

Proximity in the Age of ~~Destruction~~ Distraction: Robust Approximate Nearest Neighbor Search

Sariel Har-Peled¹ Sepideh Mahabadi²

¹University of Illinois, Shampoo-Banana

²MIT

2: k robust nearest neighbor

ANN = Approximate nearest neighbor

$\mathbf{P} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$: Points in high dimensions.

k : Parameter

Q: Given query \mathbf{q} , find ANN in \mathbf{P} , when ignoring k largest coordinates.

3: Robust distance

But without the noisy coordinates

Definition

$$\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$$

$$\pi = \text{sort}(\mathbf{x}): |x_{\pi(1)}| \geq |x_{\pi(2)}| \geq \dots \geq |x_{\pi(d)}|.$$

$$i\text{-tail of } \mathbf{x}: \mathbf{x}_{\setminus i} = (\mathbf{0}, \dots, \mathbf{0}, |x_{\pi(i+1)}|, |x_{\pi(i+2)}|, \dots, |x_{\pi(d)}|).$$

computing $\mathbf{x}_{\setminus i}$ takes $O(d)$ time.

$$\mathbf{x}, \mathbf{y} \in \mathbf{P}: \text{dist}_{\setminus k}(\mathbf{x}, \mathbf{y}) = \|(\mathbf{x} - \mathbf{y})_{\setminus k}\|$$

Definition

k -robust nearest-neighbor to $\mathbf{q} \in \mathbb{R}^d$ in \mathbf{P} :

$$\text{nn}_{\setminus k}(\mathbf{P}, \mathbf{q}) = \arg \min_{\mathbf{y} \in \mathbf{P}} \text{dist}_{\setminus k}(\mathbf{q}, \mathbf{y})$$

3: Robust distance

But without the noisy coordinates

Definition

$$\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$$

$$\pi = \text{sort}(\mathbf{x}): |x_{\pi(1)}| \geq |x_{\pi(2)}| \geq \dots \geq |x_{\pi(d)}|.$$

$$i\text{-tail of } \mathbf{x}: \mathbf{x}_{\setminus i} = (\mathbf{0}, \dots, \mathbf{0}, |x_{\pi(i+1)}|, |x_{\pi(i+2)}|, \dots, |x_{\pi(d)}|).$$

computing $\mathbf{x}_{\setminus i}$ takes $O(d)$ time.

$$\mathbf{x}, \mathbf{y} \in \mathbf{P}: \text{dist}_{\setminus k}(\mathbf{x}, \mathbf{y}) = \|(\mathbf{x} - \mathbf{y})_{\setminus k}\|$$

Definition

k -robust nearest-neighbor to $\mathbf{q} \in \mathbb{R}^d$ in \mathbf{P} :

$$\text{nn}_{\setminus k}(\mathbf{P}, \mathbf{q}) = \arg \min_{\mathbf{y} \in \mathbf{P}} \text{dist}_{\setminus k}(\mathbf{q}, \mathbf{y})$$

3: Robust distance

But without the noisy coordinates

Definition

$$\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$$

$$\pi = \text{sort}(\mathbf{x}): |x_{\pi(1)}| \geq |x_{\pi(2)}| \geq \dots \geq |x_{\pi(d)}|.$$

$$i\text{-tail of } \mathbf{x}: \mathbf{x}_{\setminus i} = (\mathbf{0}, \dots, \mathbf{0}, |x_{\pi(i+1)}|, |x_{\pi(i+2)}|, \dots, |x_{\pi(d)}|).$$

computing $\mathbf{x}_{\setminus i}$ takes $O(d)$ time.

$$\mathbf{x}, \mathbf{y} \in \mathbf{P}: \text{dist}_{\setminus k}(\mathbf{x}, \mathbf{y}) = \|(\mathbf{x} - \mathbf{y})_{\setminus k}\|$$

Definition

k -robust nearest-neighbor to $\mathbf{q} \in \mathbb{R}^d$ in \mathbf{P} :

$$\text{nn}_{\setminus k}(\mathbf{P}, \mathbf{q}) = \arg \min_{\mathbf{y} \in \mathbf{P}} \text{dist}_{\setminus k}(\mathbf{q}, \mathbf{y})$$

4: Why robust NN?

- ① High dim data common.
- ② Some coordinates of \mathbf{x} and \mathbf{v} are irrelevant features.
- ③ Ignore noisy coordinates.

4: Why robust NN?

- ① High dim data common.
- ② Some coordinates of \mathbf{x} and \mathbf{v} are irrelevant features.
- ③ Ignore noisy coordinates.

4: Why robust NN?

- ① High dim data common.
- ② Some coordinates of \mathbf{x} and \mathbf{v} are irrelevant features.
- ③ Ignore noisy coordinates.

5: Results I

Reduction from Robust ANN to ANN.

- ① General reduction from the robust ANN problem to the “standard” ANN problem.
- ② L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
 - ① $\delta \in (0, 1)$: parameter
 - ② return \mathbf{x} s.t. $\text{dist}_{\|\cdot\|_{O(k/\delta)}}(\mathbf{q}, \mathbf{x}) \leq O(r/\delta)$
 - ③ $O(n^\delta)$ ANN queries in 2-ANN DS.
- ③ L_ρ : $\exists \text{NN}$ with $\text{dist} \leq r$
 - ① report point in distance $O(r(c + \frac{1}{\delta})^{1/\rho})$
 - ② ignoring $O(k(\frac{1}{\delta} + c))$ coordinates
 - ③ n^δ of $c^{1/\rho}$ -ANN queries.

5: Results I

Reduction from Robust ANN to ANN.

- 1 General reduction from the robust ANN problem to the “standard” ANN problem.
- 2 L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
 - 1 $\delta \in (0, 1)$: parameter
 - 2 return \mathbf{x} s.t. $\text{dist}_{\|_{O(k/\delta)}}(\mathbf{q}, \mathbf{x}) \leq O(r/\delta)$
 - 3 $O(n^\delta)$ ANN queries in 2-ANN DS.
- 3 L_ρ : $\exists \text{NN}$ with $\text{dist} \leq r$
 - 1 report point in distance $O(r(c + \frac{1}{\delta})^{1/\rho})$
 - 2 ignoring $O(k(\frac{1}{\delta} + c))$ coordinates
 - 3 n^δ of $c^{1/\rho}$ -ANN queries.

5: Results I

Reduction from Robust ANN to ANN.

- 1 General reduction from the robust ANN problem to the “standard” ANN problem.
- 2 L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
 - 1 $\delta \in (0, 1)$: parameter
 - 2 return \mathbf{x} s.t. $\text{dist}_{\|\cdot\|_{O(k/\delta)}}(\mathbf{q}, \mathbf{x}) \leq O(r/\delta)$
 - 3 $O(n^\delta)$ ANN queries in 2-ANN DS.
- 3 L_ρ : $\exists \text{NN}$ with $\text{dist} \leq r$
 - 1 report point in distance $O(r(c + \frac{1}{\delta})^{1/\rho})$
 - 2 ignoring $O(k(\frac{1}{\delta} + c))$ coordinates
 - 3 n^δ of $c^{1/\rho}$ -ANN queries.

5: Results I

Reduction from Robust ANN to ANN.

- ① General reduction from the robust ANN problem to the “standard” ANN problem.
- ② L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
 - ① $\delta \in (0, 1)$: parameter
 - ② return \mathbf{x} s.t. $\text{dist}_{\|_{O(k/\delta)}}(\mathbf{q}, \mathbf{x}) \leq O(r/\delta)$
 - ③ $O(n^\delta)$ ANN queries in 2-ANN DS.
- ③ L_ρ : $\exists \text{NN}$ with $\text{dist} \leq r$
 - ① report point in distance $O(r(c + \frac{1}{\delta}))^{1/\rho}$
 - ② ignoring $O(k(\frac{1}{\delta} + c))$ coordinates
 - ③ n^δ of $c^{1/\rho}$ -ANN queries.

5: Results I

Reduction from Robust ANN to ANN.

- ① General reduction from the robust ANN problem to the “standard” ANN problem.
- ② L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
 - ① $\delta \in (0, 1)$: parameter
 - ② return \mathbf{x} s.t. $\text{dist}_{\|_{O(k/\delta)}}(\mathbf{q}, \mathbf{x}) \leq O(r/\delta)$
 - ③ $O(n^\delta)$ ANN queries in 2-ANN DS.
- ③ L_ρ : $\exists \text{NN}$ with $\text{dist} \leq r$
 - ① report point in distance $O(r(c + \frac{1}{\delta})^{1/\rho})$
 - ② ignoring $O(k(\frac{1}{\delta} + c))$ coordinates
 - ③ n^δ of $c^{1/\rho}$ -ANN queries.

5: Results I

Reduction from Robust ANN to ANN.

- ① General reduction from the robust ANN problem to the “standard” ANN problem.
- ② L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
 - ① $\delta \in (0, 1)$: parameter
 - ② return \mathbf{x} s.t. $\text{dist}_{\|_{O(k/\delta)}}(\mathbf{q}, \mathbf{x}) \leq O(r/\delta)$
 - ③ $O(n^\delta)$ ANN queries in 2-ANN DS.
- ③ L_ρ : $\exists \text{NN}$ with $\text{dist} \leq r$
 - ① report point in distance $O(r(c + \frac{1}{\delta})^{1/\rho})$
 - ② ignoring $O(k(\frac{1}{\delta} + c))$ coordinates
 - ③ n^δ of $c^{1/\rho}$ -ANN queries.

6: Results II

$(1 + \varepsilon)$ -approximation

L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$

- 1 $\varepsilon, \delta \in (0, 1)$: parameters
- 2 return \mathbf{x} s.t. $\text{dist}_{\|\cdot\|_{O(\frac{k}{\delta\varepsilon})}}(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$
- 3 $\tilde{O}(n^\delta/\varepsilon)$ ANN queries in $(1 + O(\varepsilon))$ -ANN DS.

6: Results II

$(1 + \epsilon)$ -approximation

L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$

- 1 $\epsilon, \delta \in (0, 1)$: parameters
- 2 return \mathbf{x} s.t. $\text{dist}_{\|\cdot\|_{O(\frac{k}{\delta\epsilon})}}(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon)r$
- 3 $\tilde{O}(n^\delta/\epsilon)$ ANN queries in $(1 + O(\epsilon))$ -ANN DS.

6: Results II

$(1 + \varepsilon)$ -approximation

L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$

- 1 $\varepsilon, \delta \in (0, 1)$: parameters
- 2 return \mathbf{x} s.t. $\text{dist}_{\|\cdot\|_{O(\frac{k}{\delta\varepsilon})}}(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$
- 3 $\tilde{O}(n^\delta/\varepsilon)$ ANN queries in $(1 + O(\varepsilon))$ -ANN DS.

6: Results II

$(1 + \varepsilon)$ -approximation

L_1 : $\exists \mathbf{q}^* \in \mathbf{P}$ s.t. $\text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$

- 1 $\varepsilon, \delta \in (0, 1)$: parameters
- 2 return \mathbf{x} s.t. $\text{dist}_{\|\cdot\|_{O(\frac{k}{\delta\varepsilon})}}(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$
- 3 $\tilde{O}(n^\delta/\varepsilon)$ ANN queries in $(1 + O(\varepsilon))$ -ANN DS.

7: Results III

Budgeted version

- 1 Every coordinate has a cost of ignoring it.
- 2 $\exists \mathbf{q}^* \in \mathbf{P}$, s.t., $\text{dist}_{\|\cdot\|_1}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
- 3 Algorithm returns a point in distance $O(r)$
- 4 Cost of ignored coordinates is $O(1)$.
- 5 in L_1 norm.
- 6 Computing exact distance is NP-Hard even for two points.

7: Results III

Budgeted version

- 1 Every coordinate has a cost of ignoring it.
- 2 $\exists \mathbf{q}^* \in \mathbf{P}$, s.t., $\text{dist}_{\|\cdot\|_1}(\mathbf{q}, \mathbf{q}^*) \leq r \implies$
- 3 Algorithm returns a point in distance $O(r)$
- 4 Cost of ignored coordinates is $O(1)$.
- 5 in L_1 norm.
- 6 Computing exact distance is NP-Hard even for two points.

8: Results III

Data sensitive LSH queries

① known ANN DS perform better in practice.

① Many cases return *exact* NN!

② Find it early!

③ Spend more time “proving” its ANN.

② **Q:** If query is easy, why work hard?

③ **Result:** (for $\{0, 1\}^d$).

① query time: $O(d \exp(\Delta) \log n)$

② Δ is smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\Delta \frac{\text{dist}(\mathbf{q}, \mathbf{x}_i)}{r}\right) \leq 1,$$

③ works quickly on low dimensional data.

8: Results III

Data sensitive LSH queries

① known ANN DS perform better in practice.

① Many cases return *exact* NN!

② Find it early!

③ Spend more time “proving” its ANN.

② **Q:** If query is easy, why work hard?

③ **Result:** (for $\{0, 1\}^d$).

① query time: $O(d \exp(\Delta) \log n)$

② Δ is smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\Delta \frac{\text{dist}(\mathbf{q}, \mathbf{x}_i)}{r}\right) \leq 1,$$

③ works quickly on low dimensional data.

8: Results III

Data sensitive LSH queries

- 1 known ANN DS perform better in practice.
 - 1 Many cases return *exact* NN!
 - 2 Find it early!
 - 3 Spend more time “proving” its ANN.
- 2 **Q:** If query is easy, why work hard?
- 3 **Result:** (for $\{0, 1\}^d$).
 - 1 query time: $O(d \exp(\Delta) \log n)$
 - 2 Δ is smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\Delta \frac{\text{dist}(\mathbf{q}, \mathbf{x}_i)}{r}\right) \leq 1,$$

- 3 works quickly on low dimensional data.

8: Results III

Data sensitive LSH queries

① known ANN DS perform better in practice.

- ① Many cases return *exact* NN!
- ② Find it early!
- ③ Spend more time “proving” its ANN.

② **Q:** If query is easy, why work hard?

③ **Result:** (for $\{0, 1\}^d$).

- ① query time: $O(d \exp(\Delta) \log n)$
- ② Δ is smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\Delta \frac{\text{dist}(\mathbf{q}, \mathbf{x}_i)}{r}\right) \leq 1,$$

- ③ works quickly on low dimensional data.

8: Results III

Data sensitive LSH queries

① known ANN DS perform better in practice.

- ① Many cases return *exact* NN!
- ② Find it early!
- ③ Spend more time “proving” its ANN.

② **Q:** If query is easy, why work hard?

③ **Result:** (for $\{0, 1\}^d$).

- ① query time: $O(d \exp(\Delta) \log n)$
- ② Δ is smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\Delta \frac{\text{dist}(\mathbf{q}, \mathbf{x}_i)}{r}\right) \leq 1,$$

- ③ works quickly on low dimensional data.

8: Results III

Data sensitive LSH queries

- ① known ANN DS perform better in practice.
 - ① Many cases return *exact* NN!
 - ② Find it early!
 - ③ Spend more time “proving” its ANN.
- ② **Q:** If query is easy, why work hard?
- ③ **Result:** (for $\{0, 1\}^d$).
 - ① query time: $O(d \exp(\Delta) \log n)$
 - ② Δ is smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\Delta \frac{\text{dist}(q, x_i)}{r}\right) \leq 1,$$

- ③ works quickly on low dimensional data.

8: Results III

Data sensitive LSH queries

- ① known ANN DS perform better in practice.
 - ① Many cases return *exact* NN!
 - ② Find it early!
 - ③ Spend more time “proving” its ANN.
- ② **Q:** If query is easy, why work hard?
- ③ **Result:** (for $\{0, 1\}^d$).
 - ① query time: $O(d \exp(\Delta) \log n)$
 - ② Δ is smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\Delta \frac{\text{dist}(\mathbf{q}, \mathbf{x}_i)}{r}\right) \leq 1,$$

- ③ works quickly on low dimensional data.

Part I

Locality sensitive hashing revisited

9: Locality sensitive hashing (LSH)

- ① Technique due to Indyk and Motwani [Indyk and Motwani, 1998].
- ② Input: $\mathbf{P} \subseteq \{0, 1\}^d$: set of n points.
- ③ Parameters r and ϵ .
- ④ **Q:** Given query point $\mathbf{q} \in \{0, 1\}^d$:
 - ① If $d_H(\mathbf{q}, \mathbf{P}) \leq r$ then return $\mathbf{v} \in \mathbf{P}$ s.t. $d_H(\mathbf{q}, \mathbf{v}) \leq (1 + \epsilon)r$.
 - ② If $d_H(\mathbf{q}, \mathbf{P}) > (1 + \epsilon)r$ then return " $d_H(\mathbf{q}, \mathbf{P}) > r$ ".
 - ③ If $d_H(\mathbf{q}, \mathbf{P}) \in (r, (1 + \epsilon)r)$ return either of the above.

10: A random projection...

- 1 Sequence $m \equiv i_1, i_2, \dots, i_\ell \in \llbracket d \rrbracket = \{1, \dots, d\}$.
- 2 $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$
- 3 projection by m is

$$m\mathbf{x} = (x_{i_1}, \dots, x_{i_\ell}) \in \mathbb{R}^\ell.$$

- 4 $m\mathbf{P} = \{m\mathbf{x} \mid \mathbf{x} \in \mathbf{P}\}$.

10: A random projection...

- 1 Sequence $m \equiv i_1, i_2, \dots, i_\ell \in \llbracket d \rrbracket = \{1, \dots, d\}$.
- 2 $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$
- 3 projection by m is

$$m\mathbf{x} = (x_{i_1}, \dots, x_{i_\ell}) \in \mathbb{R}^\ell.$$

- 4 $m\mathbf{P} = \{m\mathbf{x} \mid \mathbf{x} \in \mathbf{P}\}$.

10: A random projection...

- ① Sequence $m \equiv i_1, i_2, \dots, i_\ell \in \llbracket d \rrbracket = \{1, \dots, d\}$.
- ② $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$
- ③ **projection** by m is

$$m\mathbf{x} = (x_{i_1}, \dots, x_{i_\ell}) \in \mathbb{R}^\ell.$$

④ $mP = \{m\mathbf{x} \mid \mathbf{x} \in P\}$.

10: A random projection...

- ① Sequence $m \equiv i_1, i_2, \dots, i_\ell \in \llbracket d \rrbracket = \{1, \dots, d\}$.
- ② $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$
- ③ **projection** by m is

$$m\mathbf{x} = (x_{i_1}, \dots, x_{i_\ell}) \in \mathbb{R}^\ell.$$

- ④ $m\mathbf{P} = \{m\mathbf{x} \mid \mathbf{x} \in \mathbf{P}\}$.

11: A splay

- 1 τ : a probability
- 2 include i th coordinate, for $i = 1, \dots, d$, in m with probability τ
- 3 \mathcal{D}_τ : Distribution of such projections.
- 4 $m = i_1, \dots, i_\ell$ and $n = j_1, \dots, j_{\ell'}$:
 $m|n = i_1, \dots, i_\ell, j_1, \dots, j_{\ell'}$.
- 5 \mathcal{D}_τ^t : distribution of sequences concatenating t sequences of \mathcal{D}_τ .

11: A splay

- 1 τ : a probability
- 2 include i th coordinate, for $i = 1, \dots, d$, in m with probability τ
- 3 \mathcal{D}_τ : Distribution of such projections.
- 4 $m = i_1, \dots, i_\ell$ and $n = j_1, \dots, j_{\ell'}$:
 $m|n = i_1, \dots, i_\ell, j_1, \dots, j_{\ell'}$.
- 5 \mathcal{D}_τ^t : distribution of sequences concatenating t sequences of \mathcal{D}_τ .

11: A splay

- ① τ : a probability
- ② include i th coordinate, for $i = 1, \dots, d$, in m with probability τ
- ③ \mathcal{D}_τ : Distribution of such projections.
- ④ $m = i_1, \dots, i_\ell$ and $n = j_1, \dots, j_{\ell'}$:
 $m|n = i_1, \dots, i_\ell, j_1, \dots, j_{\ell'}$.
- ⑤ \mathcal{D}_τ^t : distribution of sequences concatenating t sequences of \mathcal{D}_τ .

11: A splay

- ① τ : a probability
- ② include i th coordinate, for $i = 1, \dots, d$, in m with probability τ
- ③ \mathcal{D}_τ : Distribution of such projections.
- ④ $m = i_1, \dots, i_\ell$ and $n = j_1, \dots, j_{\ell'}$:
 $m|n = i_1, \dots, i_\ell, j_1, \dots, j_{\ell'}$.
- ⑤ \mathcal{D}_τ^t : distribution of sequences concatenating t sequences of \mathcal{D}_τ .

12: Preprocessing

- 1 Set $\beta = 1/(1 + \epsilon) \in (0, 1)$.
- 2 $\tau = 1/r$, $t = \beta \log n$, and $L = O(n^\beta \log n)$.
- 3 Pick randomly L sequences $M_1, \dots, M_L \in \mathcal{D}_\tau^t$.
- 4 Algorithm computes $Q_i = M_i P_i$, for $i = 1, \dots, L$.
- 5 Store Q_i in a hash table \mathcal{H}_i for $i = 1, \dots, L$.

- 6 **Answering a query.**
 - 1 Given $q \in \{0, 1\}^d$: compute $q_i = M_i q$, for $i = 1, \dots, L$.
 - 2 Retrieve list ℓ_i from \mathcal{H}_i , of all points that collide with q_i .
 - 3 Scan ℓ_1, \dots, ℓ_L .
If any point is in $\text{dist} \leq (1 + \epsilon)r$, return it.
 - 4 Otherwise, return " $d_H(q, P) > r$ ".

12: Preprocessing

- 1 Set $\beta = 1/(1 + \epsilon) \in (0, 1)$.
- 2 $\tau = 1/r$, $t = \beta \log n$, and $L = O(n^\beta \log n)$.
- 3 Pick randomly L sequences $M_1, \dots, M_L \in \mathcal{D}_\tau^t$.
- 4 Algorithm computes $Q_i = M_i P_i$, for $i = 1, \dots, L$.
- 5 Store Q_i in a hash table \mathcal{H}_i for $i = 1, \dots, L$.

- 6 **Answering a query.**
 - 1 Given $q \in \{0, 1\}^d$: compute $q_i = M_i q$, for $i = 1, \dots, L$.
 - 2 Retrieve list ℓ_i from \mathcal{H}_i , of all points that collide with q_i .
 - 3 Scan ℓ_1, \dots, ℓ_L .
If any point is in $\text{dist} \leq (1 + \epsilon)r$, return it.
 - 4 Otherwise, return " $d_H(q, P) > r$ ".

12: Preprocessing

- 1 Set $\beta = 1/(1 + \varepsilon) \in (0, 1)$.
- 2 $\tau = 1/r$, $t = \beta \log n$, and $L = O(n^\beta \log n)$.
- 3 Pick randomly L sequences $M_1, \dots, M_L \in \mathcal{D}_\tau^t$.
- 4 Algorithm computes $Q_i = M_i P_i$, for $i = 1, \dots, L$.
- 5 Store Q_i in a hash table \mathcal{H}_i for $i = 1, \dots, L$.

- 6 **Answering a query.**
 - 1 Given $\mathbf{q} \in \{0, 1\}^d$: compute $\mathbf{q}_i = M_i \mathbf{q}$, for $i = 1, \dots, L$.
 - 2 Retrieve list ℓ_i from \mathcal{H}_i , of all points that collide with \mathbf{q}_i .
 - 3 Scan ℓ_1, \dots, ℓ_L .
If any point is in $\text{dist} \leq (1 + \varepsilon)r$, return it.
 - 4 Otherwise, return " $d_H(\mathbf{q}, \mathbf{P}) > r$ ".

13: What is going on?

Collision probabilities

\mathbf{q}, \mathbf{y} : Two points.

$$\epsilon = 1 \quad x = d_H(\mathbf{q}, \mathbf{y}). \quad t = \frac{\log n}{1+\epsilon}$$

$$f(x) = \Pr[M\mathbf{q} = M\mathbf{x}] = \left(\left(1 - \frac{1}{r}\right)^x \right)^t \approx \exp\left(-\frac{xt}{r}\right)$$

1

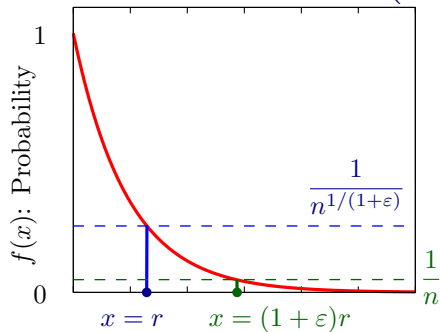
13: What is going on?

Collision probabilities

\mathbf{q}, \mathbf{y} : Two points.

$$\epsilon = 1 \quad x = d_H(\mathbf{q}, \mathbf{y}). \quad t = \frac{\log n}{1+\epsilon}$$

$$f(x) = \Pr[M\mathbf{q} = M\mathbf{x}] = \left(1 - \frac{1}{r}\right)^x \approx \exp\left(-\frac{xt}{r}\right)$$



1

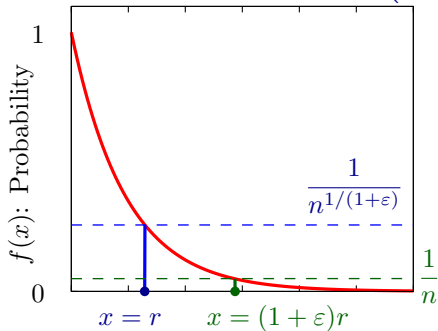
13: What is going on?

Collision probabilities

\mathbf{q}, \mathbf{y} : Two points.

$$\epsilon = 1 \quad x = d_H(\mathbf{q}, \mathbf{y}). \quad t = \frac{\log n}{1+\epsilon}$$

$$f(x) = \Pr[M\mathbf{q} = M\mathbf{x}] = \left(1 - \frac{1}{r}\right)^x \approx \exp\left(-\frac{xt}{r}\right)$$



$$\textcircled{1} f(x) = \exp\left(-\frac{x \log n}{(1+\epsilon)r}\right)$$

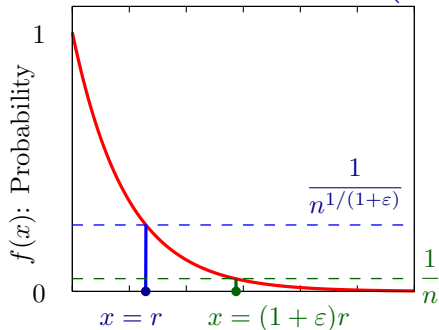
13: What is going on?

Collision probabilities

\mathbf{q}, \mathbf{y} : Two points.

$$\epsilon = 1 \quad x = d_H(\mathbf{q}, \mathbf{y}). \quad t = \frac{\log n}{1+\epsilon}$$

$$f(x) = \Pr[M\mathbf{q} = M\mathbf{x}] = \left(1 - \frac{1}{r}\right)^x \approx \exp\left(-\frac{xt}{r}\right)$$



① $f(x) = \exp\left(-\frac{x \log n}{(1+\epsilon)r}\right)$

② $f(r) = 1/n^{1/(1+\epsilon)}$

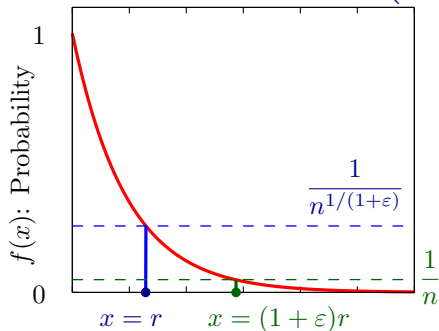
13: What is going on?

Collision probabilities

\mathbf{q}, \mathbf{y} : Two points.

$$\epsilon = 1 \quad x = d_H(\mathbf{q}, \mathbf{y}). \quad t = \frac{\log n}{1+\epsilon}$$

$$f(x) = \Pr[M\mathbf{q} = M\mathbf{x}] = \left(1 - \frac{1}{r}\right)^x \approx \exp\left(-\frac{xt}{r}\right)$$



- 1 $f(x) = \exp\left(-\frac{x \log n}{(1+\epsilon)r}\right)$
- 2 $f(r) = 1/n^{1/(1+\epsilon)}$
- 3 $f((1 + \epsilon)r) = 1/n$

14: Correctness – avoiding few bad coordinates

Lemma

K : set of r forbidden coordinates. Probability

$M = (m_1, \dots, m_t) \in \mathcal{D}_\tau^t$ does not sample any coordinate of K is $\geq n^{-\beta(1+1/(2r))} \approx 1/n^\beta = 1/n^{1/(1+\epsilon)}$.

Proof.

Probability m_i not contain any of coordinates of K is $(1 - \tau)^r$.

Desired probability is

$$\begin{aligned} \left(1 - \frac{1}{r}\right)^{rt} &= \left(1 - \frac{1}{r}\right)^{r\beta \log n} \geq \exp(-\beta \log n) \left(1 - \frac{1}{r}\right)^{\beta \log n} \\ &\geq n^{-\beta} \exp\left(-\frac{1}{2r}\beta \log n\right) = n^{-\beta - \beta/(2r)}, \end{aligned}$$

by $(1 - 1/m)^{m-1} \geq \frac{1}{e} \geq (1 - 1/m)^m$, $e^{-2x} \leq 1 - x \leq e^x$. □

15: Proof

Lemma

- (A) If $d_H(\mathbf{q}, \mathbf{P}) \leq r$. W.h.p. return $\mathbf{x} \in \mathbf{P}$ s.t.
 $d_H(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$.
- (B) If $d_H(\mathbf{q}, \mathbf{P}) \geq (1 + \varepsilon)r$. Then, $\mathbf{E}[\sum_i |\ell_i|] = O(L)$.

Proof.

- (A) By previous lemma. Probability $M_i \mathbf{x} = M_i \mathbf{q}$ is at least $1/n^\beta$. Probability all L data-structures fails, is $\leq (1 - 1/n^\beta)^L < 1/n^{O(1)}$, as $L = O(n^\beta \log n)$.
- (B) For $\mathbf{v} \in \mathbf{P}$, s.t. $d_H(\mathbf{v}, \mathbf{q}) \geq (1 + \varepsilon)r$. Probability for $M \in \mathcal{D}_\tau^t$ misses all the $\geq (1 + \varepsilon)r$ coordinates (i.e., collision!) is $\leq (1 - \tau)^{(1+\varepsilon)rt} \leq \exp\left(- (1 + \varepsilon) \frac{tr}{r}\right) = \exp(-(1 + \varepsilon)\beta \log n) = 1/n$, as $\beta = 1/(1 + \varepsilon)$. Expected size of ℓ_i is $\leq n \times (1/n) = 1$.

15: Proof

Lemma

- (A) If $d_H(\mathbf{q}, \mathbf{P}) \leq r$. W.h.p. return $\mathbf{x} \in \mathbf{P}$ s.t.
 $d_H(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$.
- (B) If $d_H(\mathbf{q}, \mathbf{P}) \geq (1 + \varepsilon)r$. Then, $\mathbf{E}[\sum_i |\ell_i|] = O(L)$.

Proof.

- (A) By previous lemma. Probability $M_i \mathbf{x} = M_i \mathbf{q}$ is at least $1/n^\beta$. Probability all L data-structures fails, is $\leq (1 - 1/n^\beta)^L < 1/n^{O(1)}$, as $L = O(n^\beta \log n)$.
- (B) For $\mathbf{v} \in \mathbf{P}$, s.t. $d_H(\mathbf{v}, \mathbf{q}) \geq (1 + \varepsilon)r$. Probability for $M \in \mathcal{D}_\tau^t$ misses all the $\geq (1 + \varepsilon)r$ coordinates (i.e., collision!) is $\leq (1 - \tau)^{(1+\varepsilon)rt} \leq \exp\left(- (1 + \varepsilon) \frac{tr}{r}\right) = \exp(-(1 + \varepsilon)\beta \log n) = 1/n$, as $\beta = 1/(1 + \varepsilon)$. Expected size of ℓ_i is $\leq n \times (1/n) = 1$.

15: Proof

Lemma

- (A) If $d_H(\mathbf{q}, \mathbf{P}) \leq r$. W.h.p. return $\mathbf{x} \in \mathbf{P}$ s.t.
 $d_H(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$.
- (B) If $d_H(\mathbf{q}, \mathbf{P}) \geq (1 + \varepsilon)r$. Then, $\mathbf{E}[\sum_i |\ell_i|] = O(L)$.

Proof.

- (A) By previous lemma. Probability $M_i \mathbf{x} = M_i \mathbf{q}$ is at least $1/n^\beta$. Probability all L data-structures fails, is $\leq (1 - 1/n^\beta)^L < 1/n^{O(1)}$, as $L = O(n^\beta \log n)$.
- (B) For $\mathbf{v} \in \mathbf{P}$, s.t. $d_H(\mathbf{v}, \mathbf{q}) \geq (1 + \varepsilon)r$. Probability for $M \in \mathcal{D}_\tau^t$ misses all the $\geq (1 + \varepsilon)r$ coordinates (i.e., collision!) is $\leq (1 - \tau)^{(1+\varepsilon)rt} \leq \exp\left(- (1 + \varepsilon) \frac{tr}{r}\right) = \exp(- (1 + \varepsilon)\beta \log n) = 1/n$, as $\beta = 1/(1 + \varepsilon)$. Expected size of ℓ_i is $\leq n \times (1/n) = 1$.

15: Proof

Lemma

- (A) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \leq r$. W.h.p. return $\mathbf{x} \in \mathbf{P}$ s.t.
 $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$.
- (B) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \geq (1 + \varepsilon)r$. Then, $\mathbf{E}[\sum_i |\ell_i|] = O(L)$.

Proof.

- (A) By previous lemma. Probability $M_i \mathbf{x} = M_i \mathbf{q}$ is at least $1/n^\beta$. Probability all L data-structures fails, is $\leq (1 - 1/n^\beta)^L < 1/n^{O(1)}$, as $L = O(n^\beta \log n)$.
- (B) For $\mathbf{v} \in \mathbf{P}$, s.t. $\mathbf{d}_H(\mathbf{v}, \mathbf{q}) \geq (1 + \varepsilon)r$. Probability for $M \in \mathcal{D}_\tau^t$ misses all the $\geq (1 + \varepsilon)r$ coordinates (i.e., collision!) is $\leq (1 - \tau)^{(1+\varepsilon)rt} \leq \exp\left(- (1 + \varepsilon) \frac{tr}{r}\right) = \exp(-(1 + \varepsilon)\beta \log n) = 1/n$, as $\beta = 1/(1 + \varepsilon)$. Expected size of ℓ_i is $\leq n \times (1/n) = 1$.

15: Proof

Lemma

- (A) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \leq r$. W.h.p. return $\mathbf{x} \in \mathbf{P}$ s.t.
 $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$.
- (B) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \geq (1 + \varepsilon)r$. Then, $\mathbf{E}[\sum_i |\ell_i|] = O(L)$.

Proof.

- (A) By previous lemma. Probability $M_i \mathbf{x} = M_i \mathbf{q}$ is at least $1/n^\beta$. Probability all L data-structures fails, is $\leq (1 - 1/n^\beta)^L < 1/n^{O(1)}$, as $L = O(n^\beta \log n)$.
- (B) For $\mathbf{v} \in \mathbf{P}$, s.t. $\mathbf{d}_H(\mathbf{v}, \mathbf{q}) \geq (1 + \varepsilon)r$. Probability for $M \in \mathcal{D}_\tau^t$ misses all the $\geq (1 + \varepsilon)r$ coordinates (i.e., collision!) is $\leq (1 - \tau)^{(1+\varepsilon)rt} \leq \exp\left(- (1 + \varepsilon) \frac{tr}{r}\right) = \exp(-(1 + \varepsilon)\beta \log n) = 1/n$, as $\beta = 1/(1 + \varepsilon)$. Expected size of ℓ_i is $\leq n \times (1/n) = 1$.

15: Proof

Lemma

(A) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \leq r$. W.h.p. return $\mathbf{x} \in \mathbf{P}$ s.t.

$$\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r.$$

(B) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \geq (1 + \varepsilon)r$. Then, $\mathbf{E}[\sum_i |\ell_i|] = O(L)$.

Proof.

(A) By previous lemma. Probability $M_i \mathbf{x} = M_i \mathbf{q}$ is at least $1/n^\beta$.

Probability all L data-structures fails, is

$$\leq (1 - 1/n^\beta)^L < 1/n^{O(1)}, \text{ as } L = O(n^\beta \log n).$$

(B) For $\mathbf{v} \in \mathbf{P}$, s.t. $\mathbf{d}_H(\mathbf{v}, \mathbf{q}) \geq (1 + \varepsilon)r$. Probability for $M \in \mathcal{D}_\tau^t$ misses all the $\geq (1 + \varepsilon)r$ coordinates (i.e., collision!) is

$$\leq (1 - \tau)^{(1+\varepsilon)rt} \leq \exp\left(- (1 + \varepsilon) \frac{tr}{r}\right) =$$

$\exp(- (1 + \varepsilon)\beta \log n) = 1/n$, as $\beta = 1/(1 + \varepsilon)$. Expected size of ℓ_i is $\leq n \times (1/n) = 1$.

15: Proof

Lemma

- (A) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \leq r$. W.h.p. return $\mathbf{x} \in \mathbf{P}$ s.t.
 $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \varepsilon)r$.
- (B) If $\mathbf{d}_H(\mathbf{q}, \mathbf{P}) \geq (1 + \varepsilon)r$. Then, $\mathbf{E}[\sum_i |\ell_i|] = O(L)$.

Proof.

- (A) By previous lemma. Probability $M_i \mathbf{x} = M_i \mathbf{q}$ is at least $1/n^\beta$. Probability all L data-structures fails, is
 $\leq (1 - 1/n^\beta)^L < 1/n^{O(1)}$, as $L = O(n^\beta \log n)$.
- (B) For $\mathbf{v} \in \mathbf{P}$, s.t. $\mathbf{d}_H(\mathbf{v}, \mathbf{q}) \geq (1 + \varepsilon)r$. Probability for $M \in \mathcal{D}_\tau^t$ misses all the $\geq (1 + \varepsilon)r$ coordinates (i.e., collision!) is
 $\leq (1 - \tau)^{(1+\varepsilon)rt} \leq \exp\left(- (1 + \varepsilon) \frac{tr}{r}\right) =$
 $\exp(- (1 + \varepsilon)\beta \log n) = 1/n$, as $\beta = 1/(1 + \varepsilon)$. Expected size of ℓ_i is $\leq n \times (1/n) = 1$. □

16: LSH: Query running time

- 1 Compute $\forall i \quad M_i \mathbf{q}$. Takes $O(dt) = O(d \log n)$ time.
- 2 Repeat above L times: $O(Ld \log n) = O(n^{1/(1+\epsilon)} \log^2 n)$ time overall.
- 3 In expectation, after scanning $O(L)$ points, either encountered point in distance r , or all points are in distance larger than $(1 + \epsilon)r$.
- 4 Expected query time is $O(Ld + Ld \log n) = O(n^{1/(1+\epsilon)} \log^2 n)$.

16: LSH: Query running time

- 1 Compute $\forall i \ M_i; \mathbf{q}$. Takes $O(dt) = O(d \log n)$ time.
- 2 Repeat above L times: $O(Ld \log n) = O(n^{1/(1+\epsilon)} \log^2 n)$ time overall.
- 3 In expectation, after scanning $O(L)$ points, either encountered point in distance r , or all points are in distance larger than $(1 + \epsilon)r$.
- 4 Expected query time is $O(Ld + Ld \log n) = O(n^{1/(1+\epsilon)} \log^2 n)$.

Part II

Robust nearest neighbor

17: The good, the bad, and the truncated.

- ① technique probes coordinates, trying to detect “hidden” mass of a bad point.
- ② mass might concentrate in few coordinates (i.e., $k + 1$).
- ③ Such point still good?
- ④ Point one has to ignore many more than k coordinates is “bad”.
- ⑤ Bound the influence of a single coordinate on distance to \mathbf{q} .
- ⑥ $r = \text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{P})$.
- ⑦ Truncate coordinate of $\mathbf{x} - \mathbf{q}$ if larger than r/k .

17: The good, the bad, and the truncated.

- ① technique probes coordinates, trying to detect “hidden” mass of a bad point.
- ② mass might concentrate in few coordinates (i.e., $k + 1$).
- ③ Such point still good?
- ④ Point one has to ignore many more than k coordinates is “bad”.
- ⑤ Bound the influence of a single coordinate on distance to \mathbf{q} .
- ⑥ $r = \text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{P})$.
- ⑦ Truncate coordinate of $\mathbf{x} - \mathbf{q}$ if larger than r/k .

17: The good, the bad, and the truncated.

- ① technique probes coordinates, trying to detect “hidden” mass of a bad point.
- ② mass might concentrate in few coordinates (i.e., $k + 1$).
- ③ Such point still good?
- ④ Point one has to ignore many more than k coordinates is “bad”.
- ⑤ Bound the influence of a single coordinate on distance to \mathbf{q} .
- ⑥ $r = \text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{P})$.
- ⑦ Truncate coordinate of $\mathbf{x} - \mathbf{q}$ if larger than r/k .

17: The good, the bad, and the truncated.

- ① technique probes coordinates, trying to detect “hidden” mass of a bad point.
- ② mass might concentrate in few coordinates (i.e., $k + 1$).
- ③ Such point still good?
- ④ Point one has to ignore many more than k coordinates is “bad”.
- ⑤ Bound the influence of a single coordinate on distance to \mathbf{q} .
- ⑥ $r = \text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{P})$.
- ⑦ Truncate coordinate of $\mathbf{x} - \mathbf{q}$ if larger than r/k .

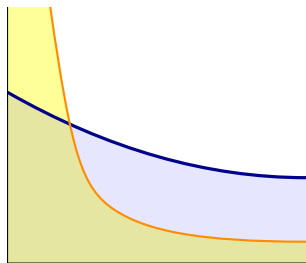
17: The good, the bad, and the truncated.

- ① technique probes coordinates, trying to detect “hidden” mass of a bad point.
- ② mass might concentrate in few coordinates (i.e., $k + 1$).
- ③ Such point still good?
- ④ Point one has to ignore many more than k coordinates is “bad”.
- ⑤ Bound the influence of a single coordinate on distance to \mathbf{q} .
- ⑥ $r = \text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{P})$.
- ⑦ Truncate coordinate of $\mathbf{x} - \mathbf{q}$ if larger than r/k .

17: The good, the bad, and the truncated.

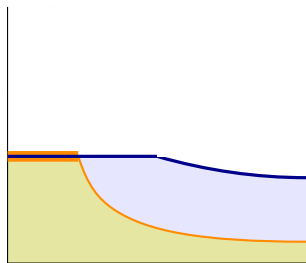
- ① technique probes coordinates, trying to detect “hidden” mass of a bad point.
- ② mass might concentrate in few coordinates (i.e., $k + 1$).
- ③ Such point still good?
- ④ Point one has to ignore many more than k coordinates is “bad”.
- ⑤ Bound the influence of a single coordinate on distance to \mathbf{q} .
- ⑥ $r = \text{dist}_{\|\cdot\|_k}(\mathbf{q}, \mathbf{P})$.
- ⑦ Truncate coordinate of $\mathbf{x} - \mathbf{q}$ if larger than r/k .

18: Truncation



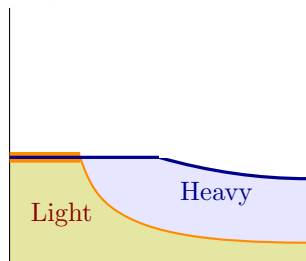
- 1 Pretend that points are truncated.
- 2 Prove good approximation (slightly refined LSH analysis).
- 3 Interpret the results for robust ANN.

18: Truncation



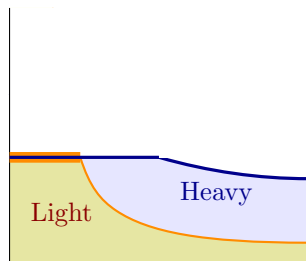
- 1 Pretend that points are truncated.
- 2 Prove good approximation (slightly refined LSH analysis).
- 3 Interpret the results for robust ANN.

18: Truncation



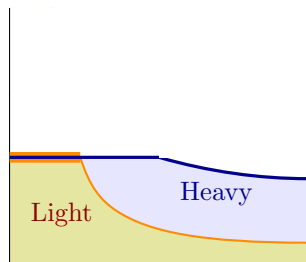
- 1 Pretend that points are truncated.
- 2 Prove good approximation (slightly refined LSH analysis).
- 3 Interpret the results for robust ANN.

18: Truncation



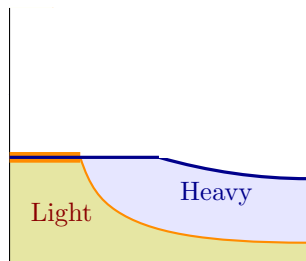
- 1 Pretend that points are truncated.
- 2 Prove good approximation (slightly refined LSH analysis).
- 3 Interpret the results for robust ANN.

18: Truncation



- 1 Pretend that points are truncated.
- 2 Prove good approximation (slightly refined LSH analysis).
- 3 Interpret the results for robust ANN.

18: Truncation



- 1 Pretend that points are truncated.
- 2 Prove good approximation (slightly refined LSH analysis).
- 3 Interpret the results for robust ANN.

19: The good, the bad, and the truncated.

- ① Algorithm is bicriterion approximation.
Returns a point where ignore slightly more coordinates than k ... but resulting point is constant approximation to NN when ignoring k coordinates.
- ② Bad point after truncation, is amenable to random probing.
- ③ Heavy points have many bad coordinates.

19: The good, the bad, and the truncated.

- ① Algorithm is bicriterion approximation.
Returns a point where ignore slightly more coordinates than k
... but resulting point is constant approximation to NN when
ignoring k coordinates.
- ② Bad point after truncation, is amenable to random probing.
- ③ Heavy points have many bad coordinates.

19: The good, the bad, and the truncated.

- ① Algorithm is bicriterion approximation.
Returns a point where ignore slightly more coordinates than k
... but resulting point is constant approximation to NN when
ignoring k coordinates.
- ② Bad point after truncation, is amenable to random probing.
- ③ Heavy points have many bad coordinates.

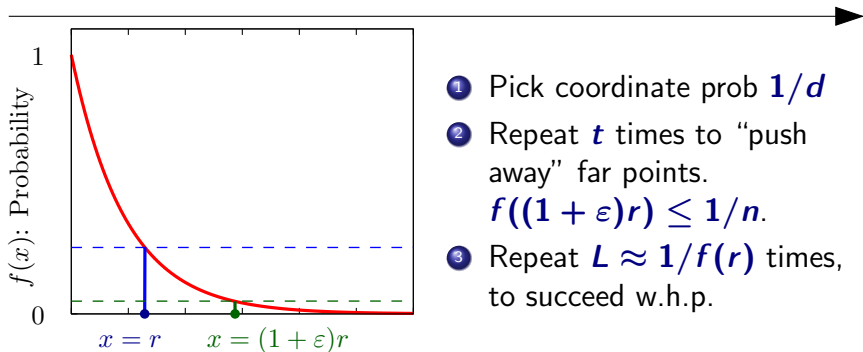
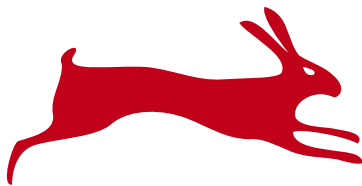
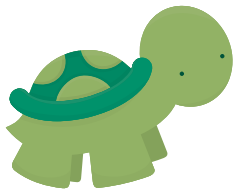
19: The good, the bad, and the truncated.

- ① Algorithm is bicriterion approximation.
Returns a point where ignore slightly more coordinates than k
... but resulting point is constant approximation to NN when
ignoring k coordinates.
- ② Bad point after truncation, is amenable to random probing.
- ③ Heavy points have many bad coordinates.

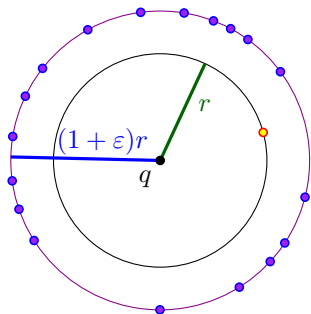
Part III

Data sensitive LSH

20: The tortoise and the hare



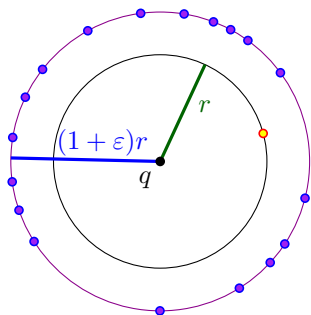
21: The worst/good case for LSH



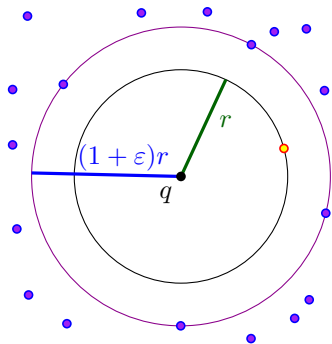
Worst case

Work hard to push away a single point...

21: The worst/good case for LSH



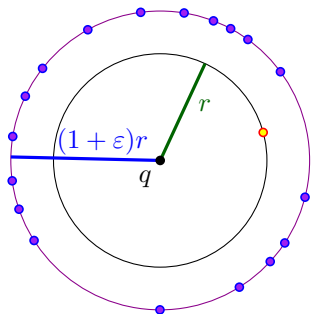
Worst case



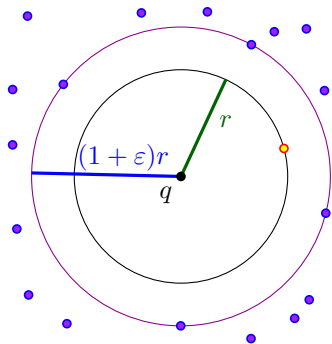
Maybe realistic case

Work hard to push away a single point...

21: The worst/good case for LSH



Worst case



Maybe realistic case

Work hard to push away a single point...

22: Data sensitive LSH.

- 1 **Idea:** interpret LSH as estimator of local density.
- 2 If data set is sparse near query point, LSH would hit NN point quickly...
- 3 ...density estimation would tell us that there are no more close by points.
- 4 Do exponential search – start with an insensitive LSH scheme (that is fast).
- 5 ... use more sensitive LSHs, till density estimation tell us we are done.
- 6 if all fails, the last LSH data-structure used is essentially the old LSH scheme.

22: Data sensitive LSH.

- 1 **Idea:** interpret LSH as estimator of local density.
- 2 If data set is sparse near query point, LSH would hit NN point quickly...
- 3 ...density estimation would tell us that there are no more close by points.
- 4 Do exponential search – start with an insensitive LSH scheme (that is fast).
- 5 ... use more sensitive LSHs, till density estimation tell us we are done.
- 6 if all fails, the last LSH data-structure used is essentially the old LSH scheme.

22: Data sensitive LSH.

- 1 **Idea:** interpret LSH as estimator of local density.
- 2 If data set is sparse near query point, LSH would hit NN point quickly...
- 3 ...density estimation would tell us that there are no more close by points.
- 4 Do exponential search – start with an insensitive LSH scheme (that is fast).
- 5 ... use more sensitive LSHs, till density estimation tell us we are done.
- 6 if all fails, the last LSH data-structure used is essentially the old LSH scheme.

22: Data sensitive LSH.

- 1 **Idea:** interpret LSH as estimator of local density.
- 2 If data set is sparse near query point, LSH would hit NN point quickly...
- 3 ...density estimation would tell us that there are no more close by points.
- 4 Do exponential search – start with an insensitive LSH scheme (that is fast).
- 5 ... use more sensitive LSHs, till density estimation tell us we are done.
- 6 if all fails, the last LSH data-structure used is essentially the old LSH scheme.

22: Data sensitive LSH.

- 1 **Idea:** interpret LSH as estimator of local density.
- 2 If data set is sparse near query point, LSH would hit NN point quickly...
- 3 ...density estimation would tell us that there are no more close by points.
- 4 Do exponential search – start with an insensitive LSH scheme (that is fast).
- 5 ... use more sensitive LSHs, till density estimation tell us we are done.
- 6 if all fails, the last LSH data-structure used is essentially the old LSH scheme.

22: Data sensitive LSH.

- 1 **Idea:** interpret LSH as estimator of local density.
- 2 If data set is sparse near query point, LSH would hit NN point quickly...
- 3 ...density estimation would tell us that there are no more close by points.
- 4 Do exponential search – start with an insensitive LSH scheme (that is fast).
- 5 ... use more sensitive LSHs, till density estimation tell us we are done.
- 6 if all fails, the last LSH data-structure used is essentially the old LSH scheme.

23: Setup

① $\mathbf{P} \subseteq \{0, 1\}^d$: n points.
 r : radius, $\epsilon > 0$.

② Preprocessing:

$$i = 1, \dots, \mathcal{N} = \frac{\log n}{1+\epsilon}:$$

① $f_i(\ell) = (1 - 1/r)^{\ell i} \leq \exp(-\ell i/r)$

② $s(i) = \frac{c \log n}{f_i(r)} = O(e^i \log n)$,

③ $\mathcal{F}_i = \{m_1, m_2, \dots, m_{s(i)} \in \mathcal{D}_{1/r}^i\}$

④ $\forall m \in \mathcal{F}_i$: store $m\mathbf{P}$ in hash table.

⑤ colliding with \mathbf{q} : $\mathcal{C}_m(\mathbf{q}) = \{\mathbf{x} \in \mathbf{P} \mid m\mathbf{q} = m\mathbf{x}\}$,

⑥ Can extract $|\mathcal{C}_m(\mathbf{q})|$ in $O(d)$ time.

23: Setup

① $\mathbf{P} \subseteq \{0, 1\}^d$: n points.
 r : radius, $\epsilon > 0$.

② Preprocessing:

$$i = 1, \dots, \mathcal{N} = \frac{\log n}{1+\epsilon}:$$

① $f_i(\ell) = (1 - 1/r)^{\ell i} \leq \exp(-\ell i/r)$

② $s(i) = \frac{c \log n}{f_i(r)} = O(e^i \log n)$,

③ $\mathcal{F}_i = \{m_1, m_2, \dots, m_{s(i)} \in \mathcal{D}_{1/r}^i\}$

④ $\forall m \in \mathcal{F}_i$: store $m\mathbf{P}$ in hash table.

⑤ colliding with \mathbf{q} : $\mathcal{C}_m(\mathbf{q}) = \{\mathbf{x} \in \mathbf{P} \mid m\mathbf{q} = m\mathbf{x}\}$,

⑥ Can extract $|\mathcal{C}_m(\mathbf{q})|$ in $O(d)$ time.

23: Setup

- 1 $\mathbf{P} \subseteq \{0, 1\}^d$: n points.
 r : radius, $\epsilon > 0$.

- 2 Preprocessing:

$$i = 1, \dots, \mathcal{N} = \frac{\log n}{1+\epsilon}:$$

- 1 $f_i(\ell) = (1 - 1/r)^{\ell i} \leq \exp(-\ell i/r)$

- 2 $s(i) = \frac{c \log n}{f_i(r)} = O(e^i \log n)$,

- 3 $\mathcal{F}_i = \{m_1, m_2, \dots, m_{s(i)} \in \mathcal{D}_{1/r}^i\}$

- 4 $\forall m \in \mathcal{F}_i$: store $m\mathbf{P}$ in hash table.

- 5 colliding with \mathbf{q} : $\mathcal{C}_m(\mathbf{q}) = \{\mathbf{x} \in \mathbf{P} \mid m\mathbf{q} = m\mathbf{x}\}$,

- 6 Can extract $|\mathcal{C}_m(\mathbf{q})|$ in $O(d)$ time.

23: Setup

① $\mathbf{P} \subseteq \{0, 1\}^d$: n points.
 r : radius, $\epsilon > 0$.

② Preprocessing:

$$i = 1, \dots, \mathcal{N} = \frac{\log n}{1+\epsilon}:$$

① $f_i(\ell) = (1 - 1/r)^{\ell i} \leq \exp(-\ell i/r)$

② $s(i) = \frac{c \log n}{f_i(r)} = O(e^i \log n)$,

③ $\mathcal{F}_i = \{m_1, m_2, \dots, m_{s(i)} \in \mathcal{D}_{1/r}^i\}$

④ $\forall m \in \mathcal{F}_i$: store $m\mathbf{P}$ in hash table.

⑤ colliding with \mathbf{q} : $\mathcal{C}_m(\mathbf{q}) = \{\mathbf{x} \in \mathbf{P} \mid m\mathbf{q} = m\mathbf{x}\}$,

⑥ Can extract $|\mathcal{C}_m(\mathbf{q})|$ in $O(d)$ time.

23: Setup

① $\mathbf{P} \subseteq \{0, 1\}^d$: n points.
 r : radius, $\epsilon > 0$.

② Preprocessing:

$$i = 1, \dots, \mathcal{N} = \frac{\log n}{1+\epsilon}:$$

① $f_i(\ell) = (1 - 1/r)^{\ell i} \leq \exp(-\ell i/r)$

② $s(i) = \frac{c \log n}{f_i(r)} = O(e^i \log n)$,

③ $\mathcal{F}_i = \{m_1, m_2, \dots, m_{s(i)} \in \mathcal{D}_{1/r}^i\}$

④ $\forall m \in \mathcal{F}_i$: store $m\mathbf{P}$ in hash table.

⑤ colliding with \mathbf{q} : $\mathcal{C}_m(\mathbf{q}) = \{\mathbf{x} \in \mathbf{P} \mid m\mathbf{q} = m\mathbf{x}\}$,

⑥ Can extract $|\mathcal{C}_m(\mathbf{q})|$ in $O(d)$ time.

23: Setup

- 1 $\mathbf{P} \subseteq \{0, 1\}^d$: n points.
 r : radius, $\epsilon > 0$.
- 2 Preprocessing:
 $i = 1, \dots, \mathcal{N} = \frac{\log n}{1+\epsilon}$:
 - 1 $f_i(\ell) = (1 - 1/r)^{\ell i} \leq \exp(-\ell i/r)$
 - 2 $s(i) = \frac{c \log n}{f_i(r)} = O(e^i \log n)$,
 - 3 $\mathcal{F}_i = \{m_1, m_2, \dots, m_{s(i)} \in \mathcal{D}_{1/r}^i\}$
 - 4 $\forall m \in \mathcal{F}_i$: store $m\mathbf{P}$ in hash table.
 - 5 colliding with \mathbf{q} : $\mathcal{C}_m(\mathbf{q}) = \{\mathbf{x} \in \mathbf{P} \mid m\mathbf{q} = m\mathbf{x}\}$,
 - 6 Can extract $|\mathcal{C}_m(\mathbf{q})|$ in $O(d)$ time.

24: Query.

- 1 Query: $\mathbf{q} \in \mathbb{H}^d$
- 2 For $i = 1 \dots$
 - 1 $X_i = \sum_{m \in \mathcal{F}_i} |C_m(\mathbf{q})|$. (# of colliding points.)
 - 2 If $X_i > 4s(i)$ then $i++$, continue to next iteration.
 - 3 If $X_i \leq 4s(i)$, then
 - 1 extract the X_i points
 - 2 Return closest point among them.
 - 4 If $i = \mathcal{N}$ then do standard LSH query.

24: Query.

- 1 Query: $\mathbf{q} \in \mathbb{H}^d$
- 2 For $i = 1 \dots$
 - 1 $X_i = \sum_{m \in \mathcal{F}_i} |\mathcal{C}_m(\mathbf{q})|$. (# of colliding points.)
 - 2 If $X_i > 4s(i)$ then $i++$, continue to next iteration.
 - 3 If $X_i \leq 4s(i)$, then
 - 1 extract the X_i points
 - 2 Return closest point among them.
 - 4 If $i = \mathcal{N}$ then do standard LSH query.

24: Query.

- 1 Query: $\mathbf{q} \in \mathbb{H}^d$
- 2 For $i = 1 \dots$
 - 1 $X_i = \sum_{m \in \mathcal{F}_i} |\mathcal{C}_m(\mathbf{q})|$. (# of colliding points.)
 - 2 If $X_i > 4s(i)$ then $i++$, continue to next iteration.
 - 3 If $X_i \leq 4s(i)$, then
 - 1 extract the X_i points
 - 2 Return closest point among them.
 - 4 If $i = \mathcal{N}$ then do standard LSH query.

24: Query.

- 1 Query: $\mathbf{q} \in \mathbb{H}^d$
- 2 For $i = 1 \dots$
 - 1 $X_i = \sum_{m \in \mathcal{F}_i} |\mathcal{C}_m(\mathbf{q})|$. (# of colliding points.)
 - 2 If $X_i > 4s(i)$ then $i++$, continue to next iteration.
 - 3 If $X_i \leq 4s(i)$, then
 - 1 extract the X_i points
 - 2 Return closest point among them.
 - 4 If $i = \mathcal{N}$ then do standard LSH query.

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = s(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + s(i)))$.

If $X_i = O(s(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = s(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + s(i)))$.

If $X_i = O(s(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = \mathbf{s}(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + \mathbf{s}(i)))$.

If $X_i = O(\mathbf{s}(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = s(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + s(i)))$.

If $X_i = O(s(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = s(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + s(i)))$.

If $X_i = O(s(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = s(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + s(i)))$.

If $X_i = O(s(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = s(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + s(i)))$.

If $X_i = O(s(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

25: Intuition

Expected number of collisions with \mathbf{q} , for a $m \in \mathcal{F}_i$ is

$$\gamma_i = \mathbf{E} \left[|\mathcal{C}_m(\mathbf{q})| \right] = \sum_{\mathbf{x} \in \mathcal{P}} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \leq \sum_{\mathbf{x} \in \mathcal{P}} \exp\left(-\frac{\mathbf{d}_H(\mathbf{q}, \mathbf{x})i}{r}\right)$$

A convolution over the input set...

Expected number collisions with \mathbf{q} (with all projections \mathcal{F}_i):

$$\Gamma_i = \mathbf{E}[X_i] = s(i)\gamma_i.$$

If scan collision lists in i th level: query time $O(d(X_i + s(i)))$.

If $X_i = O(s(i))$, then query time is small.

Intuitively: Bigly i stronger DS push points away from \mathbf{q} .

If lucky, need to push little.

Standard LSH: ends up immediately with \mathcal{F}_N .

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= s(i) \sum_{\mathbf{x} \in P} f_i(d_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= s(i) \sum_{\mathbf{x} \in P} f_i(d_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= s(i) \sum_{\mathbf{x} \in P} f_i(d_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= \mathbf{s}(i) \sum_{\mathbf{x} \in P} \mathbf{f}_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= s(i) \sum_{\mathbf{x} \in P} f_i(d_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= s(i) \sum_{\mathbf{x} \in P} f_i(d_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= s(i) \sum_{\mathbf{x} \in P} f_i(\mathbf{d}_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

26: Example: If the data is locally low dimensional

Γ_i depends on how the data looks like near the query.

Assume locally near query point, data is a uniform low dimensional point set.

Assume # points in distance ℓ from $\mathbf{q} \leq \#(\ell) = O((\ell/r)^k)$.

k : small, r : distance to NN to \mathbf{q} .

$$\begin{aligned}\Gamma_i &= s(i) \sum_{\mathbf{x} \in P} f_i(d_H(\mathbf{q}, \mathbf{x})) \\ &\leq O(e^i \log n) \sum_{j=1}^{\infty} \#(jr) \exp(-ji) \\ &= \sum_{j=1}^{\infty} O(j^k \exp(i - ji) \log n).\end{aligned}$$

Setting $i = k$:

$$\Gamma_k \leq \sum_{j=1}^{\infty} O\left(\left(\frac{j}{e^{j-1}}\right)^k \log n\right) \leq \sum_{j=1}^{\infty} O\left(\frac{j}{2^j} \log n\right) \leq O(\log n)$$

Expectation: Stop $O(k)$ rounds. Query time: $O(d \log n)$.

27: Analysis

Lemma

If \exists point with distance $\leq r$ from \mathbf{q} , then algorithm computes, w.h.p., point \mathbf{x} s.t., $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon)r$.

Proof.

$\mathbf{q}^* \in \mathbf{P}$: nearest neighbor to \mathbf{q} .

For DS_{*i*}: probability \mathbf{q}^* not collide with \mathbf{q}
 $\leq (1 - f_i(r))^{s(i)} \leq \exp(-f_i(r)s(i)) = \exp(-O(\log n)) \leq 1/n^{O(1)}$.

Since $s(i) = O(e^i \log n)$

$i < \mathcal{N}$: Algorithm stops and scans all lists.

$i = \mathcal{N}$: Just reg LSH. (**Approximation only in this case!**)

27: Analysis

Lemma

If \exists point with distance $\leq r$ from \mathbf{q} , then algorithm computes, w.h.p., point \mathbf{x} s.t., $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon)r$.

Proof.

$\mathbf{q}^* \in \mathbf{P}$: nearest neighbor to \mathbf{q} .

For DS_i : probability \mathbf{q}^* not collide with \mathbf{q}
 $\leq (1 - f_i(r))^{s(i)} \leq \exp(-f_i(r)s(i)) = \exp(-O(\log n)) \leq 1/n^{O(1)}$.

Since $s(i) = O(e^i \log n)$

$i < \mathcal{N}$: Algorithm stops and scans all lists.

$i = \mathcal{N}$: Just reg LSH. (Approximation only in this case!)

27: Analysis

Lemma

If \exists point with distance $\leq r$ from \mathbf{q} , then algorithm computes, w.h.p., point \mathbf{x} s.t., $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon)r$.

Proof.

$\mathbf{q}^* \in \mathbf{P}$: nearest neighbor to \mathbf{q} .

For DS_i : probability \mathbf{q}^* not collide with \mathbf{q}
 $\leq (1 - f_i(r))^{s(i)} \leq \exp(-f_i(r)s(i)) = \exp(-O(\log n)) \leq 1/n^{O(1)}$.

Since $s(i) = O(e^i \log n)$

$i < \mathcal{N}$: Algorithm stops and scans all lists.

$i = \mathcal{N}$: Just reg LSH. (Approximation only in this case!)

27: Analysis

Lemma

If \exists point with distance $\leq r$ from \mathbf{q} , then algorithm computes, w.h.p., point \mathbf{x} s.t., $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon)r$.

Proof.

$\mathbf{q}^* \in \mathbf{P}$: nearest neighbor to \mathbf{q} .

For DS_i : probability \mathbf{q}^* not collide with \mathbf{q}
 $\leq (1 - f_i(r))^{s(i)} \leq \exp(-f_i(r)s(i)) = \exp(-O(\log n)) \leq 1/n^{O(1)}$.

Since $s(i) = O(e^i \log n)$

$i < \mathcal{N}$: Algorithm stops and scans all lists.

$i = \mathcal{N}$: Just reg LSH. (Approximation only in this case!)

27: Analysis

Lemma

If \exists point with distance $\leq r$ from \mathbf{q} , then algorithm computes, w.h.p., point \mathbf{x} s.t., $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon)r$.

Proof.

$\mathbf{q}^* \in \mathbf{P}$: nearest neighbor to \mathbf{q} .

For DS_i : probability \mathbf{q}^* not collide with \mathbf{q}
 $\leq (1 - f_i(r))^{s(i)} \leq \exp(-f_i(r)s(i)) = \exp(-O(\log n)) \leq 1/n^{O(1)}$.

Since $s(i) = O(e^i \log n)$

$i < \mathcal{N}$: Algorithm stops and scans all lists.

$i = \mathcal{N}$: Just reg LSH. (Approximation only in this case!)

27: Analysis

Lemma

If \exists point with distance $\leq r$ from \mathbf{q} , then algorithm computes, w.h.p., point \mathbf{x} s.t., $\mathbf{d}_H(\mathbf{q}, \mathbf{x}) \leq (1 + \epsilon)r$.

Proof.

$\mathbf{q}^* \in \mathbf{P}$: nearest neighbor to \mathbf{q} .

For DS_i : probability \mathbf{q}^* not collide with \mathbf{q}
 $\leq (1 - f_i(r))^{s(i)} \leq \exp(-f_i(r)s(i)) = \exp(-O(\log n)) \leq 1/n^{O(1)}$.

Since $s(i) = O(e^i \log n)$

$i < \mathcal{N}$: Algorithm stops and scans all lists.

$i = \mathcal{N}$: Just reg LSH. (**Approximation only in this case!**)

28: Value of $s(i)$

Lemma

For any $i > 0$, we have $s(i) = O(e^i \log n)$.

Proof.

$$f_i(r) = (1 - 1/r)^{(r-1)ri/(r-1)} \geq \exp\left(-r \frac{i}{r-1}\right) = \exp\left(-i - \frac{i}{r-1}\right)$$

since $\forall m > 0: (1 - 1/m)^{m-1} \geq 1/e \geq (1 - 1/m)^m$.

Can assume that $r > \log n \geq i$ (duplicate coordinates)...

$$\implies f_i(r) = \Theta(\exp(-i)).$$

$$\implies s(i) = O\left(\frac{\log n}{f_i(r)}\right) = O(e^i \log n). \quad \square$$

28: Value of $s(i)$

Lemma

For any $i > 0$, we have $s(i) = O(e^i \log n)$.

Proof.

$$f_i(r) = (1 - 1/r)^{(r-1)ri/(r-1)} \geq \exp\left(-r \frac{i}{r-1}\right) = \exp\left(-i - \frac{i}{r-1}\right)$$

since $\forall m > 0: (1 - 1/m)^{m-1} \geq 1/e \geq (1 - 1/m)^m$.

Can assume that $r > \log n \geq i$ (duplicate coordinates)...

$$\implies f_i(r) = \Theta(\exp(-i)).$$

$$\implies s(i) = O\left(\frac{\log n}{f_i(r)}\right) = O(e^i \log n). \quad \square$$

28: Value of $s(i)$

Lemma

For any $i > 0$, we have $s(i) = O(e^i \log n)$.

Proof.

$$f_i(r) = (1 - 1/r)^{(r-1)ri/(r-1)} \geq \exp\left(-r \frac{i}{r-1}\right) = \exp\left(-i - \frac{i}{r-1}\right)$$

since $\forall m > 0: (1 - 1/m)^{m-1} \geq 1/e \geq (1 - 1/m)^m$.

Can assume that $r > \log n \geq i$ (duplicate coordinates)...

$$\implies f_i(r) = \Theta(\exp(-i)).$$

$$\implies s(i) = O\left(\frac{\log n}{f_i(r)}\right) = O(e^i \log n). \quad \square$$

28: Value of $s(i)$

Lemma

For any $i > 0$, we have $s(i) = O(e^i \log n)$.

Proof.

$$f_i(r) = (1 - 1/r)^{(r-1)ri/(r-1)} \geq \exp\left(-r \frac{i}{r-1}\right) = \exp\left(-i - \frac{i}{r-1}\right)$$

since $\forall m > 0: (1 - 1/m)^{m-1} \geq 1/e \geq (1 - 1/m)^m$.

Can assume that $r > \log n \geq i$ (duplicate coordinates)...

$$\implies f_i(r) = \Theta(\exp(-i)).$$

$$\implies s(i) = O\left(\frac{\log n}{f_i(r)}\right) = O(e^i \log n). \quad \square$$

28: Value of $s(i)$

Lemma

For any $i > 0$, we have $s(i) = O(e^i \log n)$.

Proof.

$$f_i(r) = (1 - 1/r)^{(r-1)ri/(r-1)} \geq \exp\left(-r \frac{i}{r-1}\right) = \exp\left(-i - \frac{i}{r-1}\right)$$

since $\forall m > 0: (1 - 1/m)^{m-1} \geq 1/e \geq (1 - 1/m)^m$.

Can assume that $r > \log n \geq i$ (duplicate coordinates)...

$$\implies f_i(r) = \Theta(\exp(-i)).$$

$$\implies s(i) = O\left(\frac{\log n}{f_i(r)}\right) = O(e^i \log n). \quad \square$$

29: Worst case query time

Lemma

For a query point \mathbf{q} , the worst case query time is $O(dn^{1/(1+\epsilon)} \log n)$, with high probability.

30: Data-dependent query time

Lemma

$$\mathbf{P} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{H}^d$$

$$\ell_i = \mathbf{d}_H(\mathbf{q}, \mathbf{x}_i), \text{ for } i = 1, \dots, n,$$

Δ : smallest value s.t.

$$\sum_{i=1}^n \exp\left(-\frac{\Delta \ell_i}{r}\right) \leq 1.$$

\implies expected query time $O(de^\Delta \log n)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^N Y_j s(j)\right]\right)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^{\mathcal{N}} Y_j s(j)\right]\right)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^N Y_j s(j)\right]\right)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^N Y_j s(j)\right]\right)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^N Y_j s(j)\right]\right)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^N Y_j s(j)\right]\right)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^N Y_j s(j)\right]\right)$.

31: Proof

Condition implies: $\gamma_j \leq \gamma_\Delta \leq 1$, for any $j \geq \Delta$.

$$\implies j \geq \Delta: \Gamma_j = \mathbf{E}[X_j] = \gamma_j s(j) \leq s(j).$$

Markov's inequality, $\Pr[X_j \leq 4s(j)] \geq 1 - 1/4$.

Then the algorithm would terminate in this iteration.

$Y_j = 1$: indicator the algorithm reached j th iteration.

Algorithm files in iterations $\Delta, \Delta + 1, \dots, j - 1$.

$$\implies \Pr[Y_j = 1] \leq 1/4^{j-\Delta}.$$

j th iteration (if happens): $O(ds(j))$ time.

Overall expected running time is $O\left(\mathbf{E}\left[d \sum_{j=1}^{\mathcal{N}} Y_j s(j)\right]\right)$.

32: Proof continued...

$$\begin{aligned} & O\left(d \sum_{j=1}^{\Delta} s(j) + d \sum_{j=\Delta+1}^{\mathcal{N}} \frac{s(j)}{4^{j-\Delta}}\right) \\ &= O\left(de^{\Delta} \log n + de^{\Delta} \sum_{j=\Delta+1}^{\mathcal{N}} \frac{e^{j-\Delta} \log n}{4^{j-\Delta}}\right) \\ &= O(de^{\Delta} \log n), \end{aligned}$$

using the bound $s(j) = O(e^j \log n)$.



32: Proof continued...

$$\begin{aligned} & O\left(d \sum_{j=1}^{\Delta} s(j) + d \sum_{j=\Delta+1}^{\mathcal{N}} \frac{s(j)}{4^{j-\Delta}}\right) \\ &= O\left(de^{\Delta} \log n + de^{\Delta} \sum_{j=\Delta+1}^{\mathcal{N}} \frac{e^{j-\Delta} \log n}{4^{j-\Delta}}\right) \\ &= O(de^{\Delta} \log n), \end{aligned}$$

using the bound $s(j) = O(e^j \log n)$.



32: Proof continued...

$$\begin{aligned} & O\left(d \sum_{j=1}^{\Delta} s(j) + d \sum_{j=\Delta+1}^{\mathcal{N}} \frac{s(j)}{4^{j-\Delta}}\right) \\ &= O\left(de^{\Delta} \log n + de^{\Delta} \sum_{j=\Delta+1}^{\mathcal{N}} \frac{e^{j-\Delta} \log n}{4^{j-\Delta}}\right) \\ &= O(de^{\Delta} \log n), \end{aligned}$$

using the bound $s(j) = O(e^j \log n)$.



32: Proof continued...

$$\begin{aligned} & O\left(d \sum_{j=1}^{\Delta} s(j) + d \sum_{j=\Delta+1}^{\mathcal{N}} \frac{s(j)}{4^{j-\Delta}}\right) \\ &= O\left(de^{\Delta} \log n + de^{\Delta} \sum_{j=\Delta+1}^{\mathcal{N}} \frac{e^{j-\Delta} \log n}{4^{j-\Delta}}\right) \\ &= O(de^{\Delta} \log n), \end{aligned}$$

using the bound $s(j) = O(e^j \log n)$.



33: The result

$\mathbf{P} \subseteq \mathbb{H}^d = \{0, 1\}^d$: set of n points,
 $\epsilon > 0$ and r .

Preprocess in $O(dn^{1+1/(1+\epsilon)} \log n)$ time/ space.

Given $\mathbf{q} \in \mathbb{H}^d$, decide (approximately) if $d_H(\mathbf{q}, \mathbf{P}) \leq r$, in
 $O(dn^{1/(1+\epsilon)} \log n)$.

If query is “easy”, DS returns *exact* nearest neighbor.

If $d_H(\mathbf{q}, \mathbf{P}) \leq r$, and $\exists \Delta < (\log n)/(1 + \epsilon)$ s.t.
 $\sum_{\mathbf{x} \in \mathbf{P}} \exp(-d_H(\mathbf{q}, \mathbf{x})\Delta/r) \leq 1$

\implies Exact NN in $O(d \exp(\Delta) \log n)$ expected time.

33: The result

$\mathbf{P} \subseteq \mathbb{H}^d = \{0, 1\}^d$: set of n points,
 $\epsilon > 0$ and r .

Preprocess in $O(dn^{1+1/(1+\epsilon)} \log n)$ time/ space.

Given $\mathbf{q} \in \mathbb{H}^d$, decide (approximately) if $d_H(\mathbf{q}, \mathbf{P}) \leq r$, in
 $O(dn^{1/(1+\epsilon)} \log n)$.

If query is “easy”, DS returns *exact* nearest neighbor.

If $d_H(\mathbf{q}, \mathbf{P}) \leq r$, and $\exists \Delta < (\log n)/(1 + \epsilon)$ s.t.
 $\sum_{\mathbf{x} \in \mathbf{P}} \exp(-d_H(\mathbf{q}, \mathbf{x})\Delta/r) \leq 1$

\implies Exact NN in $O(d \exp(\Delta) \log n)$ expected time.

33: The result

$\mathbf{P} \subseteq \mathbb{H}^d = \{0, 1\}^d$: set of n points,
 $\epsilon > 0$ and r .

Preprocess in $O(dn^{1+1/(1+\epsilon)} \log n)$ time/ space.

Given $\mathbf{q} \in \mathbb{H}^d$, decide (approximately) if $d_H(\mathbf{q}, \mathbf{P}) \leq r$, in
 $O(dn^{1/(1+\epsilon)} \log n)$.

If query is “easy”, DS returns *exact* nearest neighbor.

If $d_H(\mathbf{q}, \mathbf{P}) \leq r$, and $\exists \Delta < (\log n)/(1 + \epsilon)$ s.t.
 $\sum_{\mathbf{x} \in \mathbf{P}} \exp(-d_H(\mathbf{q}, \mathbf{x})\Delta/r) \leq 1$

\implies Exact NN in $O(d \exp(\Delta) \log n)$ expected time.

33: The result

$\mathbf{P} \subseteq \mathbb{H}^d = \{0, 1\}^d$: set of n points,
 $\epsilon > 0$ and r .

Preprocess in $O(dn^{1+1/(1+\epsilon)} \log n)$ time/ space.

Given $\mathbf{q} \in \mathbb{H}^d$, decide (approximately) if $d_H(\mathbf{q}, \mathbf{P}) \leq r$, in
 $O(dn^{1/(1+\epsilon)} \log n)$.

If query is “easy”, DS returns *exact* nearest neighbor.

If $d_H(\mathbf{q}, \mathbf{P}) \leq r$, and $\exists \Delta < (\log n)/(1 + \epsilon)$ s.t.
 $\sum_{\mathbf{x} \in \mathbf{P}} \exp(-d_H(\mathbf{q}, \mathbf{x})\Delta/r) \leq 1$

\implies Exact NN in $O(d \exp(\Delta) \log n)$ expected time.

33: The result

$\mathbf{P} \subseteq \mathbb{H}^d = \{0, 1\}^d$: set of n points,
 $\varepsilon > 0$ and r .

Preprocess in $O(dn^{1+1/(1+\varepsilon)} \log n)$ time/ space.

Given $\mathbf{q} \in \mathbb{H}^d$, decide (approximately) if $d_H(\mathbf{q}, \mathbf{P}) \leq r$, in
 $O(dn^{1/(1+\varepsilon)} \log n)$.

If query is “easy”, DS returns *exact* nearest neighbor.

If $d_H(\mathbf{q}, \mathbf{P}) \leq r$, and $\exists \Delta < (\log n)/(1 + \varepsilon)$ s.t.

$$\sum_{\mathbf{x} \in \mathbf{P}} \exp(-d_H(\mathbf{q}, \mathbf{x})\Delta/r) \leq 1$$

\implies Exact NN in $O(d \exp(\Delta) \log n)$ expected time.

34: Remarks

- 1 If data is k dimensional (bounded growth)
Then query time $O(d \log n)$.
Known: **[Datar, Immorlica, Indyk, and Mirrokni, 2004]**.
New DS more general. Handles this case with no modification.
- 2 Fine tuning the LSH scheme to the hardness of the given data is not a new idea.
 - 1 **[Andoni, Datar, Immorlica, Indyk, and Mirrokni, 2006]**
 - 2 Fine tune LSH construction parameters for set of queries, to optimize overall query time.
 - 3 New DS: adapts parameters on the fly during query process, depending on query.

34: Remarks

- 1 If data is k dimensional (bounded growth)
Then query time $O(d \log n)$.
Known: **[Datar, Immorlica, Indyk, and Mirrokni, 2004]**.
New DS more general. Handles this case with no modification.
- 2 Fine tuning the LSH scheme to the hardness of the given data is not a new idea.
 - 1 **[Andoni, Datar, Immorlica, Indyk, and Mirrokni, 2006]**
 - 2 Fine tune LSH construction parameters for set of queries, to optimize overall query time.
 - 3 New DS: adapts parameters on the fly during query process, depending on query.

34: Remarks

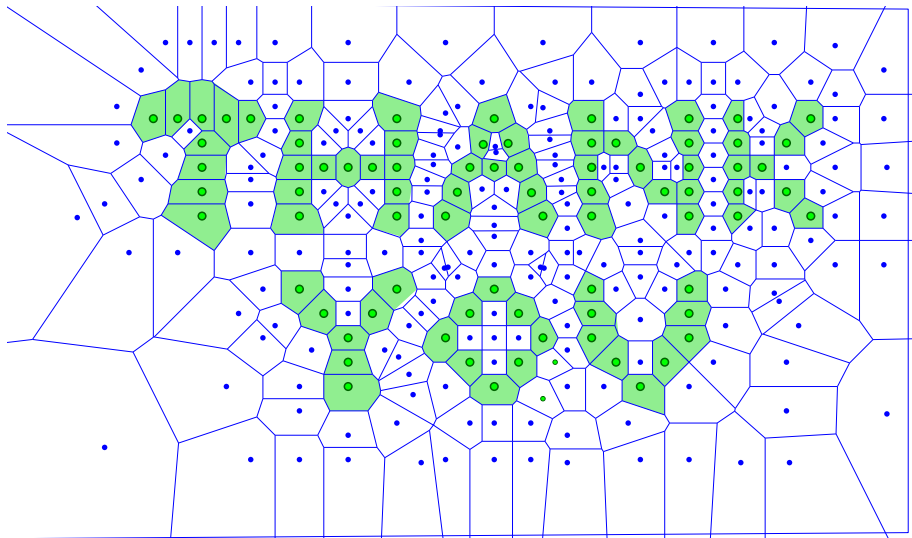
- 1 If data is k dimensional (bounded growth)
Then query time $O(d \log n)$.
Known: **[Datar, Immorlica, Indyk, and Mirrokni, 2004]**.
New DS more general. Handles this case with no modification.
- 2 Fine tuning the LSH scheme to the hardness of the given data is not a new idea.
 - 1 **[Andoni, Datar, Immorlica, Indyk, and Mirrokni, 2006]**
 - 2 Fine tune LSH construction parameters for set of queries, to optimize overall query time.
 - 3 New DS: adapts parameters on the fly during query process, depending on query.

34: Remarks

- ① If data is k dimensional (bounded growth)
Then query time $O(d \log n)$.
Known: **[Datar, Immorlica, Indyk, and Mirrokni, 2004]**.
New DS more general. Handles this case with no modification.
- ② Fine tuning the LSH scheme to the hardness of the given data is not a new idea.
 - ① **[Andoni, Datar, Immorlica, Indyk, and Mirrokni, 2006]**
 - ② Fine tune LSH construction parameters for set of queries, to optimize overall query time.
 - ③ New DS: adapts parameters on the fly during query process, depending on query.

34: Remarks

- ① If data is k dimensional (bounded growth)
Then query time $O(d \log n)$.
Known: **[Datar, Immorlica, Indyk, and Mirrokni, 2004]**.
New DS more general. Handles this case with no modification.
- ② Fine tuning the LSH scheme to the hardness of the given data is not a new idea.
 - ① **[Andoni, Datar, Immorlica, Indyk, and Mirrokni, 2006]**
 - ② Fine tune LSH construction parameters for set of queries, to optimize overall query time.
 - ③ New DS: adapts parameters on the fly during query process, depending on query.



A. Andoni, M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni. Locality-sensitive hashing using stable distributions. In T. Darrell, P. Indyk, and G. Shakhnarovich, editors, *Nearest-Neighbor Methods in Learning and Vision: Theory and Practice*, pages 61–72. MIT Press, 2006. URL <https://graphics.stanford.edu/courses/cs468-06-fall/Papers/13%201sh06.pdf>.

M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni. Locality-sensitive hashing scheme based on p -stable distributions. In *Proc. 20th Annu. Sympos. Comput. Geom. (SoCG)*, pages 253–262, 2004. doi: [10.1145/997817.997857](https://doi.org/10.1145/997817.997857).

P. Indyk and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proc. 30th Annu. ACM Sympos. Theory Comput. (STOC)*, pages 604–613, 1998. doi: [10.1145/276698.276876](https://doi.org/10.1145/276698.276876).

