## Near neighbor search via Locality Sensitive Hashing (LSH)
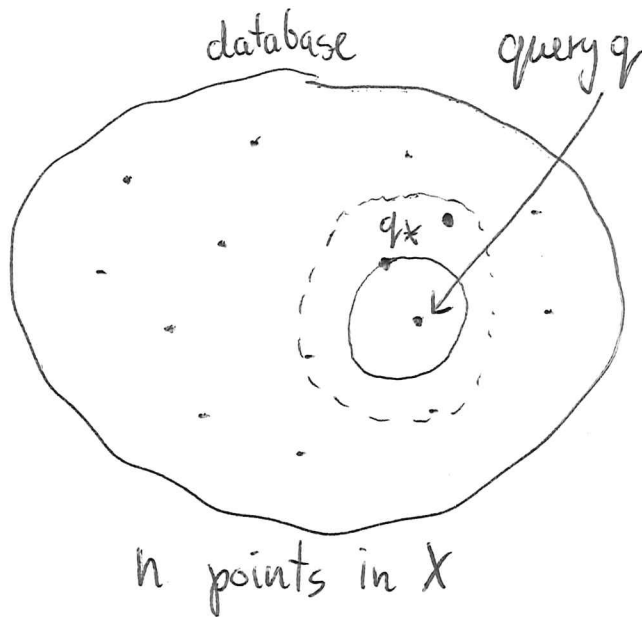
Scenario: new point arrives, find something similar in your database



database     query $q$

$q*$

n points in $X$

Ideal: return the closest point $q*$

With techniques today: return something approximately closest, at distance up to

$$c \cdot dist(q, q*)$$

↑ some fixed constant

"Approximate Nearest Neighbor Search"

12 - 1

# Simplifying the task: Near~~est~~ Neighbor Search

Task: if there is a point at distance $r$,
    return something up to distance $r' = c \cdot r$

Often: To solve "near<u>est</u>" it suffices to consider
    small number of instances of "near"
    with different values of $r$

---

Naïve solution:
- Assumption: all $n$ points live in $d$-dimensional space $X$
- compare all points to the query point: $O(nd)$ time
- Space: $O(nd)$ as well

      Want something faster!

---

# Locality sensitive hash function family $H$

$H$ is an $(r, r', p_1, p_2)$ – locality sensitive hash
                 function family if

for each $u, v \in X$:
$$dist(u, v) \le r \implies Pr\left[h(u) = h(v)\right] \ge p_1$$
and
$$dist(u, v) > r' \implies Pr\left[h(u) = h(v)\right] \le p_2$$

$\boxed{12-2}$    $(0 < r < r')$   $(0 < p_2 < p_1 < 1)$

Example: $X = \{0, 1\}^d$   $\text{dist}(u, v) = \|u - v\|_1$
                                Hamming distance

$$h_i(u_1 u_2 \ldots u_d) = u_i$$
       $\uparrow$
   $i \in [d]$

$H = \{h_i : i \in [d]\}$ is a $(0.2\,d, 0.4\,d, 0.8, 0.6)$–locality
                sensitive hash function family

---

Intuition:

- close points likely to be mapped to the same bucket
- far points likely to be mapped to different buckets

How to use this?

Lots of far from $q$ points can still
be mapped to the same bucket!

---

Step 1: Avoiding too many unwanted collisions

New hash function $g$ concatenates results of many
random $h \in H$

$$g : X \to Y^k$$
$$g(x) = \left( h_1^*(x), h_2^*(x), \ldots, h_k^*(x) \right)$$
                    $\nwarrow$  $\nearrow$      $\underline{\quad}$  $\nearrow$

$\boxed{12-3}$

drawn independently from $H$

Question: What is $\Pr[g(u) = g(v)]$ if $\text{dist}(u,v) \geqslant r'$?

Answer: $\leqslant p_2^k$

Want this to be $\sim \frac{1}{n}$ so expected collisions with far points at most constant

Set $k = \left\lceil \frac{\log n}{\log(1/p_2)} \right\rceil$

Probability a close point in the same bucket?

$$\geqslant p_1^k \geqslant p_1^{\left(\frac{\log n}{\log(1/p_2)} + 1\right)} = p_1 \cdot n^{-\left(\frac{\log(1/p_1)}{\log(1/p_2)}\right)}$$

$$= p_1 \cdot n^{-\rho} \qquad \rho = \frac{\log(1/p_1)}{\log(1/p_2)}$$

$$\rho \in (0, 1)$$

[This could be small!]

Step 2: Repeat hashing for many independent $g$

Probability of $q$ in the same bucket as a given close point $\geqslant p_1 \cdot n^{-\rho}$ Homework!

12-4

Repeat $\Omega\left(1/(p_1 \cdot n^{-\rho})\right) = \Omega\left(\frac{1}{p_1} \cdot n^{\rho}\right)$ times

to "find" the close point, with constant probability

$\underbrace{\quad\quad\quad\quad\quad}$
""

end up in the same bucket
as $q$

---

## Full approximate near neighbor data structure

Preprocessing: - select $k' = \Theta\left(\frac{1}{P_1} \cdot n^S\right)$ hash functions

$g_i$, each a "concatenation of results"

of $k = \left\lceil \frac{\log n}{\log(1/p_2)} \right\rceil$ independently selected

hash functions from $H$

- create $k'$ hash # tables with all
  $n$ points in our database

  Time: $O(k' \cdot k \cdot d \cdot n)$ $\boxed{\begin{array}{l} = O\left(n^{1+S} \, d \, \log n\right) \\ = O\left(nd + n^{1+S}\right) \end{array}}$

  Space: $\nearrow O(nd + k'n)$

    for typical $d$ dimensional
    data

    for constant $p_1$ & $p_2$

$\boxed{12 - 5}$

# Query q

- go over all $k'$ hash tables
- compute the distance of all points hashed to the same buckets as $q$
- stop when you find a point at distance at most $r_2$ from $q$

## Expected time:

$$O\left(\underbrace{k' \cdot k \cdot d}_{\substack{\text{computing} \\ \text{hash functions}}} + \underbrace{k' \cdot O(1) \cdot d}_{\substack{\text{expected distance} \\ \text{computation for} \\ \text{points at distance} \\ > r' \text{ from } q}} + \overset{\substack{\text{distance computation} \\ \text{for the point} \\ \text{close to } q}}{d}\right)$$

$$\boxed{= O\left(n^{\text{S}} d \log n\right)}$$
for constant $p_1$ & $p_2$

vs. naïve $O(nd)$