Approx
# Nearest Neighbor Search in High Dims

**Problem**   Given $n$ pts $P$ in $\mathbb{R}^d$,
build data structure s.t.
given query pt $q \in \mathbb{R}^d$, can find $p \in P$
minimizing $d(p,q)$.

Euclidean: $\sqrt{\sum\limits_{i=1}^{d} (p_i - q_i)^2}$
$(L_2)$

$L_1$: $\sum\limits_{i=1}^{d} |p_i - q_i|$

$L_\infty$: $\max\limits_{i=1}^{d} |p_i - q_i|$

**Known:**   $d=2$:   $O(n)$ space, $O(\log n)$ time by PL...
& Voronoi diag.

larger const $d$:

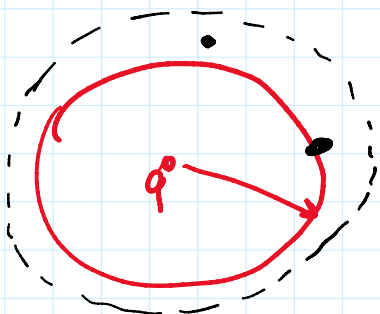$L_1 / L_\infty$:   $O(n \log^d n)$ space, $O(\log^d n)$ query time

$L_2$:   $O(n^{\lceil d/2 \rceil})$ space, $O(\log n)$ time
or $O(n)$ space, $\tilde{O}(n^{1 - \frac{1}{\lceil d/2 \rceil}})$ time

terrible for large $d$!

Relax problem: given $q \in \mathbb{R}^d$, find $p \in P$ s.t.

$$d(p,q) \leq c \cdot \min\limits_{p' \in P} d(p',q)$$

approximation factor

$q$

Suffice to solve approx <u>decision</u> problem, for <u>fixed radius r</u>:

given $q \in \mathbb{R}^d$,

return some pt $p \in P$ with $d(p,q) \leq cr$

or declare all pts $p \in P$ with $d(p,q) > r$.

(original prob can be solved by approx binary search

extra factor $\log U$)

try all $r = (1+\varepsilon)^i$, $(i = 0, \ldots, \log_{1+\varepsilon} U)$

---

## <span style="color:blue">Method 0 : Grid</span>

form uniform grid of side length $\varepsilon r / \sqrt{d}$
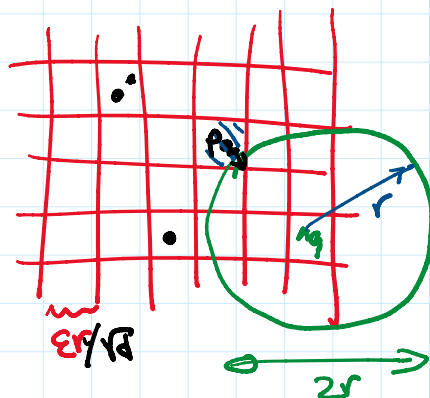
<span style="color:red">dictionary →</span> Store $S$ = all nonempty grid cells

query(q):

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

<span style="color:red">by hashing → in O(1) time per grid cell</span> is in $S$

$\varepsilon r / \sqrt{d}$

$2r$

$\leq r + \varepsilon r = (1+\varepsilon) r$

$\Rightarrow 1+\varepsilon$ approx factor

query time $O(\text{\# grid cells intersecting ball}(q,r))$

$$\leq O\left(\left(\frac{2r}{\varepsilon r / \sqrt{d}}\right)^d\right)$$

$$= \boxed{O\left(\left(\frac{2\sqrt{d}}{\varepsilon}\right)^d\right)}$$

<span style="color:red">← improves to $(O(1))^d$</span>

Space $\boxed{O(dn)}$

$\left(\dfrac{O(1)}{\varepsilon}\right)^d$

Alternative: to reduce query time

store $S$ = all grid cells intersecting $\bigcup\limits_{p \in P} ball(p,r)$

space $\boxed{\left(\dfrac{O(1)}{\varepsilon}\right)^d n}$

query time $\boxed{O(d)}$

Ⓢ: how to avoid exponential dependence on $d$?

## Method 1: Dimension Reduction (Indyk-Motwani '98)

### Johnson-Lindenstrauss Lemma ('84)

For $n$ pts $P$ in $\mathbb{R}^d$, for Euclidean, $^{(or\ L_2)}$

$\exists$ mapping $f: P \to \mathbb{R}^k$, s.t. $\forall p, q \in P$,

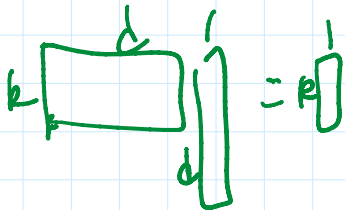$$d(p,q) \le d(f(p), f(q)) \le (1+\varepsilon)\, d(p,q)$$

with $k = \boxed{O\left(\dfrac{1}{\varepsilon^2} \log n\right)}$

(Pf: take a random projection $f$))
Chernoff bd ...

Chernoff bd ...

$(a^{b\log n} = n^{b\log a})$


Combine with grid

space $\left(\dfrac{O(1)}{\varepsilon}\right)^{O\left(\frac{1}{\varepsilon}\log n\right)}$ $n$

$= n^{O\left(\frac{1}{\varepsilon^2}\log\frac{1}{\varepsilon}\right)}$

query $O(dk) = O\left(\frac{1}{\varepsilon^2}d\log n\right)$

no exponential in $d$!

# Method 2: Locality-Sensitive Hashing (LSH)
(Indyk, Motwani '98)

idea- if $p,q$ are close, more likely that $h(p)=h(q)$

Consider special case of <u>Hamming space</u> $\{0,1\}^d$

given $p, q \in \{0,1\}^d$,

define $d(p,q) = $ # bit positions that are diff.

eg. $p = (10110 1)$   $d(p,q)=2$.
    $q = (10010 1)$

Will solve approx decis problem for fixed $r$, appox factor $c$.

let $k$ be param.

define hash fn $h: \{0,1\}^d \to \{0,1\}^k$
by projecting to $k$ rand dims

eg. $p = 10\overset{\downarrow}{1}11\overset{\downarrow}{0}01\overset{\downarrow}{1}$

$h(p) = 010$

for each $p \in P$
put $p$ in bucket for $h(p)$.

Query(q): if each $p \in P$ is bucket for $h(q)$
if $d(p,q) \leq cr$ return $p$
return no

---

**Prop** $\forall p, q \in \{0,1\}^d$,

$$Pr\left[ h(p) = h(q) \right] = \left( 1 - \frac{d(p,q)}{d} \right)^k$$

eg. $p = 10\textcircled{0}11\textcircled{0}001$     $h(p) = 010$
$q = 10\textcircled{0}11\textcircled{1}001$     $h(q) = 010$

$d(p,q) = 2$.