

# Beyond Locality-Sensitive Hashing

Alexandr Andoni<sup>1</sup>   Piotr Indyk<sup>2</sup>   Huy L. Nguyễn<sup>3</sup>  
**Ilya Razenshteyn<sup>2</sup>**

<sup>1</sup>Microsoft Research SVC

<sup>2</sup>MIT, CSAIL

<sup>3</sup>Princeton

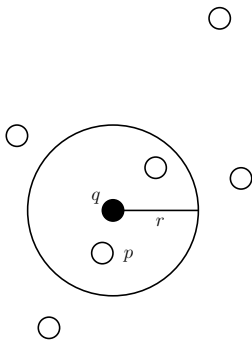
Brown ICERM Theory Seminar

# The Near Neighbor Problem

- Let  $P$  be an  $n$ -point subset of a metric  $(X, D)$ ,  $r > 0$
- For  $q \in X$  find any  $p \in P$  with  $D(p, q) \leq r$

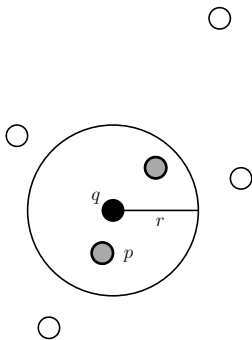
# The Near Neighbor Problem

- Let  $P$  be an  $n$ -point subset of a metric  $(X, D)$ ,  $r > 0$
- For  $q \in X$  find any  $p \in P$  with  $D(p, q) \leq r$



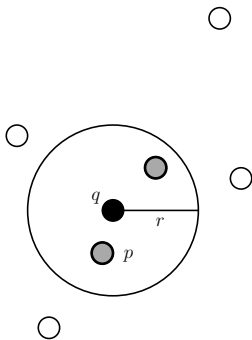
# The Near Neighbor Problem

- Let  $P$  be an  $n$ -point subset of a metric  $(X, D)$ ,  $r > 0$
- For  $q \in X$  find any  $p \in P$  with  $D(p, q) \leq r$



# The Near Neighbor Problem

- Let  $P$  be an  $n$ -point subset of a metric  $(X, D)$ ,  $r > 0$
- For  $q \in X$  find any  $p \in P$  with  $D(p, q) \leq r$
- Hard, if  $(X, D)$  is *high-dimensional* (space or query time is exponential in the dimension)



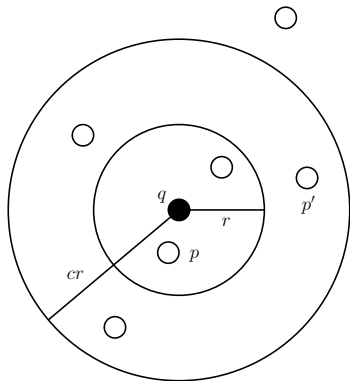
- Machine Learning / Data Mining
- Spell checkers
- Computational Geometry
- Databases
- Clustering
- You name it ...

# The Approximate Near Neighbor Problem (ANN)

- Let  $P$  be an  $n$ -point subset of a metric  $(X, D)$ ,  $r > 0$ ,  $c > 1$
- For  $q \in X$  find any  $p' \in P$  with  $D(p', q) \leq cr$ , provided that there exists  $p \in P$  with  $D(p, q) \leq r$

# The Approximate Near Neighbor Problem (ANN)

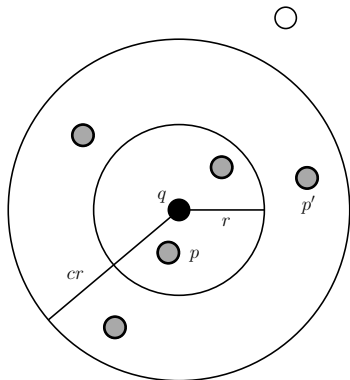
- Let  $P$  be an  $n$ -point subset of a metric  $(X, D)$ ,  $r > 0$ ,  $c > 1$
- For  $q \in X$  find any  $p' \in P$  with  $D(p', q) \leq cr$ , provided that there exists  $p \in P$  with  $D(p, q) \leq r$





# The Approximate Near Neighbor Problem (ANN)

- Let  $P$  be an  $n$ -point subset of a metric  $(X, D)$ ,  $r > 0$ ,  $c > 1$
- For  $q \in X$  find any  $p' \in P$  with  $D(p', q) \leq cr$ , provided that there exists  $p \in P$  with  $D(p, q) \leq r$



- **Exponential dependence on the dimension:**

(Arya, Mount 1993), (Meister 1993), (Clarkson 1994),  
(Arya, Mount, Netanyahu, Silverman, We, 1998), (Kleinberg, 1997),  
(Har-Peled 2002)

- **Polynomial dependence on the dimension:**

(Indyk, Motwani 1998), (Kushilevitz, Ostrovsky, Rabani 1998),  
(Indyk 1998), (Indyk 2001), (Gionis, Indyk, Motwani 1999),  
(Charikar 2002), (Datar, Immorlica, Indyk, Mirrokni 2004),  
(Chakrabarti, Regev 2004), (Panigrahy 2006), (Ailon, Chazelle 2006),  
(Andoni, Indyk 2006), (Indyk, Kapralov 2013), (Nguyễn 2013),  
(Andoni 2014)

# Two regimes

What is known about  $c$ -approximate ANN in  $\mathbb{R}^d$ ?

approximation	space	query time
$c = 1 + \varepsilon$	$(dn)^{O_\varepsilon(1)}$	$(d \log n / \varepsilon)^{O(1)}$
$c \rightarrow \infty$ $\rho = \rho(c) \rightarrow 0$	$d^{O(1)} \cdot n^{1+\rho+o(1)}$	$d^{O(1)} \cdot n^{\rho+o(1)}$

This talk is exclusively about the second regime

# Locality-Sensitive Hashing (LSH)

- The goal: solve ANN with polynomial in the dimension space and query time, near-linear in  $n$  space, and sublinear in  $n$  query time

# Locality-Sensitive Hashing (LSH)

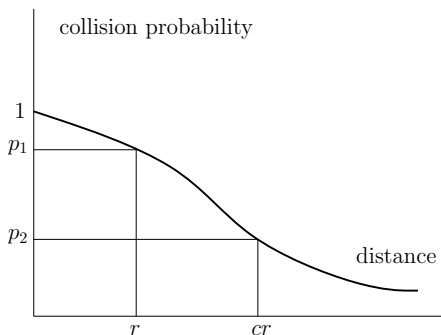
- The goal: solve ANN with polynomial in the dimension space and query time, near-linear in  $n$  space, and sublinear in  $n$  query time
- The only known technique: **Locality-Sensitive Hashing (LSH)** (Indyk, Motwani 1998)

# Locality-Sensitive Hashing (LSH)

- The goal: solve ANN with polynomial in the dimension space and query time, near-linear in  $n$  space, and sublinear in  $n$  query time
- The only known technique: **Locality-Sensitive Hashing (LSH)** (Indyk, Motwani 1998)
- A hash family  $\mathcal{H}$  on  $(X, D)$  is  $(r, cr, p_1, p_2)$ -sensitive, if for every  $p, q \in X$ :
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$

# Locality-Sensitive Hashing (LSH)

- The goal: solve ANN with polynomial in the dimension space and query time, near-linear in  $n$  space, and sublinear in  $n$  query time
- The only known technique: **Locality-Sensitive Hashing (LSH)** (Indyk, Motwani 1998)
- A hash family  $\mathcal{H}$  on  $(X, D)$  is  $(r, cr, p_1, p_2)$ -sensitive, if for every  $p, q \in X$ :
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$



- Let  $\mathcal{H}$  be a “reasonable”  $(r, cr, p_1, p_2)$ -sensitive family



- Let  $\mathcal{H}$  be a “reasonable”  $(r, cr, p_1, p_2)$ -sensitive family
- Define “quality” of  $\mathcal{H}$  as

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)}$$

- Let  $\mathcal{H}$  be a “reasonable”  $(r, cr, p_1, p_2)$ -sensitive family
- Define “quality” of  $\mathcal{H}$  as

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)}$$

- Then, can solve ANN with roughly  $O(n^{1+\rho} + nd)$  space and  $O(d \cdot n^\rho)$  query time (Indyk, Motwani 1998)

- Let  $\mathcal{H}$  be a “reasonable”  $(r, cr, p_1, p_2)$ -sensitive family
- Define “quality” of  $\mathcal{H}$  as

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)}$$

- Then, can solve ANN with roughly  $O(n^{1+\rho} + nd)$  space and  $O(d \cdot n^\rho)$  query time ([Indyk, Motwani 1998](#))
- Example:  $\{0, 1\}^d$  with Hamming distance; Let  $\mathcal{H} = \{h_1, \dots, h_d\}$ , where  $h_i(x) = x_i$ ; One can check that  $\rho \leq 1/c$

- Let  $\mathcal{H}$  be a “reasonable”  $(r, cr, p_1, p_2)$ -sensitive family
- Define “quality” of  $\mathcal{H}$  as

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)}$$

- Then, can solve ANN with roughly  $O(n^{1+\rho} + nd)$  space and  $O(d \cdot n^\rho)$  query time (Indyk, Motwani 1998)
- Example:  $\{0, 1\}^d$  with Hamming distance; Let  $\mathcal{H} = \{h_1, \dots, h_d\}$ , where  $h_i(x) = x_i$ ; One can check that  $\rho \leq 1/c$

11101110  
10111101

# Known LSH constructions

(Indyk, Motwani 1998), (Andoni, Indyk 2006),  
(Motwani, Naor, Panigrahy 2007), (O'Donnell, Wu, Zhou 2011),  
(Indyk, Kapralov 2013), (Nguyễn 2013), (Andoni 2014)

Bounds on  $\rho = \ln(1/p_1)/\ln(1/p_2)$  for various spaces:

Space	Upper bound	Lower bound
$l_1$	$\rho \leq 1/c$	$\rho \geq 1/c - o(1)$
$l_p$ $1 < p < 2$	$\rho \leq 1/c^p + o(1)$	$\rho \geq 1/c^p - o(1)$
$l_2$	$\rho \leq 1/c^2 + o(1)$	$\rho \geq 1/c^2 - o(1)$

# Known LSH constructions

(Indyk, Motwani 1998), (Andoni, Indyk 2006),  
(Motwani, Naor, Panigrahy 2007), (O'Donnell, Wu, Zhou 2011),  
(Indyk, Kapralov 2013), (Nguyễn 2013), (Andoni 2014)

Bounds on  $\rho = \ln(1/p_1)/\ln(1/p_2)$  for various spaces:

Space	Upper bound	Lower bound
$\ell_1$	$\rho \leq 1/c$	$\rho \geq 1/c - o(1)$
$\ell_p$ $1 < p < 2$	$\rho \leq 1/c^p + o(1)$	$\rho \geq 1/c^p - o(1)$
$\ell_2$	$\rho \leq 1/c^2 + o(1)$	$\rho \geq 1/c^2 - o(1)$

This work: ANN in space  $O(n^{1+\tau} + nd)$  and time  $O(dn^\tau)$ , where

- $\tau \leq \frac{7}{8c} + O\left(\frac{1}{c^{3/2}}\right) + o(1)$  for  $\ell_1$
- $\tau \leq \frac{7}{8c^2} + O\left(\frac{1}{c^3}\right) + o(1)$  for  $\ell_2$

The first improvement upon (Indyk, Motwani 1998) for  $\ell_1$  and  
(Andoni, Indyk 2006) for  $\ell_2$ !

# The main idea

- LSH is oblivious, but can we construct a hash family that would depend on data (think (Fredman, Komlós, Szemerédi 1984))?

# The main idea

- LSH is oblivious, but can we construct a hash family that would depend on data (think (Fredman, Komlós, Szemerédi 1984))?
- $\mathcal{H}$  is  $(r, cr, p_1, p_2)$ -sensitive, if for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$



# The main idea

- LSH is oblivious, but can we construct a hash family that would depend on data (think (Fredman, Komlós, Szemerédi 1984))?
- $\mathcal{H}$  is  $(r, cr, p_1, p_2)$ -sensitive, if for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$
- **Too strong!** Enough to satisfy these for  $p \in P$  and  $q \in X$ . Can exploit the geometry of  $P$  to construct a better family

# The main idea

- LSH is oblivious, but can we construct a hash family that would depend on data (think (Fredman, Komlós, Szemerédi 1984))?
- $\mathcal{H}$  is  $(r, cr, p_1, p_2)$ -sensitive, if for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$
- **Too strong!** Enough to satisfy these for  $p \in P$  and  $q \in X$ . Can exploit the geometry of  $P$  to construct a better family
- Parallels with practice!
  - PCA trees (Sproull 1991), (McNames 2001), (Verma, Kpotufe, Dasgupta 2009)
  - Spectral Hashing (Weiss, Torralba, Fergus 2008)
  - Semantic Hashing (Salakhutdinov, Hinton 2009)
  - WTA Hashing (Yagnik, Strelow, Ross, Lin 2011)

## The main idea (contd)

- From now on, looking at the Euclidean case and trying to improve upon  $\rho \leq 1/c^2$  (Andoni, Indyk 2006)

## The main idea (contd)

- From now on, looking at the Euclidean case and trying to improve upon  $\rho \leq 1/c^2$  (Andoni, Indyk 2006)
- Partition  $P$  into low-diameter clusters (of diameter  $O(cr)$ )
- Improve upon  $1/c^2$  for the low-diameter case

# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$

# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$
- Can achieve

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

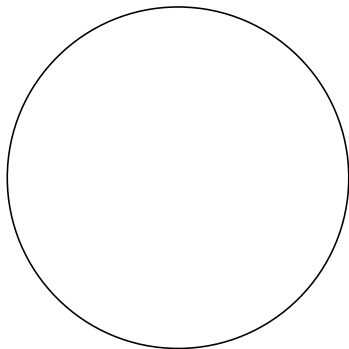
using “ball carving” (similar to [\(Karger, Motwani, Sudan 1998\)](#))

# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$
- Can achieve

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

using “ball carving” (similar to [\(Karger, Motwani, Sudan 1998\)](#))

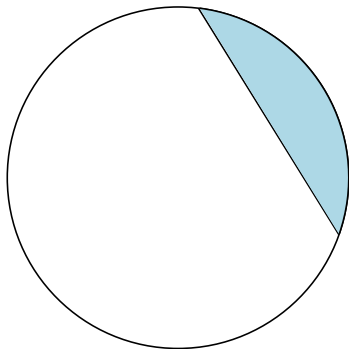


# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$
- Can achieve

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

using “ball carving” (similar to [\(Karger, Motwani, Sudan 1998\)](#))



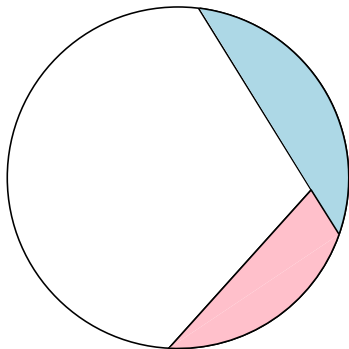


# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$
- Can achieve

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

using “ball carving” (similar to [\(Karger, Motwani, Sudan 1998\)](#))

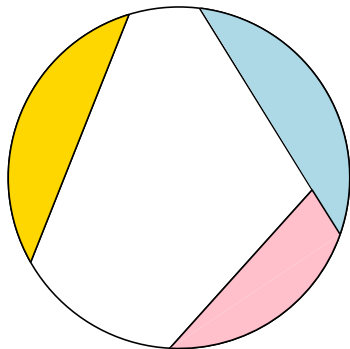


# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$
- Can achieve

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

using “ball carving” (similar to [\(Karger, Motwani, Sudan 1998\)](#))

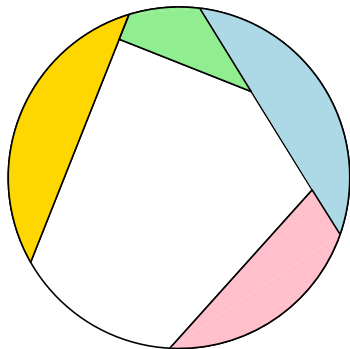


# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$
- Can achieve

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

using “ball carving” (similar to [\(Karger, Motwani, Sudan 1998\)](#))

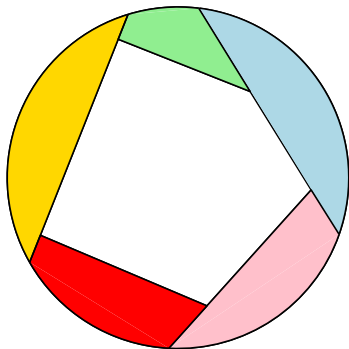


# The low-diameter case

- All points and queries are on a sphere of radius  $O(cr)$
- Can achieve

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

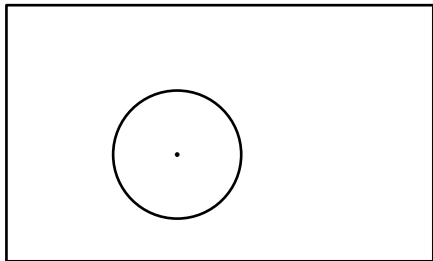
using “ball carving” (similar to [\(Karger, Motwani, Sudan 1998\)](#))



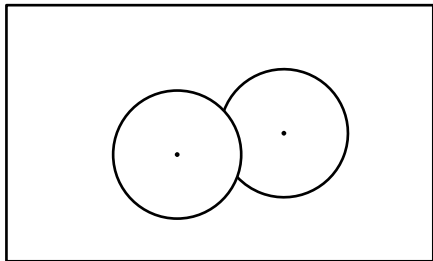


- Choose radius  $R$
- Sample a point  $x$  wrt to a measure  $\mu$ , “carve” a ball  $B(x, R)$ , repeat until everything is covered

# Ball carving

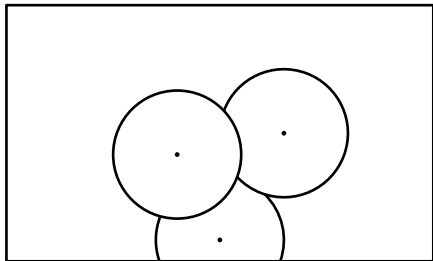


- Choose radius  $R$
- Sample a point  $x$  wrt to a measure  $\mu$ , “carve” a ball  $B(x, R)$ , repeat until everything is covered



- Choose radius  $R$
- Sample a point  $x$  wrt to a measure  $\mu$ , “carve” a ball  $B(x, R)$ , repeat until everything is covered

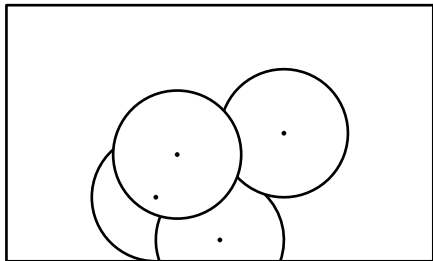
# Ball carving



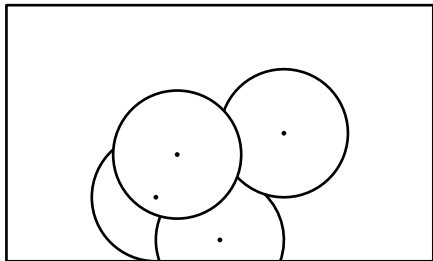
- Choose radius  $R$
- Sample a point  $x$  wrt to a measure  $\mu$ , “carve” a ball  $B(x, R)$ , repeat until everything is covered



# Ball carving

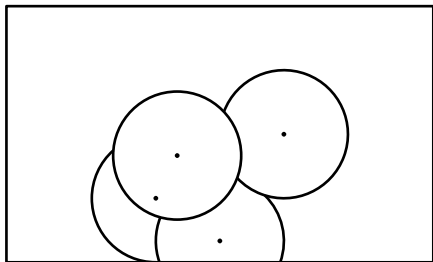


- Choose radius  $R$
- Sample a point  $x$  wrt to a measure  $\mu$ , “carve” a ball  $B(x, R)$ , repeat until everything is covered



- Choose radius  $R$
- Sample a point  $x$  wrt to a measure  $\mu$ , “carve” a ball  $B(x, R)$ , repeat until everything is covered
- Probability of collision?

$$\Pr[x \text{ and } y \text{ collide}] = \frac{\mu(B(x, R) \cap B(y, R))}{\mu(B(x, R) \cup B(y, R))}$$



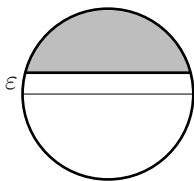
- Choose radius  $R$
- Sample a point  $x$  wrt to a measure  $\mu$ , “carve” a ball  $B(x, R)$ , repeat until everything is covered
- Probability of collision?

$$\Pr[x \text{ and } y \text{ collide}] = \frac{\mu(B(x, R) \cap B(y, R))}{\mu(B(x, R) \cup B(y, R))}$$

- Number of balls  $\approx \mu(X)/\mu(B(x, R))$

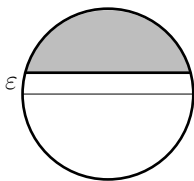
- Points and queries lie on a sphere of radius  $\eta cr$ , for  $\eta = O(1)$  (say,  $\eta = 10$ )

## Ball carving for spheres



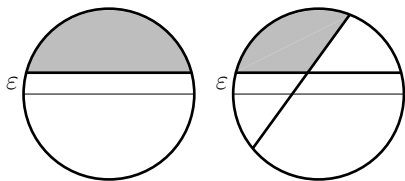
- Points and queries lie on a sphere of radius  $\eta cr$ , for  $\eta = O(1)$  (say,  $\eta = 10$ )
- What are the balls? Spherical caps!  $\{x \mid x_1 \geq \epsilon \eta cr\}$

## Ball carving for spheres



- Points and queries lie on a sphere of radius  $\eta cr$ , for  $\eta = O(1)$  (say,  $\eta = 10$ )
- What are the balls? Spherical caps!  $\{x \mid x_1 \geq \varepsilon \eta cr\}$
- In order to have  $2^{o(d)} = n^{o(1)}$  (can achieve  $d = O(\log n)$ ) balls one should choose  $\varepsilon = o(1)$  (due to measure concentration), but one can not set  $\varepsilon = 0$

# Ball carving for spheres



- Points and queries lie on a sphere of radius  $\eta cr$ , for  $\eta = O(1)$  (say,  $\eta = 10$ )
- What are the balls? Spherical caps!  $\{x \mid x_1 \geq \epsilon \eta cr\}$
- In order to have  $2^{o(d)} = n^{o(1)}$  (can achieve  $d = O(\log n)$ ) balls one should choose  $\epsilon = o(1)$  (due to measure concentration), but one can not set  $\epsilon = 0$
- Computing the intersection of two caps is non-trivial. . .

# Gaussian LSH

- Sample  $w \sim N(0, 1)^d$
- Carve  $\{u \mid \|u\| = \eta cr, \langle u, w \rangle \geq \varepsilon \sqrt{d} \cdot \eta cr\}$ , repeat
- The same as in (Karger, Motwani, Sudan 1998)

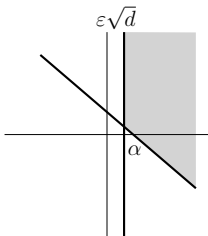


# Gaussian LSH

- Sample  $w \sim N(0, 1)^d$
- Carve  $\{u \mid \|u\| = \eta cr, \langle u, w \rangle \geq \varepsilon\sqrt{d} \cdot \eta cr\}$ , repeat
- The same as in (Karger, Motwani, Sudan 1998)
- Probability of collision?

$$\Pr_{X, Y \sim N(0, 1)}[X \geq \varepsilon\sqrt{d} \wedge \cos \alpha \cdot X - \sin \alpha \cdot Y \geq \varepsilon\sqrt{d}],$$

where  $\alpha$  is the angle between points of interest



- Easy computations show that this family for  $\varepsilon = d^{-1/4}$  is  $(r, cr, p_1, p_2)$ -sensitive with
  - $p_1, p_2 = \exp(-o(d))$  ( $= n^{-o(1)}$ )
  - one has

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \approx \left(1 - \frac{1}{4\eta^2}\right) \cdot \frac{1}{c^2} = \frac{1 - \Omega_\eta(1)}{c^2}$$

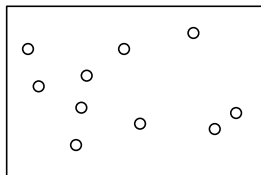
- Easy computations show that this family for  $\varepsilon = d^{-1/4}$  is  $(r, cr, p_1, p_2)$ -sensitive with
  - $p_1, p_2 = \exp(-o(d))$  ( $= n^{-o(1)}$ )
  - one has

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \approx \left(1 - \frac{1}{4\eta^2}\right) \cdot \frac{1}{c^2} = \frac{1 - \Omega_\eta(1)}{c^2}$$

- Simplification and improvement of the result from [\(Andoni, Indyk 2006\)](#) for the case, when everything lies in a ball of radius  $O(cr)$ !

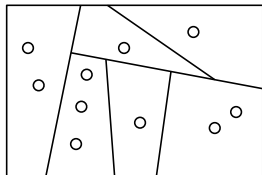
# From LSH to ANN: the basic reduction

- Let  $\mathcal{H}$  be a  $(r, cr, p_1, p_2)$ -sensitive family: for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$



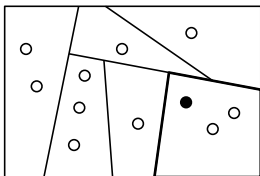
# From LSH to ANN: the basic reduction

- Let  $\mathcal{H}$  be a  $(r, cr, p_1, p_2)$ -sensitive family: for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$
- Hash the dataset  $P$  using a concatenation of  $k$  functions from  $\mathcal{H}$ :  
 $x \mapsto (h_1(x), h_2(x), \dots, h_k(x))$



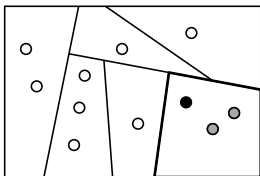
# From LSH to ANN: the basic reduction

- Let  $\mathcal{H}$  be a  $(r, cr, p_1, p_2)$ -sensitive family: for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$
- Hash the dataset  $P$  using a concatenation of  $k$  functions from  $\mathcal{H}$ :  
 $x \mapsto (h_1(x), h_2(x), \dots, h_k(x))$
- Locate a query  $q$  and enumerate all points from the bucket



# From LSH to ANN: the basic reduction

- Let  $\mathcal{H}$  be a  $(r, cr, p_1, p_2)$ -sensitive family: for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$
- Hash the dataset  $P$  using a concatenation of  $k$  functions from  $\mathcal{H}$ :  
 $x \mapsto (h_1(x), h_2(x), \dots, h_k(x))$
- Locate a query  $q$  and enumerate all points from the bucket
- The optimal choice of  $k$  leads to the need in  $n^\rho$  hash tables
- Overall:  $n^{1+\rho}$  space,  $n^\rho$  query time



- For every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$ .



- For every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$ .
- $\mathcal{H}^{\otimes k}$ : sample  $k$  independent functions  $h_1, h_2, \dots, h_k$  according to  $\mathcal{H}$ , then hash  $p \in X$  into  $(h_1(p), h_2(p), \dots, h_k(p))$ .
- if  $\mathcal{H}$  is  $(r, cr, p_1, p_2)$ -sensitive, then  $\mathcal{H}^{\otimes k}$  is  $(r, cr, p_1^k, p_2^k)$ -sensitive

- For every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \geq p_1$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}}[h(p) = h(q)] \leq p_2$ .
- $\mathcal{H}^{\otimes k}$ : sample  $k$  independent functions  $h_1, h_2, \dots, h_k$  according to  $\mathcal{H}$ , then hash  $p \in X$  into  $(h_1(p), h_2(p), \dots, h_k(p))$ .
- if  $\mathcal{H}$  is  $(r, cr, p_1, p_2)$ -sensitive, then  $\mathcal{H}^{\otimes k}$  is  $(r, cr, p_1^k, p_2^k)$ -sensitive
- Choose  $k$  s.t. for every  $p, q \in X$ 
  - if  $D(p, q) \leq r$ , then  $\Pr_{h \sim \mathcal{H}^{\otimes k}}[h(p) = h(q)] \geq n^{-\rho - o(1)}$ ;
  - if  $D(p, q) \geq cr$ , then  $\Pr_{h \sim \mathcal{H}^{\otimes k}}[h(p) = h(q)] \leq 1/n$ ,where  $\rho = \ln(1/p_1) / \ln(1/p_2)$ .

- $\mathcal{H}^{\otimes k}$  is  $(r, cr, n^{-(\rho+o(1))}, n^{-1})$ -sensitive

- $\mathcal{H}^{\otimes k}$  is  $(r, cr, n^{-(\rho+o(1))}, n^{-1})$ -sensitive
- Hash  $P$  wrt to a random function from  $h \sim \mathcal{H}^{\otimes k}$
- For a query  $q$  look at the bin  $h(q)$

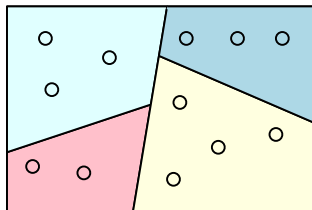
- $\mathcal{H}^{\otimes k}$  is  $(r, cr, n^{-(\rho+o(1))}, n^{-1})$ -sensitive
- Hash  $P$  wrt to a random function from  $h \sim \mathcal{H}^{\otimes k}$
- For a query  $q$  look at the bin  $h(q)$
- By Markov inequality, on average collides with at most one  $p \in P$  s.t.  $D(p, q) \geq cr$
- If  $p \in P$  with  $D(p, q) \leq r$ , then the probability of collision is at least  $n^{-(\rho+o(1))}$

## From LSH to ANN: the proof II

- $\mathcal{H}^{\otimes k}$  is  $(r, cr, n^{-(\rho+o(1))}, n^{-1})$ -sensitive
- Hash  $P$  wrt to a random function from  $h \sim \mathcal{H}^{\otimes k}$
- For a query  $q$  look at the bin  $h(q)$
- By Markov inequality, on average collides with at most one  $p \in P$  s.t.  $D(p, q) \geq cr$
- If  $p \in P$  with  $D(p, q) \leq r$ , then the probability of collision is at least  $n^{-(\rho+o(1))}$
- Consider  $n^{\rho+o(1)}$  independent hash tables (each uses its own  $h \sim \mathcal{H}^{\otimes k}$ )
- Space  $\approx n^{1+\rho+o(1)}$ , query time  $\approx n^{\rho+o(1)}$

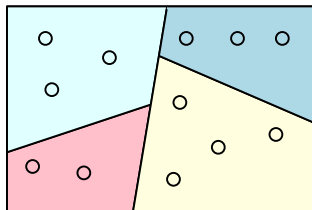
# From LSH to ANN: two-level hashing

- Partition space somewhat coarsely (using smaller  $k$  than before)



# From LSH to ANN: two-level hashing

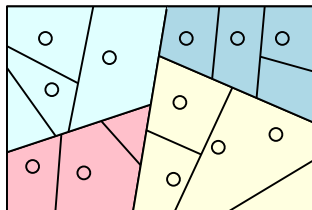
- Partition space somewhat coarsely (using smaller  $k$  than before)
- Argue that every part has a low diameter (aim at  $O(cr)$ )





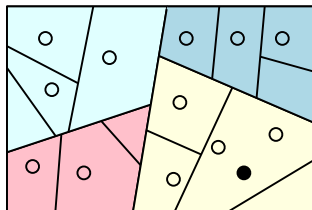
# From LSH to ANN: two-level hashing

- Partition space somewhat coarsely (using smaller  $k$  than before)
- Argue that every part has a low diameter (aim at  $O(cr)$ )
- Use the better family for the low-diameter case to partition space finer



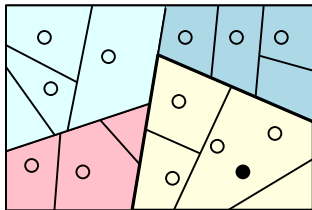
# From LSH to ANN: two-level hashing

- Partition space somewhat coarsely (using smaller  $k$  than before)
- Argue that every part has a low diameter (aim at  $O(cr)$ )
- Use the better family for the low-diameter case to partition space finer
- “Outer” (data-independent,  $\rho \leq 1/c^2$ ) + “inner” (data-dependent, “low-diameter” family) hash tables



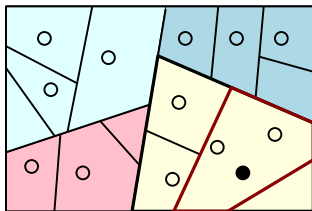
# From LSH to ANN: two-level hashing

- Partition space somewhat coarsely (using smaller  $k$  than before)
- Argue that every part has a low diameter (aim at  $O(cr)$ )
- Use the better family for the low-diameter case to partition space finer
- “Outer” (data-independent,  $\rho \leq 1/c^2$ ) + “inner” (data-dependent, “low-diameter” family) hash tables



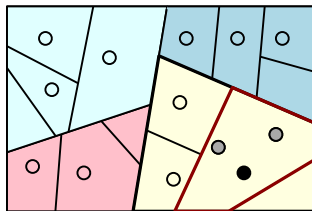
# From LSH to ANN: two-level hashing

- Partition space somewhat coarsely (using smaller  $k$  than before)
- Argue that every part has a low diameter (aim at  $O(cr)$ )
- Use the better family for the low-diameter case to partition space finer
- “Outer” (data-independent,  $\rho \leq 1/c^2$ ) + “inner” (data-dependent, “low-diameter” family) hash tables



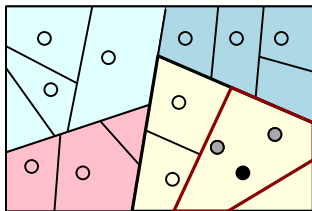
# From LSH to ANN: two-level hashing

- Partition space somewhat coarsely (using smaller  $k$  than before)
- Argue that every part has a low diameter (aim at  $O(cr)$ )
- Use the better family for the low-diameter case to partition space finer
- “Outer” (data-independent,  $\rho \leq 1/c^2$ ) + “inner” (data-dependent, “low-diameter” family) hash tables



# From LSH to ANN: two-level hashing

- Partition space somewhat coarsely (using smaller  $k$  than before)
- Argue that every part has a low diameter (aim at  $O(cr)$ )
- Use the better family for the low-diameter case to partition space finer
- “Outer” (data-independent,  $\rho \leq 1/c^2$ ) + “inner” (data-dependent, “low-diameter” family) hash tables
- Get  $\rho \leq (1 - \Omega(1))/c^2$ , since in inner tables we get better relation between  $p_1$  and  $p_2$ !



# The proof idea

- Let  $\tau > 1$  be some parameter
- Let  $\mathcal{H}_1$  be the LSH family for  $\ell_2$  that achieves  $\rho \lesssim 1/c^2$ , choose  $k$  s.t. the probabilities of collision for  $\mathcal{H}_1^{\otimes k}$  are as follows

Distance	$r$	$cr$	$\tau cr$
Collision probability	$n^{-2}/(\tau c)^2$	$n^{-2}/\tau^2$	$n^{-2}$

# The proof idea

- Let  $\tau > 1$  be some parameter
- Let  $\mathcal{H}_1$  be the LSH family for  $\ell_2$  that achieves  $\rho \lesssim 1/c^2$ , choose  $k$  s.t. the probabilities of collision for  $\mathcal{H}_1^{\otimes k}$  are as follows

Distance	$r$	$cr$	$\tau cr$
Collision probability	$n^{-2}/(\tau c)^2$	$n^{-2}/\tau^2$	$n^{-2}$

- Whp all bins have diameter at most  $\tau cr$  by Markov inequality



# The proof idea

- Let  $\tau > 1$  be some parameter
- Let  $\mathcal{H}_1$  be the LSH family for  $\ell_2$  that achieves  $\rho \lesssim 1/c^2$ , choose  $k$  s.t. the probabilities of collision for  $\mathcal{H}_1^{\otimes k}$  are as follows

Distance	$r$	$cr$	$\tau cr$
Collision probability	$n^{-2/(\tau c)^2}$	$n^{-2/\tau^2}$	$n^{-2}$

- Whp all bins have diameter at most  $\tau cr$  by Markov inequality
- For each bin use (a power of) Gaussian LSH using any point as a center (do not care about points that are more than  $\tau cr + r$  from the center)

Distance	$r$	$cr$
Collision probability	$n^{(-1+2/\tau^2)(1-\Omega_\tau(1))}/c^2$	$n^{-1+2/\tau^2}$

# The proof idea

- Let  $\tau > 1$  be some parameter
- Let  $\mathcal{H}_1$  be the LSH family for  $\ell_2$  that achieves  $\rho \lesssim 1/c^2$ , choose  $k$  s.t. the probabilities of collision for  $\mathcal{H}_1^{\otimes k}$  are as follows

Distance	$r$	$cr$	$\tau cr$
Collision probability	$n^{-2/(\tau c)^2}$	$n^{-2/\tau^2}$	$n^{-2}$

- Whp all bins have diameter at most  $\tau cr$  by Markov inequality
- For each bin use (a power of) Gaussian LSH using any point as a center (do not care about points that are more than  $\tau cr + r$  from the center)

Distance	$r$	$cr$
Collision probability	$n^{(-1+2/\tau^2)(1-\Omega_\tau(1))}/c^2$	$n^{-1+2/\tau^2}$

- Due to the independence the collision probabilities for one scale multiply (kind of), as a result obtain  $\rho \leq (1 - \Omega_\tau(1))/c^2$

Why is this family data-dependent?

- Use a point from  $P$  as a center for an inner hash table
- If  $q \in X$  is far from the center of the outer bin, then we can not handle it (but we do not care about this case)

# Smallest enclosing balls

- Jung's theorem: any set of *diameter*  $D$  lies in a ball of *radius*  $D/\sqrt{2}$
- For each bin find a smallest enclosing ball and hash wrt its center

# Smallest enclosing balls

- Jung's theorem: any set of *diameter*  $D$  lies in a ball of *radius*  $D/\sqrt{2}$
- For each bin find a smallest enclosing ball and hash wrt its center
- Careful analysis leads to

$$\rho \leq \frac{7}{8c^2} + O\left(\frac{1}{c^3}\right) + o_c(1)$$

- Can embed  $\ell_1$  into  $\ell_2$ -squared, which gives an algorithm with

$$\rho \leq \frac{7}{8c} + O\left(\frac{1}{c^{3/2}}\right) + o_c(1)$$

for  $\ell_1$  (in particular, Hamming distance for binary strings)

- Can embed  $\ell_1$  into  $\ell_2$ -squared, which gives an algorithm with

$$\rho \leq \frac{7}{8c} + O\left(\frac{1}{c^{3/2}}\right) + o_c(1)$$

for  $\ell_1$  (in particular, Hamming distance for binary strings)

- Instead of two-level hashing can consider many levels; preliminary computations give

$$\rho \leq \frac{1}{2c^2 \ln 2} + O\left(\frac{1}{c^3}\right) + o_c(1)$$

for the Euclidean case (and the similar result for  $\ell_1$  and Hamming)

- Can embed  $\ell_1$  into  $\ell_2$ -squared, which gives an algorithm with

$$\rho \leq \frac{7}{8c} + O\left(\frac{1}{c^{3/2}}\right) + o_c(1)$$

for  $\ell_1$  (in particular, Hamming distance for binary strings)

- Instead of two-level hashing can consider many levels; preliminary computations give

$$\rho \leq \frac{1}{2c^2 \ln 2} + O\left(\frac{1}{c^3}\right) + o_c(1)$$

for the Euclidean case (and the similar result for  $\ell_1$  and Hamming)

- Using this multilevel partitioning can improve known constructions for spanners for subsets of  $\ell_1$  and  $\ell_2$   
(upon [\(Har-Peled, Indyk, Sidiropoulos 2013\)](#))



# Conclusions and open problems

- Able to overcome the LSH barrier for the case of  $\ell_1$  and  $\ell_2$  using data-dependent hashing
- Can one improve our bounds?

# Conclusions and open problems

- Able to overcome the LSH barrier for the case of  $\ell_1$  and  $\ell_2$  using data-dependent hashing
- Can one improve our bounds?
- For a certain random instance can achieve  $1/(2c)$  and  $1/(2c^2)$ , which is tight for the data-dependent hashing by (Motwani, Naor, Panigrahy 2007)

# Conclusions and open problems

- Able to overcome the LSH barrier for the case of  $\ell_1$  and  $\ell_2$  using data-dependent hashing
- Can one improve our bounds?
- For a certain random instance can achieve  $1/(2c)$  and  $1/(2c^2)$ , which is tight for the data-dependent hashing by [\(Motwani, Naor, Panigrahy 2007\)](#)
- Can one get these exponents for the general case?

# Conclusions and open problems

- Able to overcome the LSH barrier for the case of  $\ell_1$  and  $\ell_2$  using data-dependent hashing
- Can one improve our bounds?
- For a certain random instance can achieve  $1/(2c)$  and  $1/(2c^2)$ , which is tight for the data-dependent hashing by [\(Motwani, Naor, Panigrahy 2007\)](#)
- Can one get these exponents for the general case?
- Can one improve the bound for this random instance further? (Looks hard!)