-2} \cdots a_0 b_0$

e.g.  $x = 2, y = 1 \rightarrow \begin{array}{l} 010 \\ 001 \end{array} \xrightarrow{\text{shuffle}} 001001 \rightarrow 9$

Notes
- shuffle$(p)$ corresponds to position in Z-order
- shuffle maps a quadtree cell to an interval of contiguous integers in 1D.

Lemma  can compare shuffle$(p)$ with shuffle$(q)$ in $O(1)$ time.
Pf:  $p = (x, y), \quad q = (x', y')$

   if $msb(x \oplus x') > msb(y \oplus y')$
     compare $x$ with $x'$
   else compare $y$ with $y'$

e.g. $\begin{array}{ll} x = 1101 & y = 1011 \\ x' = 1011 & y' = 1001 \\ x \oplus x' = 0110 & y \oplus y' = 00100 \end{array}$

$\oplus$ bitwise or
$msb$ most significant 1 bit

to check $msb(a) > msb(b)$:
  check $a > b$ and $a \oplus b > b$.

e.g. $\begin{array}{l} a = 01XXX \\ b = 00001X \\ a \oplus b = 01XXX \end{array}$    $\begin{array}{l} a = 00001X \\ b = 01XXX \\ a \oplus b = 01XXX \end{array}$    $\begin{array}{l} a = 01XXX \\ b = 01XXX \\ a \oplus b = 00XXX \end{array}$

       yes               no           no

Preproc alg'm:
  for each of $d+1$ shifts       $\Rightarrow$ $O(n \log n)$ time
    store pts sorted along Z-order       $O(n)$ space

query(q):
    for each of $d+1$ shifts
      find pred & succ of q     $\Rightarrow$ $\boxed{O(\log n)\text{ time}}$
    return closest pt found
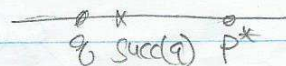
Analysis:
    Suppose q is good.
    $p^*$, q are in quadtree
      cell of side length $< 4(d+1) r^*$
  $\Rightarrow$ so is succ or pred of q.
  $\Rightarrow$ approx factor $O(1)$.

Z-order

Rmk: immediately improvable to $O(\log\log U)$ or $O(\sqrt{\log n})$
                by 1D data structures

    can be made dynamic!

Refining approx factor:

Option 1:   let r = factor-c approx.
    divide ball(q, cr) into $O\left(\frac{1}{\varepsilon^d}\right)$ quadtree cells
               of side length $\leq \frac{\varepsilon r}{\sqrt{d}}$
    test if each subcell is empty
           maps to 1D search.
  $\Rightarrow$ query time $\boxed{O\left(\frac{1}{\varepsilon^d}\log n\right)}$ or $O\left(\frac{1}{\varepsilon^d}\log\log U\right)$ ...

Best $\varepsilon$-dependency? still open
Arya, Mount et al. '94    $O(n)$ space, $O\left(\frac{1}{\varepsilon^d}\log n\right)$ time
Clarkson '94 / C '97    $\tilde{O}\left(\frac{1}{\varepsilon^{d/2}} n\right)$ space,  $O\left(\frac{1}{\varepsilon^{d/2}}\log n\right)$ time
Arya, da Fonseca, Mount '12-'16   $\tilde{O}\left(\frac{1}{\varepsilon^{d/3}} n\right)$ space,  $\dot{O}\left(\frac{1}{\varepsilon^{d/3}}\log n\right)$ time
  "    "    "   '17 / C '17   $\tilde{O}\left(\frac{1}{\varepsilon^{d/4}} n\right)$  "   $O\left(\frac{1}{\varepsilon^{d/4}}\log n\right)$ "

A7

Option 2: (C.-Har-Peled-Jones '18)

idea - keep query alg'm same
but use multiple orderings

increase fan-out to $\left(\frac{1}{\epsilon}\right)^d$ (i.e. compress $\log \frac{1}{\epsilon}$ levels into one)

reorder children

$2^c$

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

| 0 | 2 | 1 | 3 |
|---|---|---|---|
| 4 | 6 | 5 | 7 |
| 8 | 10 | 9 | 11 |
| 12 | 14 | 13 | 15 |

| 8 | 2 | 4 | 6 |
|---|---|---|---|
| 5 | 7 | 1 | 3 |
| 8 | 10 | 12 | 14 |
| 13 | 15 | 9 | 11 |

...

$\forall$ children $u_i, u_j$,

$\exists$ ordering s.t. $u_i, u_j$ adjacent

# orderings $O\left(\left(\frac{1}{\epsilon}\right)^d \log \frac{1}{\epsilon}\right)$

$\underbrace{\quad\quad}_{\substack{\text{\# levels} \\ \text{compressed}}}$

[ many appl'ns ... ]