# Approximate Near Neighbor Search: Locality-sensitive Hashing

**COMPCSI 753: Algorithms for Massive Data**

Ninh Pham

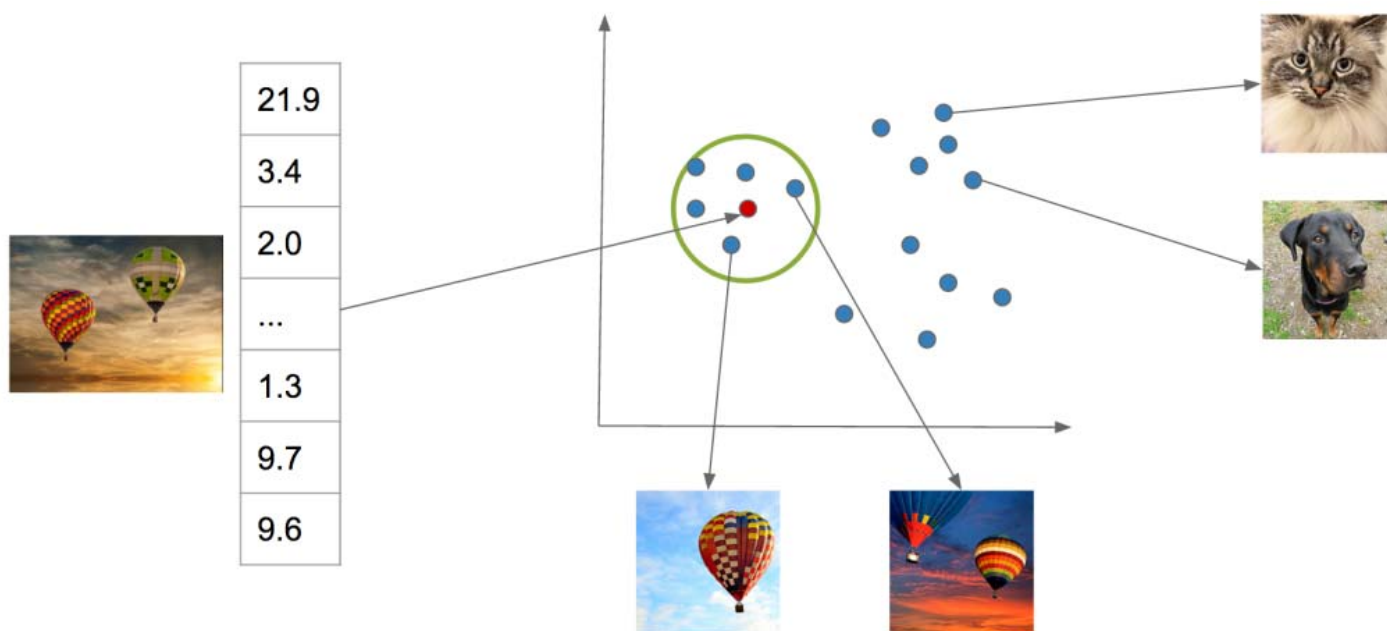University of Auckland

Auckland, Aug 10, 2020

1

# Outline

- Popular similarity/distance measure
  - Definitions
  - Applications

- Locality-sensitive hashing framework for approximate near neighbor search
  - LSH definition
  - LSH framework for near neighbor search

# A common metaphor

- We can present many problems as finding "similar" items or finding near-neighbors in high dimensions.

- Near neighbors search:



Source: https://code.flickr.net/2017/03/07/introducing-similarity-search-at-flickr/

3

# Distance/similarity measure

- For each application, we need to define a similarity/distance measure to distinguish between
  - Similar and dissimilar objects.
  - Near and far apart objects.

- Popular similarity/distance measures:
  - Jaccard similarity/distance
  - Cosine similarity/distance
  - Hamming distance
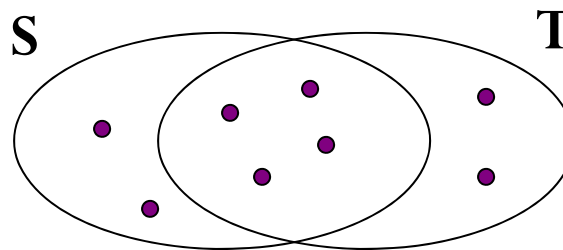  - Euclidean distance (L2)
  - Manhattan distance (L1)

# Jaccard similarity/distance

- Jaccard similarity of two sets **S** and **T**:

$$J(S,T) = \frac{|S \cap T|}{|S \cup T|}$$

- Jaccard distance of two sets **S** and **T**: $1 - J(S,T)$

- Example:

**S**  **T**

4 items in intersection

8 items in union

J(S,T) = 1/2

1 - J(S,T) = 1/2

# Applications

- Texts/Documents:
  - Present text/documents as a set of shingles
  - Usage: Detect plagiarism, mirror pages, articles from same source

- Recommender system data:
  - Online-purchases: users as sets, items as set elements
  - Netflix: users as sets, movies rated as set elements
  - Usages: Collaborative filtering recommender systems need to compute the similarity between users-users, items-items for accurate recommendation

# Applications

- Market basket data:
  - Items as sets, transaction IDs as set elements.
  - Usage: Association rules $X \rightarrow Y$ if $J(X, Y) \geq s$.

- Example:
  - Bread = {1, 2, 4, 5}
  - Milk = {1, 3, 4, 5}
  - Diapers = {2, 3, 4, 5}
  - Beer = {2, 3, 4}
  - Eggs = {2}
  - Cola = {3, 5}
  - J(Diapers, Beer) = 3/4

| ID | Items |
|----|-------|
| 1 | {Bread, Milk} |
| 2 | {Bread, Diapers, Beer, Eggs} |
| 3 | {Milk, Diapers, Beer, Cola} |
| 4 | {Bread, Milk, Diapers, Beer} |
| 5 | {Bread, Milk, Diapers, Cola} |
| … | … |

market basket transactions

{Diapers, Beer}      Example of a frequent itemset

{Diapers} → {Beer}   Example of an association rule

7

# MinHash [Broder'97]

- Permute the Boolean matrix with a random permutation $\pi$

- MinHash function:

$h_\pi(S)$ = the index of the first row with value **1** of **S** (in the permuted order $\pi$)

$$h_\pi(S) = \min_\pi \pi(S)$$

- $Pr_\pi[h_\pi(S) = h_\pi(T)] = J(S,T)$

|   | **A** | **B** | **C** | **D** |
|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 0 |
| 2 | 1 | 1 | 0 | 1 |
| 3 | 0 | 1 | 0 | 1 |
| 4 | 0 | 0 | 0 | 1 |
| 5 | 1 | 0 | 0 | 1 |
| 6 | 1 | 1 | 1 | 0 |
| 7 | 1 | 0 | 1 | 0 |

Shingles

Documents

# Cosine similarity/distance

- Given two points $\mathbf{x} = (\mathbf{x_1}, \dots, \mathbf{x_d})$ and $\mathbf{y} = (\mathbf{y_1}, \dots, \mathbf{y_d})$ in high-dimensional sphere, i.e. $\|x\|_2^2 = \|y\|_2^2 = 1$
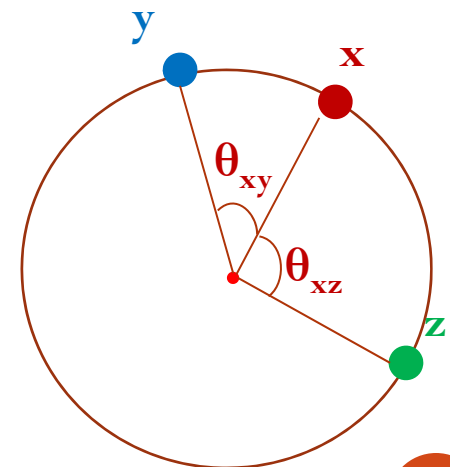
$$\|x\|_2^2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_d^2} = 1$$

- Cosine similarity:

$$\cos(\theta_{xy}) = \langle x, y \rangle = \sum_{i=1}^{d} x_i y_i$$

where $0 \leq \theta_{xy} \leq \pi$

- Cosine distance: $1 - \cos(\theta_{xy})$

# Applications

- Texts/Documents:
  - Present texts/documents as a high-dimensional points using term frequency–inverse document frequency (tf-idf).
  - L2 normalization, i.e. normalizing $\mathbf{x}$ by $(\mathbf{x_1}, \dots, \mathbf{x_d})/\|\boldsymbol{x}\|_2^2$.
  - Usage: Information retrieval (text search, topic modeling) and most of MinHash applications.
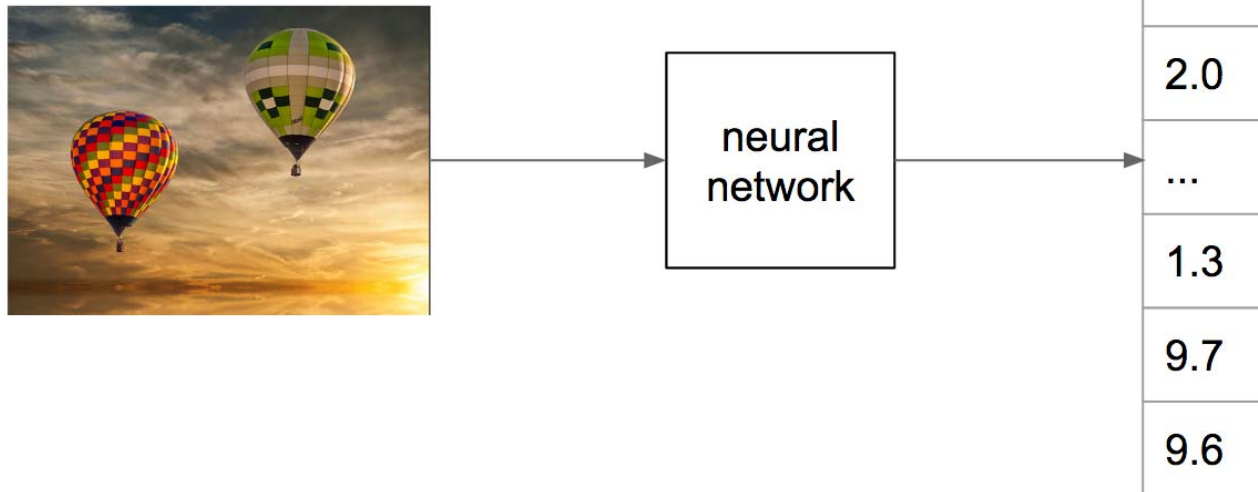
- A tf-idf definition for a corpus of $\mathbf{n}$ documents:

  - $\mathbf{tf(t, d)} = \dfrac{\text{number of times term } \mathbf{t} \text{ occur in document } \mathbf{d}}{\text{total number of terms in document } \mathbf{d}}$

  - $\mathbf{idf(t)} = \mathbf{log_2}\left(\dfrac{\mathbf{n}}{\mathbf{1} + \text{number of documents where } \mathbf{t} \text{ appears}}\right)$

  - $\mathbf{tf\text{-}idf(t, d) = tf(t, d) * idf(t)}$

# Applications

- Images:
  - Present an image as a high-dimensional vector, e.g. scale-invariant feature transform (SIFT), deep learning ("image2vec"),…
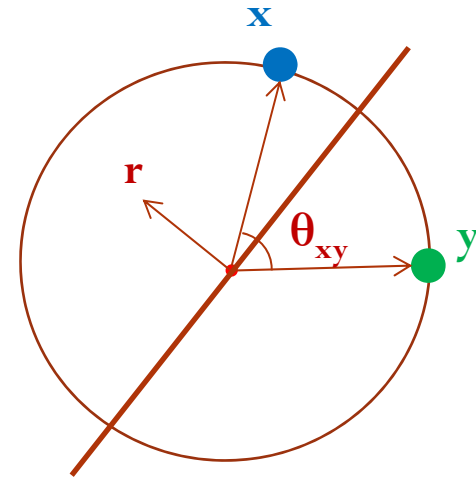  - Usage: content-based image retrieval



Source: https://code.flickr.net/2017/03/07/introducing-similarity-search-at-flickr/

11

# SimHash [Charikar'02]

- Generate a random Gaussian vector $\mathbf{r} = (\mathbf{r_1}, \dots, \mathbf{r_d})$ where $\mathbf{r_i}$ is drawn from $\mathbf{N(0, 1)}$

- SimHash function:

$$\mathbf{h_r(x)} = \begin{cases} \mathbf{1 \text{ if } \langle x, r \rangle \geq 0} \\ \\ \mathbf{0 \text{ if } \langle x, r \rangle < 0} \end{cases}$$

- $\mathbf{Pr_r [h_r(x) = h_r(y)]} = 1 - \theta_{xy} / \pi$
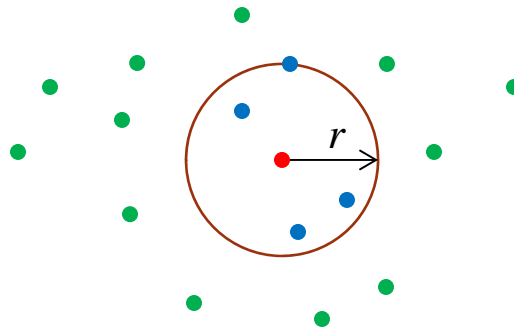
# Outline

- Popular similarity/distance measure

- Locality-sensitive hashing framework for approximate near neighbor search

13

# r-near neighbor search

- **r**-near neighbor search (**r**NNS) problem:
    - Given a data set $S \subset R^d$, a distance function $d(x, y)$, a distance threshold $r$, and a query $q \in R^d$, return some point $x \in S$ where $d(x, q) \le r$.
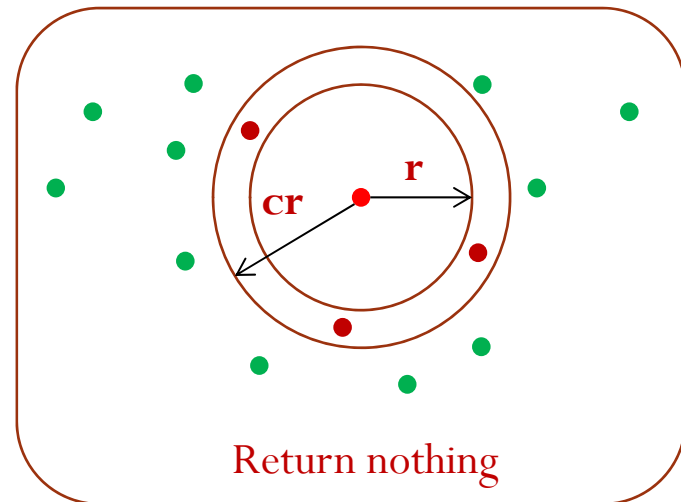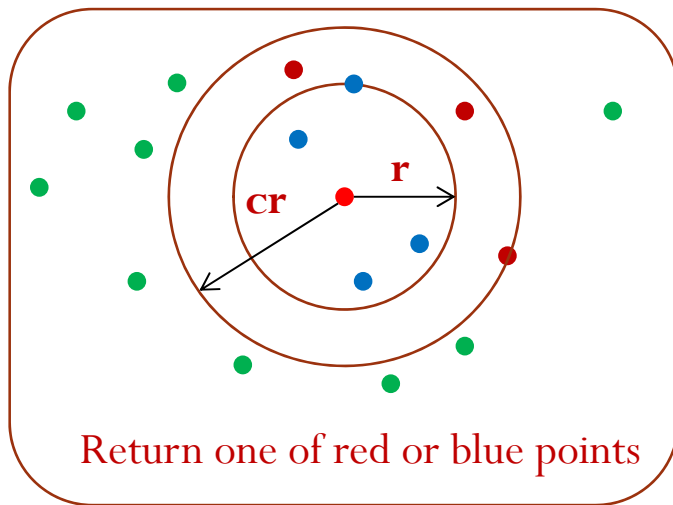


- Other variants:
    - **k-**nearest neighbors search
    - **r**-near neighbor reporting

# Problem relaxation
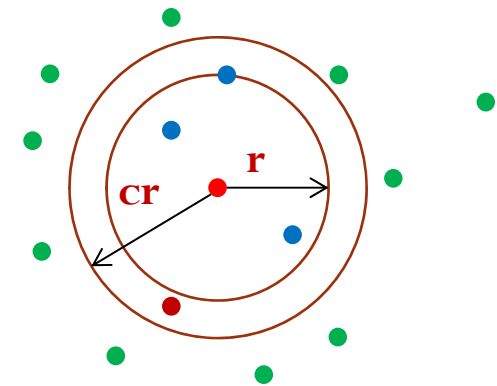
- Randomized **c**-approximation **r**NNS:
  - Given a data set $S \subset R^d$, a distance function $d(x, y)$, parameters $r > 0$, $c > 1$, $\delta > 0$, and a query $q \in R^d$, with probability $1 - \delta$
    - If exists $x \in S$ and $d(x, q) \leq r$, return some point $x' \in S$ where $d(x', q) \leq cr$.
    - Else, return nothing.



Return one of red or blue points

Return nothing

# LSH definition

- Definition [Indyk & Motwani'98]:

  - Given a distance function $d: U \times U \rightarrow R$ and positive values $r, c, p_1, p_2$ where $p_1 > p_2$, $c > 1$. A family of functions $H$ is called $(r, c, p_1, p_2)$-sensitive if for uniformly chosen $h \in H$ and all $x, y \in U$:

  - If $d(x, y) \leq r$ then $Pr [h(x) = h(y)] \geq p_1$; (similar/near points)

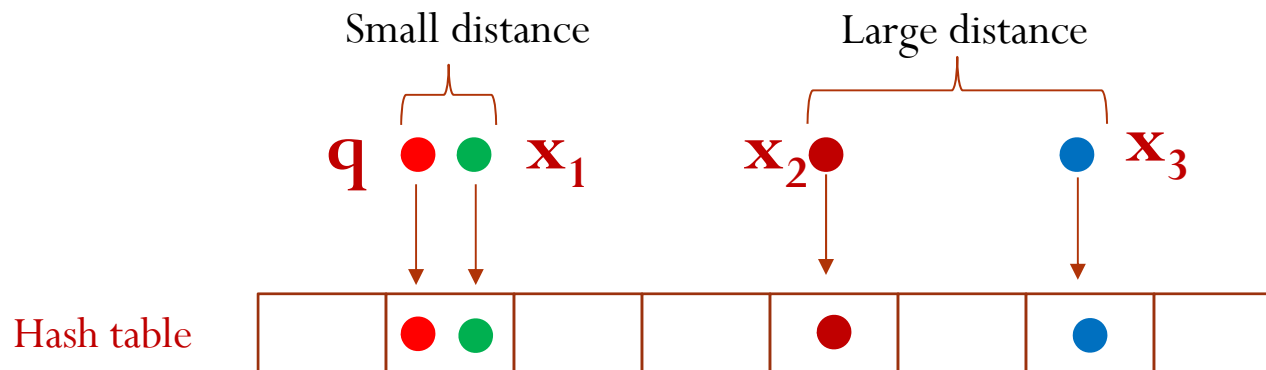  - If $d(x, y) \geq cr$ then $Pr [h(x) = h(y)] \leq p_2$. (dissimilar/far away points)

# Exercise

- MinHash is locality-sensitive hashing
  - $\Pr_\pi[h_\pi(S) = h_\pi(T)] = J(S,T)$
  - What are the values of $p_1$, $p_2$ given $r$ and $cr$?

- SimHash is locality-sensitive hashing
  - $\Pr_r [h_r(x) = h_r(y)] = 1 - \theta_{xy} / \pi$
  - What are the values of $p_1$, $p_2$ given $r$ and $cr$?

# Locality-sensitive hashing framework

- Claim:
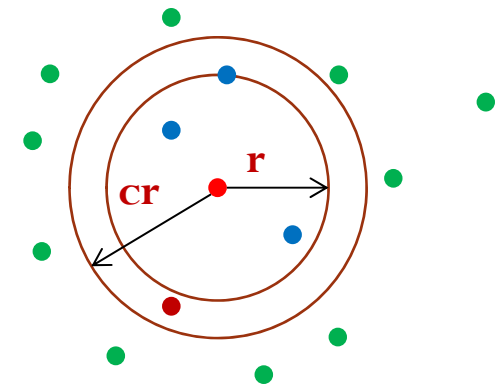    - We can solve the randomized **c**-approximation **r**-near neighbor search in sublinear time using locality-sensitive hashing.
    - Sublinear time: $O(n^\rho), 0 < \rho < 1$ for any query and data set of size **n**.

- Intuition: Candidate points are points within the same bucket of the query point.

Small distance          Large distance

q ● ● $x_1$     $x_2$●     ● $x_3$

Hash table

# Recap: LSH definition

- Definition [Indyk & Motwani'98]:

  - Given a distance function $d: U \times U \to R$ and positive values $r$, $c$, $p_1$, $p_2$ where $p_1 > p_2$, $c > 1$. A family of functions $H$ is called $(r, c, p_1, p_2)$-sensitive if for uniformly chosen $h \in H$ and all $x, y \in U$:

  - If $d(x, y) \leq r$ then $Pr\,[h(x) = h(y)] \geq p_1$; (similar/near points)

  - If $d(x, y) \geq cr$ then $Pr\,[h(x) = h(y)] \leq p_2$. (dissimilar/far away points)

How to make it as small as possible, i.e. $p_2 = 1/n.$

19

# Concatenate several hash functions

- Define a hash function $g(x) = (h_1(x), h_2(x), \ldots, h_k(x))$ where $h_i$, $1 \leq i \leq k$, are chosen independently and randomly from a LSH family $H$.

  - If $d(x, y) \leq r$ then $\mathbf{Pr}\ [g(x) = g(y)] \geq (p_1)^k$ (similar/near points)

  - If $d(x, y) \geq cr$ then $\mathbf{Pr}\ [g(x) = g(y)] \leq (p_2)^k$ (dissimilar/far away points)

- Set $k = \log(n) / \log(1/p_2)$, we have
  - $(p_2)^k = (p_2)^{\log(n)\ /\ \log(1/p2)} = 1/n$
  - $(p_1)^k = (p_1)^{\log(n)/\log(1/p2)} = 1/n^\rho$ where $0 < \rho = \dfrac{\log(1/p1)}{\log(1/p2)} < 1$

# LSH framework intuition

- Similar points:
    - If $d(x, y) \leq r$ then $\Pr [g(x) = g(y)] \geq 1/n^{\rho}$
    - If there is a similar point, i.e. $d(x, q) \leq r$, in expectation, we need to check $n^{\rho}$ points to find a similar point in our bucket.

- Dissimilar points:
    - If $d(x, y) \geq cr$ then $\Pr [g(x) = g(y)] \leq 1/n$
    - Since we have $n$ points, in expectation, only $1$ dissimilar point is in the same bucket of query.

# LSH framework intuition

- Similar points:
    - If $d(x, y) \leq r$ then $\mathbf{Pr}\ [g(x) = g(y)] \geq 1/n^\rho$
    - If there is a similar point, i.e. $d(x, q) \leq r$, in expectation, we need to check $n^\rho$ points to find a similar point in our bucket.

        How to get $n^\rho$ points?

- Dissimilar points:
    - If $d(x, y) \geq cr$ then $\mathbf{Pr}\ [g(x) = g(y)] \leq 1/n$
    - Since we have $n$ points, in expectation, only $1$ dissimilar point is in the same bucket of query.

        No worries on dissimilar points

# LSH framework intuition

- Similar points:

  - If $d(x, y) \leq r$ then $\Pr[g(x) = g(y)] \geq 1/n^\rho$

  - If there is a similar point, i.e. $d(x, q) \leq r$, in expectation, we need to check $n^\rho$ points to find a similar point in our bucket.

    Use $L = n^\rho$ hash tables
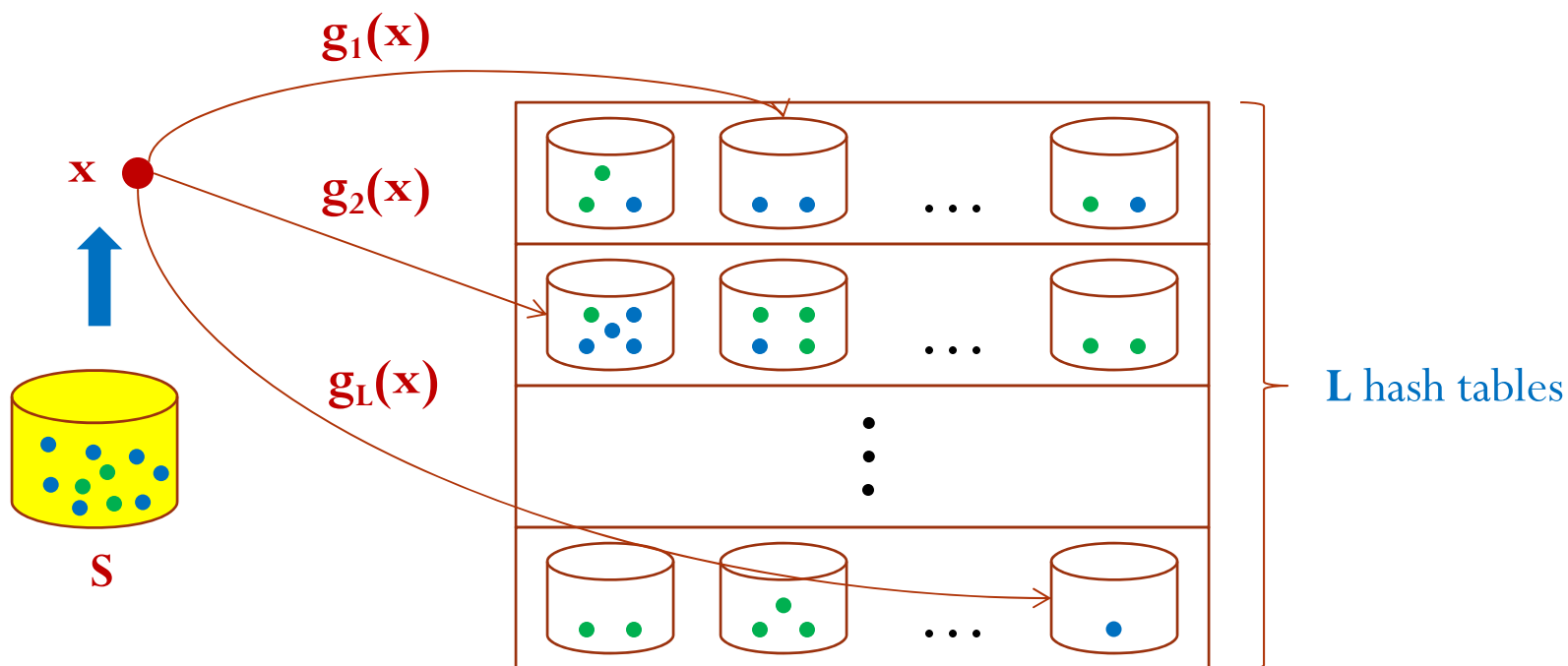
- Dissimilar points:

  - If $d(x, y) \geq cr$ then $\Pr[g(x) = g(y)] \leq 1/n$

  - Since we have $n$ points, in expectation, only $1$ dissimilar point is in the same bucket of query.

    In expectation, less than $L = n^\rho$ points

23

# LSH framework: Hash tables construction



$g_1(x)$

$g_2(x)$

$g_L(x)$

**x**

**S**

**L** hash tables

Space usage: $O(nL + dn)$

$g(x) = (h_1(x), h_2(x), \ldots, h_k(x))$ where $h_i \in H$

Prob. collision of near points: $1/n^\rho$

Prob. collision of far away points: $1/n$

# LSH framework: Hash tables construction

- Preprocessing:

  > 1) Choose L hash function $g_j$, $1 \leq j \leq L$, by setting $g_j = (h_{1,j}, h_{2,j}, \ldots, h_{k,j})$ where $h_{i,j}$, $1 \leq i \leq k$, are chosen independently and randomly from the LSH family H.
  >
  > 2) Construct L hash tables, where each hash table $T_j$, $1 \leq j \leq L$, contains the data set S hashed using the hash function $g_j$.
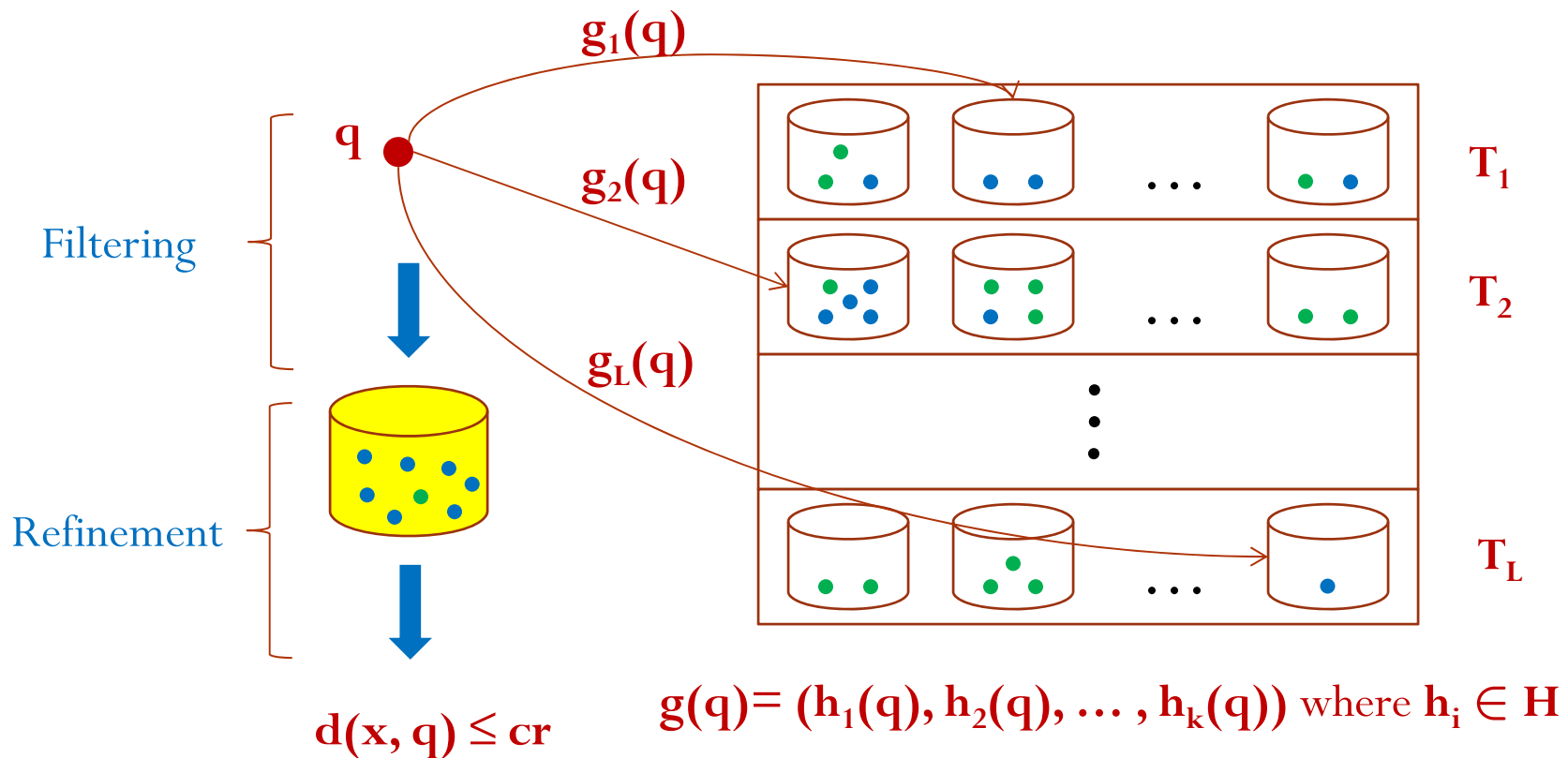
- Complexity:
  - Space usage: $\mathbf{O(nL + dn)}$
  - Time: $\mathbf{O(dknL)}$ with the assumption that the evaluation of $\mathbf{h(x)}$ takes $\mathbf{O(d)}$ time.

# LSH framework: Hash tables lookup

Important:

We interrupt search after finding first $\mathbf{L' = 2n^{\rho}}$ points, including duplicates.



$\mathbf{g(q) = (h_1(q), h_2(q), \dots, h_k(q))}$ where $\mathbf{h_i \in H}$

# LSH framework: Hash tables lookup

- Query for a point **q**:

  For each j = 1, 2, …, L
  1) Retrieve the points from the bucket $g_j(q)$ in the $T_j$ hash table.
  2) For each of the reported point x
      2.1) Compute the distance d(x, q).
      2.2) If d(x, q) ≤ cr, return x and stop.
  3) Stop as soon as the number of reported points is more than L'.

- Complexity:
    - Time: **O(dkL + dL')** with the assumption that the evaluation of **h(q)** takes **O(d)** time.

  Compute hash values     Compute actual distance with **L'** points

27

# Analysis

- Parameter settings:
    - $k = \log(n) / \log(1/p_2)$
    - $0 < \rho = \dfrac{\log(1/p1)}{\log(1/p2)} < 1$
    - $L = n^\rho$ hash tables
    - Check $L' = 2L = 2n^\rho$ points for computing actual distance

    > Prob. collision of near points: $1/n^\rho$
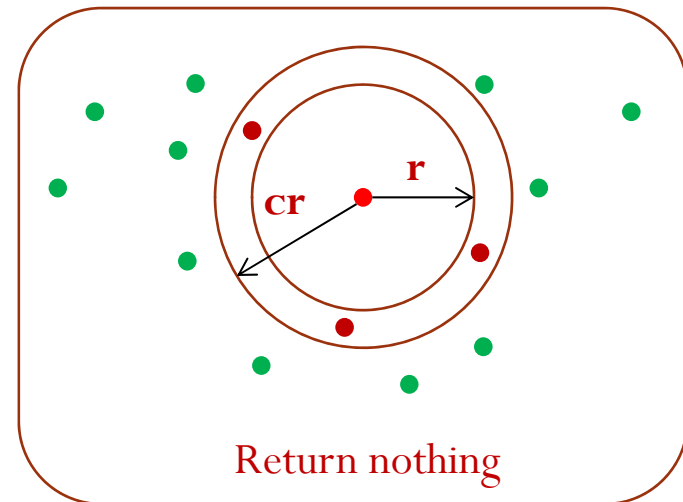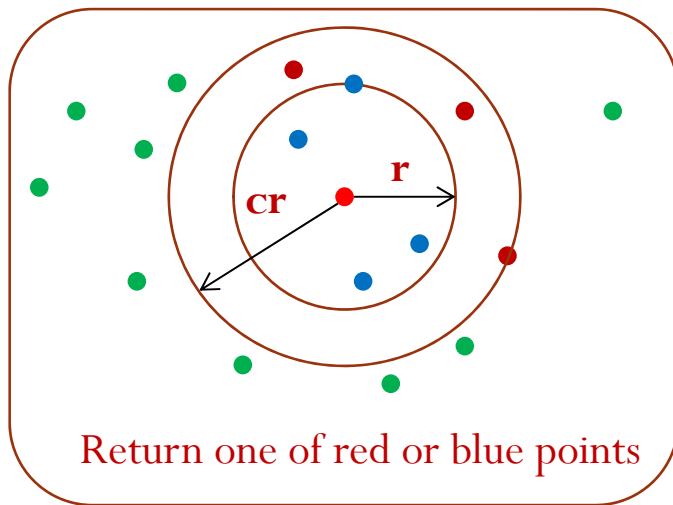    > Prob. collision of far away points: $1/n$

- Complexity:
    - Preprocessing:
        - Space usage: $O(nL + dn)$
        - Time: $O(dkn^{1+\rho})$ with the assumption that the evaluation of $h(x)$ takes $O(d)$ time.
    - Query:
        - Sublinear time: $O(dkn^\rho + dn^\rho)$ with the assumption that the evaluation of $h(q)$ takes $O(d)$ time.

28

# Recap: Problem relaxation

- Randomized **c**-approximation **r**NNS:
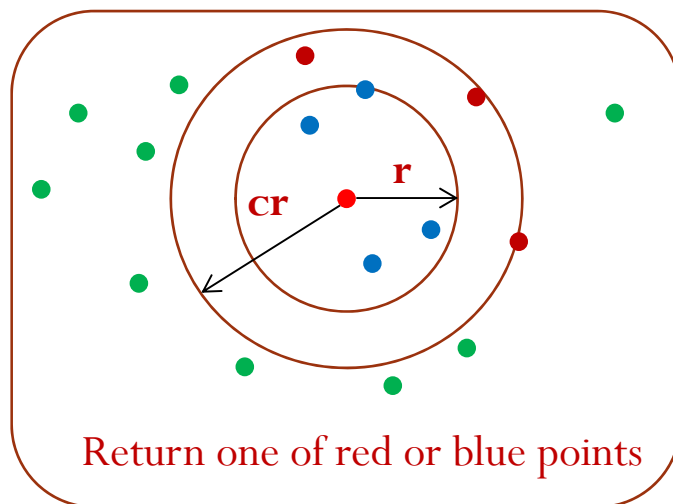  - Given a data set $S \subset R^d$, a distance function $d(x, y)$, parameters $r > 0$, $c > 1$, $\delta > 0$, and a query $q \in R^d$, with probability $1 - \delta$
    - If exists $x \in S$ and $d(x, q) \leq r$, return some point $x' \in S$ where $d(x', q) \leq cr$.
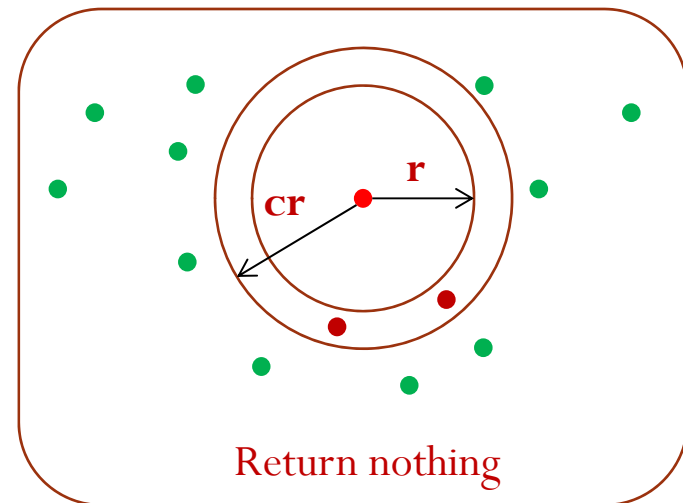    - Else, return nothing.



Return one of red or blue points

Return nothing

# Correctness

- We need to show that, with probability **1 – δ**:

  - Case 1: If exists $\mathbf{x} \in \mathbf{S}$ and $\mathbf{d(x, q)} \leq \mathbf{r}$, there exists a hash table $\mathbf{T_j}$ such that $\mathbf{g_j(x) = g_j(q)}$ for some $\mathbf{1 \leq j \leq L}$. This means that we can find $\mathbf{x}$.

  - Case 2: The total number of collisions from far away points is at most $\mathbf{2L}$. This means that after checking $\mathbf{2L}$ points, we are certain that all points are far away from $\mathbf{q}$ more than $\mathbf{cr}$.

Case 1



Return one of red or blue points

Case 2



Return nothing

30

# Correctness: Case 1

- Observation: If $\mathbf{d(x, q) \leq r, then}$

  - Prob. $\mathbf{x}$ colliding with $\mathbf{q}$ in any hash table: $\mathbf{\geq 1/n^{\rho}}$

  - Prob. $\mathbf{x}$ not colliding with $\mathbf{q}$ in any hash table: $\mathbf{\leq 1 - 1/n^{\rho}}$

  - Prob. $\mathbf{x}$ not colliding with $\mathbf{q}$ in all $\mathbf{L}$ hash tables: $\mathbf{\leq (1 - 1/n^{\rho})^{L}}$

  - Prob. $\mathbf{x}$ not colliding with $\mathbf{q}$ in all $\mathbf{L}$ hash tables: $\mathbf{1/e}$ with $\mathbf{L = n^{\rho}}$

  - Prob. $\mathbf{x}$ colliding with $\mathbf{q}$ in at least $\mathbf{1}$ hash tables: $\mathbf{1 - 1/e}$ with $\mathbf{L = n^{\rho}}$

  $(1 - 1/x)^x \approx 1/e$

- With prob. $\mathbf{1\text{-}1/e}$, any near point $\mathbf{x}$ will collide with $\mathbf{q}$ in at least $\mathbf{1}$ hash table. This means that we can find $\mathbf{x}$.
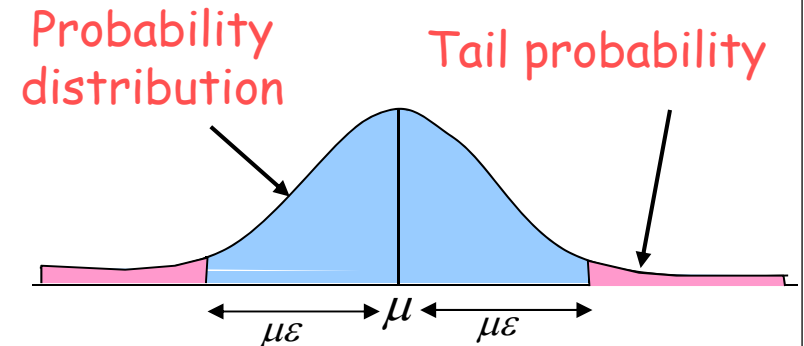
31

# Correctness: Case 2

- Observation:
  - Maximum number of far away points: **n**
  - If $d(x, q) \geq cr$, then $\Pr[h(x) = h(q)] \leq 1/n$. Hence, in expectation, we have at most **1** far away point collide with **q** in each hash table.
  - **E [number of far away points] ≤ L.**



Probability distribution    Tail probability

$$\Pr[X \geq a] \leq \frac{E[X]}{a} \text{ for any } a > 0.$$

- Using Markov's inequality:
  - Let **X** be the number of far away points collided with **q** in **L** tables.
  - $\Pr[X \geq 2L] \leq \dfrac{E[X]}{2L} \leq 1/2$

- With prob. **1/2**, after checking all first **L' = 2L** points far away from **q**, we are certain that all points are far away. Hence, return nothing.

32

# Correctness

- Probability both case 1 and case 2 hold is at least:

$$1 - ((1 - 1/2) + (1 - (1 - 1/e))) = 1/2 - 1/e$$

Prob. case 2 fails     Prob. case 1 fails

- Theorem: There exists a sublinear algorithm to answer the $c$-approximate $r$NNS with probability $1/2 - 1/e$.

- Boosting the accuracy: By repeating the LSH algorithm $O(1/\delta)$ times, we can amplify the probability of success in at least $1$ trial to $1 - \delta$ for any $\delta > 0$.

# Practical algorithm

- Query for a point **q**:

> For each j = 1, 2, ..., L
> 1) Retrieve the points from the bucket $g_j(q)$ in the $T_j$ hash table.
> 2) For each of the reported point x
>     2.1) Compute the distance d(x, q).
>     2.2) If d(x, q) ≤ r, return x.
>     2.3) Update the k-nearest neighbors with the returned x.

- We cannot have sublinear theoretical guarantee but the accuracy is much higher. In practice, it can run in sublinear time in many of data sets.