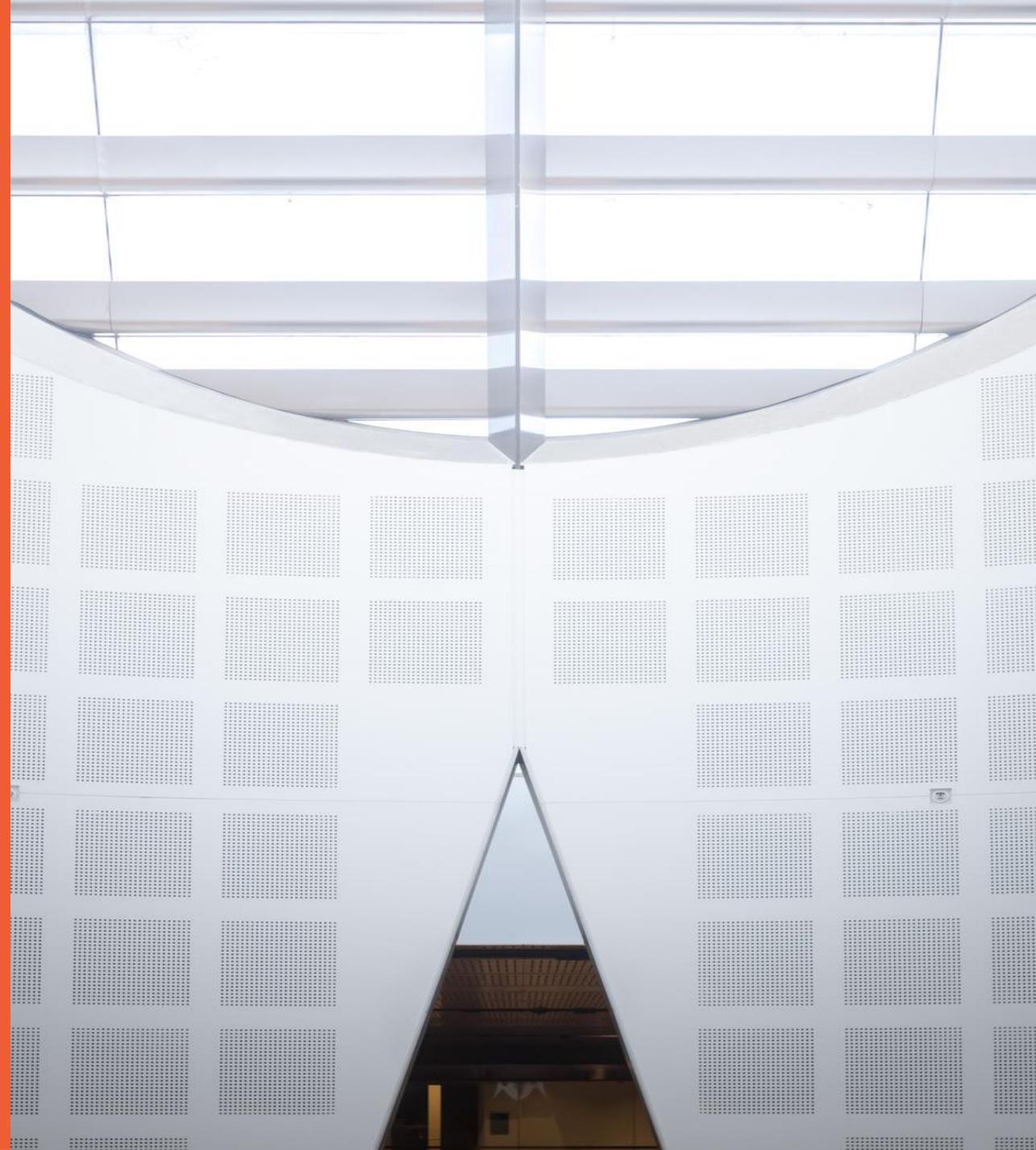COMPx270: Randomised and Advanced Algorithms
Lecture 7: Nearest Neighbours and dimensionality reduction

Clément Canonne

School of Computer Science
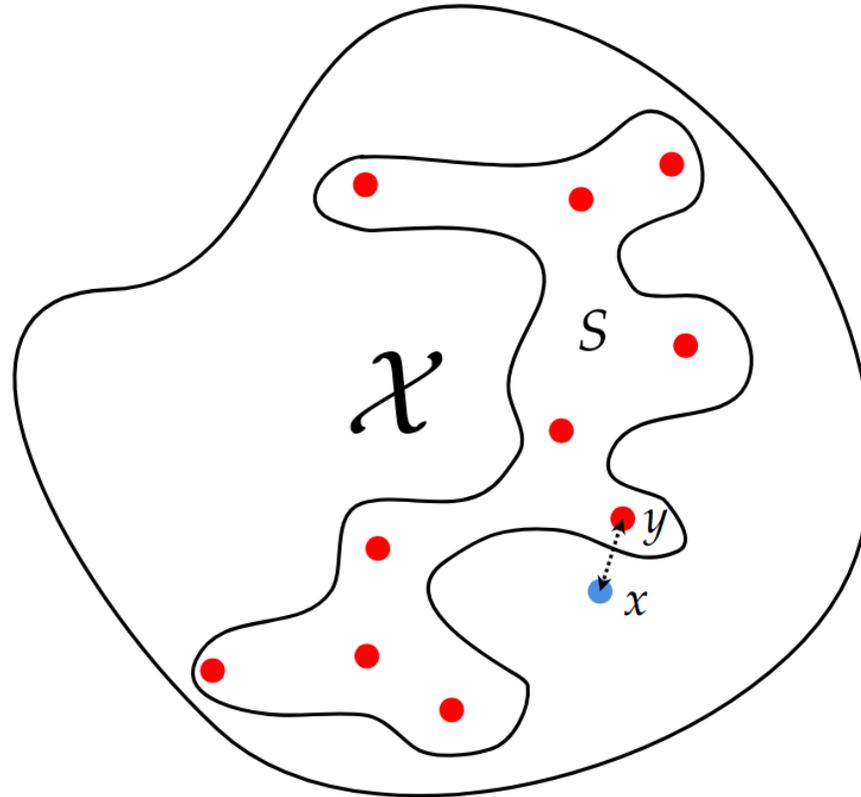
THE UNIVERSITY OF
SYDNEY

# A question 🕷

You have n pictures, each 4096x4096 pixels, of venomous spiders. Someone finds a spider in their kitchen and sends you a photo, asking which type of spider it is and if it is venomous, **because they just have been bitten.**

How long will it take you?

# Nearest Neighbour Search



$|\mathcal{X}| = m$

$\mathcal{X} = \{0,1\}^d$

or

$\mathcal{X} = \mathbb{R}^d$

+ metric $d$

($d(x,y) \equiv$ how similar $x, y$ are)

# Nearest Neighbour Search

$d: \mathcal{X} \times \mathcal{X} \to [0, \infty)$

① $d(x,y) = 0 \iff x = y$

② $d(x,y) = d(y,x)$

③ $d(x,y) \leq d(x,z) + d(z,y)$

$\forall x, y, z$

On $x$: Find $y \in S$ s.t. $y = \underset{y' \in S}{\mathrm{argmin}}\ d(x,y)$

① SPACE    Would like $O(n\ d)$.

↑ # of points    ↑ dimension

② QUERY TIME    Ideally $O(1)$, but... $O(n)$ good

Examples

- $\mathcal{X} = \{0,1\}^d$    Hamming distance: $d(x,y) = \#(\text{bits } x \text{ and } y \text{ differ on})$
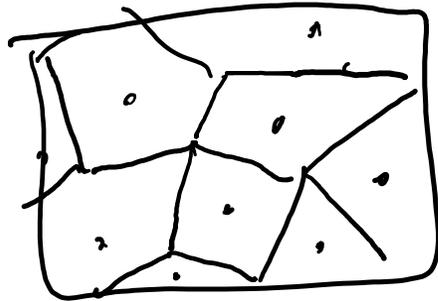- $\mathcal{X} = \mathbb{R}^d$,    Euclidean distance: $d(x,y) = \sqrt{\sum_{i=1}^{d} (x_i - y_i)^2} = \|x-y\|_2$ $(\ell_2)$
- $\mathcal{X} = \mathbb{R}^d$,    Manhattan distance: $d(x,y) = \sum_{i=1}^{d} |x_i - y_i| = \|x-y\|_1$ $(\ell_1)$

# Lists? Voronoi? K-d trees? Hash tables?

- List: query time $O(nd)$, space $O(nd)$
- (For $\{0,1\}^d$): query time $O(2^d)$, Space $O(2^d)$
- Voronoi



query time $O(nd)$, Space $O(n^{\lceil d/2 \rceil})$

- Hash tables

$d \gg 1, n \gg 1$

↑ large   ↑ even larger
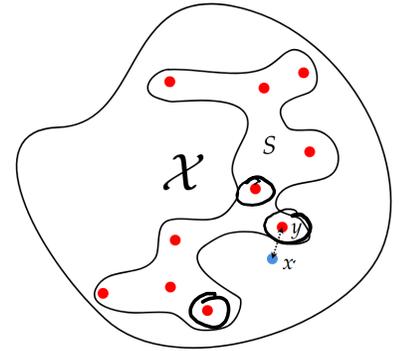
$1 \ll d \ll n \le 2^{O(d)}$

# Bad news...

NN:   we don't know how.

Either query time or space is

$$\Omega\left(\min\left(2^d, nd\right)\right)$$

(for everything we know, <u>even</u> randomized algorithms)

# Approximate Nearest Neighbour Search

$\mathrm{QUERY}(x)$: given an element $x \in \mathcal{X}$, return an element $y \in S$ sort-of-minimising $\mathrm{dist}(x,y)$, that is, $\mathrm{dist}(x,y) \leq C \cdot \min_{y' \in S} \mathrm{dist}(x,y')$.

# Dimensionality Reduction: the JL Lemma (Euclidean space)

$$\left(\mathbb{R}^d, \|\cdot\|_2\right) \xrightarrow{\ \Phi\ } \left(\mathbb{R}^k, \|\cdot\|_2\right) \qquad k \ll d$$

$$\|\phi(x) - \phi(y)\|_2 \overset{C}{\approx} \|x-y\|_2$$

Solve NN on $\mathbb{R}^k$

$\downarrow (O(nk), O(nk))$

Solve ANN on $\mathbb{R}^d$

<u>JL Lemma</u>  Can do this with

$$C = 1 + \varepsilon$$
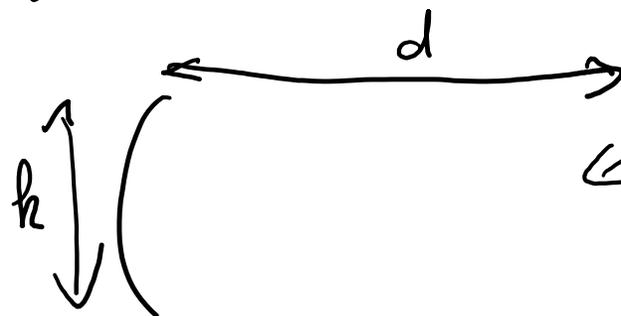
and $k = O\left(\dfrac{\log n}{\varepsilon^2}\right)$

st. if $|S| = n$, $\forall x, y \in S$  $\|\phi(x) - \phi(y)\|_2 \overset{\approx}{=} (1 \pm \varepsilon)\|x-y\|_2$

What is $\phi: \mathbb{R}^d \to \mathbb{R}^k$ ? It's linear.

$M \in \mathbb{R}^{d \times k}$  $\quad k \big\downarrow \Bigg( \xleftarrow{\quad d \quad} \Bigg)$  $M_{ij} \sim N(0,1) \in \mathbb{R}^k$ (with prob. $\geq \frac{99}{100}$)

$\cdot \frac{1}{\sqrt{k}}$  $\phi(x) = Mx$
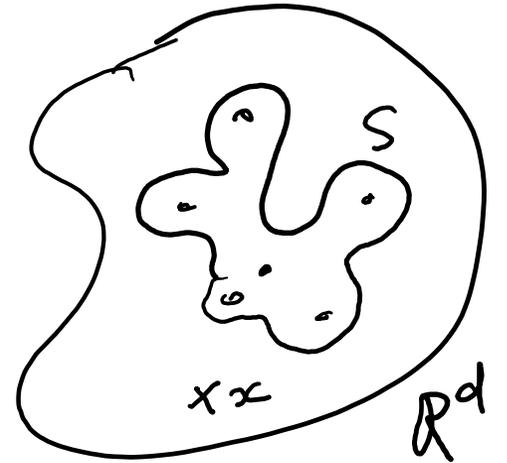
# JL Lemma and ANN

$$\phi: \mathbb{R}^d \longrightarrow \mathbb{R}^k \qquad k = O\left(\frac{\log(n+1)}{\varepsilon^2}\right)$$

Apply to $T = S \cup \{x\}$

Cost: space $O(nk) = O\left(n\frac{\log n}{\varepsilon^2}\right)$

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

Good, but neither is $o(n)$ ...

# Beyond JL Lemma: Hashing!

Spoiler: can do query time $O(n^\rho d)$ (expected)

space $O(n^{1+\rho} d)$

for some $0 < \rho < 1$

(Hamming/$\ell_i$: $\rho \approx \frac{1}{c}$

Euclidean $\rho \approx \frac{1}{c^2}$)

$\downarrow$ sublinear!

$\downarrow$ "nearly" linear

ANN:

find $y$ st

$d(x, y) \leq \boxed{C} \cdot \min_{y'} d(x, y')$

# Locality-Sensitive Hashing

**Definition 36.1.** Let $0 \leq q < p \leq 1, r > 0, C > 1$, and $(\mathcal{X}, \text{dist})$ be a metric space. Then a family of functions $\mathcal{H}$ from $\mathcal{X}$ to $\mathcal{Y}$ is a $(r, C, p, q)$-*Locality Sensitive Hash family* (LSH) if, for every $x, x' \in \mathcal{X}$,
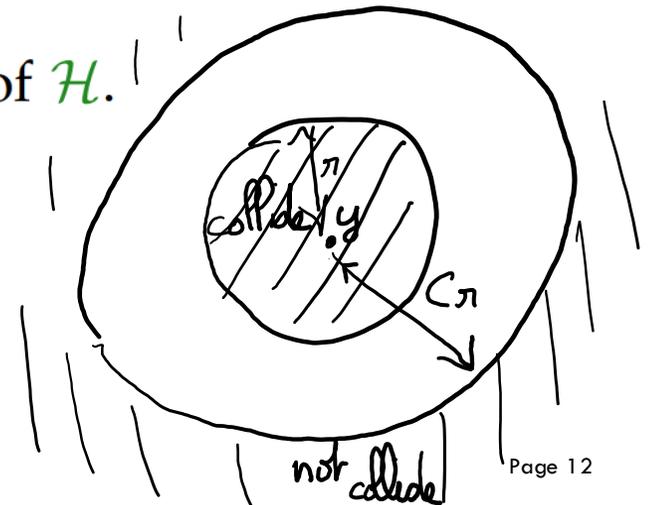
- If $\text{dist}(x, x') \leq r$, then $\Pr_{h \sim \mathcal{H}}[h(x) = h(x')] \geq p$;   ← *WANT collision*

- If $\text{dist}(x, x') \geq Cr$, then $\Pr_{h \sim \mathcal{H}}[h(x) = h(x')] \leq q$;   ← *Do not want collision*

and we say $\rho := \frac{\log(1/p)}{\log(1/q)} < 1$ is the *sensitivity parameter* of $\mathcal{H}$.
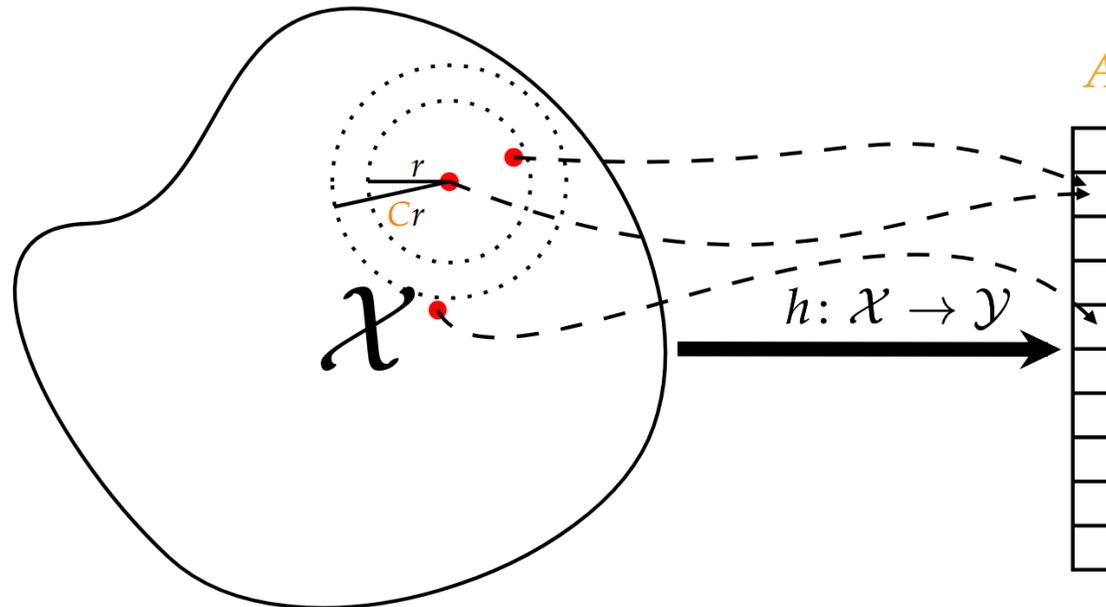
*collide y*

*not collide*

# Locality-Sensitive Hashing

**Definition 36.1.** Let $0 \leq q < p \leq 1$, $r > 0$, $C > 1$, and $(\mathcal{X}, \text{dist})$ be a metric space. Then a family of functions $\mathcal{H}$ from $\mathcal{X}$ to $\mathcal{Y}$ is a $(r, C, p, q)$-*Locality Sensitive Hash family* (LSH) if, for every $x, x' \in \mathcal{X}$,

- If $\text{dist}(x, x') \leq r$, then $\Pr_{h \sim \mathcal{H}}[h(x) = h(x')] \geq p$;

- If $\text{dist}(x, x') \geq Cr$, then $\Pr_{h \sim \mathcal{H}}[h(x) = h(x')] \leq q$;

and we say $\rho := \frac{\log(1/p)}{\log(1/q)} > 1$ is the *sensitivity parameter* of $\mathcal{H}$.

# Locality-Sensitive Hashing: "Baby version"

$\text{QUERY}_r(x)$: given an element $x \in \mathcal{X}$, return an element $y \in S$, or $\perp$, such that:

- If there exists $y^* \in S$ such that $\text{dist}(x, y^*) \leq r$, then, with probability at least $9/10$, $\text{QUERY}_r(x)$ returns an element $y \in S$ such that $\text{dist}(x, y^*) \leq C \cdot r$;

- If $\text{dist}(x, y) > C \cdot r$ for *every* $y \in S$, then, with probability 1, $\text{QUERY}_r(x)$ returns $\perp$.

- Otherwise, any output in $S \cup \{\perp\}$ is allowed.

$\text{QUERY}_r(x)$: given an element $x \in \mathcal{X}$, return an element $y \in S$, or $\perp$, such that:

- If there exists $y^* \in S$ such that $\text{dist}(x, y^*) \leq r$, then, with probability at least 9/10, $\text{QUERY}_r(x)$ returns an element $y \in S$ such that $\text{dist}(x, y^*) \leq C \cdot r$;

- If $\text{dist}(x, y) > C \cdot r$ for *every* $y \in S$, then, with probability 1, $\text{QUERY}_r(x)$ returns $\perp$.

- Otherwise, any output in $S \cup \{\perp\}$ is allowed.

$r$ is fixed

$C > 1$ fixed

$p, q$ given

$0 < q < p < 1$

$\rho = \dfrac{\log(1/p)}{\log(1/q)}$

**Claim.** From $\mathcal{H}$, can get $\mathcal{H}^{(\ell)}$

$(\ell \geq 1)$ s.t. $\mathcal{H}^{(\ell)}$ is a $(r, C, p^\ell, q^\ell)$-LSH family $\left(\text{and } |\mathcal{H}^{(\ell)}| = |\mathcal{H}|^\ell\right)$

**Pf:**
$$h(x) = \left(h_1(x), h_2(x), \dots, h_\ell(x)\right) \in \mathcal{Y}^\ell \qquad \left(\mathcal{H}: \mathcal{X} \to \mathcal{Y}\right)$$

Suppose $d(x, x') \leq r$

$$\Pr_{h \sim \mathcal{H}^{(\ell)}}[h(x) = h(x')] = \Pr_{h_1, h_2, \dots, h_\ell}\left[(h_1(x), \dots, h_\ell(x)) = (h_1(x'), \dots, h_\ell(x'))\right]$$

$$= \underbrace{\Pr_{h_1}[h_1(x) = h_1(x')]}_{\geq p} \cdots \underbrace{\Pr_{h_\ell}[h_\ell(x) = h_\ell(x')]}_{\geq p} \geq p^\ell$$

$\mathcal{H}$ is a $(r, C, p, q)$-LSH

Suppose $d(x, x') \geq C \cdot r$

$$= \underbrace{\phantom{xxxxxxxxxxx}}_{\leq q} \underbrace{\phantom{xxxxxxxxxxx}}_{\leq q} \leq q^\ell$$
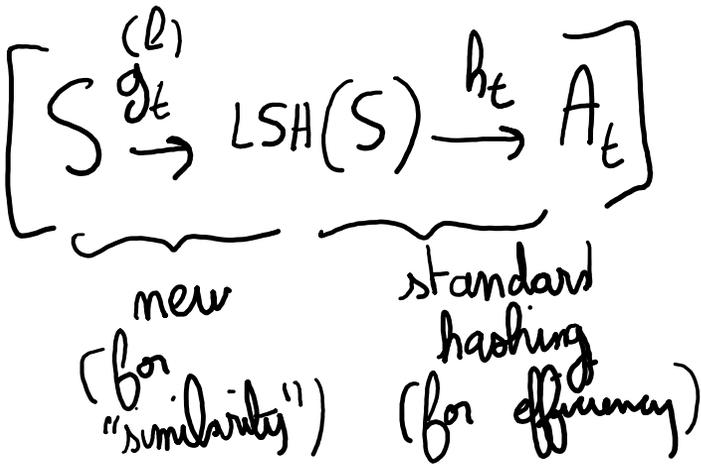
$\textcircled{\ell}$
$\textcircled{k}$

Get $k$ hash tables $A_1, -, A_k$
using good standard hash functions
(not LSH) + chaining $h_1, -, h_k$

In each $A_t$, insert the hashes of $S$ by $g_t^{(\ell)}$

where $g_1^{(\ell)}, -, g_k^{(\ell)} \sim \mathcal{GP}^{(\ell)}$

$\left[ S \xrightarrow{g_t^{(\ell)}} LSH(S) \xrightarrow{h_t} A_t \right]$

new
(for "similarity")

standard hashing
(for efficiency)

PREPROCESS

$\forall x \in S$
$\forall 1 \leq t \leq k$
$A_t . \text{INSERT}( g_t^{(\ell)}(x) )$

QUERY

$\forall 1 \leq t \leq k$
$L_t \leftarrow A_t . \text{LOOKUP}( g_t^{(\ell)}(x) )$
$\forall y \in L_t$
if $d(x,y) \leq Cr$,
return $y$

return $\perp$

$A_1$         $A_k$

Hope If $x, y$ are close
at least **one** $t$ of the $k$ LSH
$g_1^{(\ell)}, -, g_k^{(\ell)}$ will make
them collide, and so
$y \in L_t$ that $t$

# Locality-Sensitive Hashing: "Baby version" (3/4)

Space: $k$ Hash tables, each $n$ elements of size $O(d)$

$\rightarrow O(knd)$

$k$ LSH hash functions

one LSH function from $\mathcal{H}^{(\ell)}$

$k \times \ell \times O(d) = O(k\ell d)$

one LSH function from $\mathcal{H}$

$\Big\} O(knd + k\ell d)$

Query time
- Evaluate $g_t^{(\ell)}(x)$ $\forall \ 1 \leq t \leq k$ : $O(k\ell)$

- $\mathbb{E}\left[\sum_{t=1}^{k} |L_t| \cdot O(d)\right] \approx k \cdot (nq^{\ell}) \cdot O(d) = O(k d n q^{\ell})$

only far ones (bad collisions)

$\Big\} O(k\ell + k d n q^{\ell})$

# Locality-Sensitive Hashing: "Baby version" (4/4)

## Correctness

When unlucky : there is $y^* \in S$ st

$$d(x, y^*) \leq r$$

but $g_t^{(\ell)}(x) \neq g_t^{(\ell)}(y^*) \quad \forall \ 1 \leq t \leq k$

$$\Pr[\text{unlucky}] = (1 - p^\ell)^k \leq \frac{1}{10} \quad \text{④}$$

WANT

---

$k, \ell$ ?

WANT: ④ $+ \ n g^\ell \leq 1$ (for query time)

set $k = \Theta(n^\ell)$ $\longleftarrow$ set $\ell = O\left(\frac{\log n}{\log \frac{1}{q}}\right)$

Space $O(n^{1+\ell} d)$

$\mathbb{E}[\text{query time}] = O(n^\ell d)$

# Locality-Sensitive Hashing: "Baby version" (👶)

# Locality-Sensitive Hashing: "They grow up so fast" ( 👱 )

"Binary search"

$\hookrightarrow$ solves ANN from baby version closing only a $\log d$ factor.

# Locality-Sensitive Hashing: But... do they exist?

Hamming $\{0,1\}^d$

Given $C, r$

If $d(x,x') \leq r$,

$$\Pr_{h}[h(x) = h(x')] \geq 1 - \frac{r}{d}$$

If $d(x,x') \geq C \cdot r$

$$\Pr_{h}[h(x) = h(x')] \leq 1 - \frac{Cr}{d}$$

$h_1(x) = x_1 \in \{0,1\}$

$h_2(x) = x_2$

$\vdots$

$h_d(x) = x_d$

$\mathcal{H} = \{h_i : i \in [d]\}$

$$\rho = \frac{\log(1/p)}{\log(1/q)} = \frac{\log\left(1 - \frac{r}{d}\right)}{\log\left(1 - \frac{Cr}{d}\right)} \approx \frac{1}{C}$$

# Locality-Sensitive Hashing: But... do they exist?