

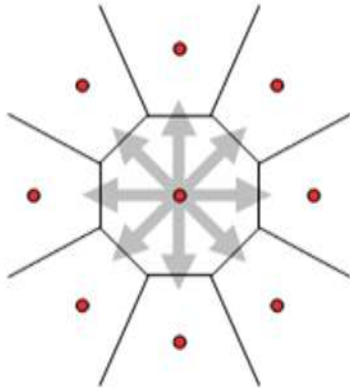
هشپنگ حساس به محل

سپیده آقاملائی

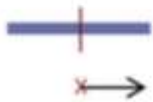
مشکلات ابعاد بالا

2

□ d همسایه ی مستقل در هر جهت



□ سرعت رشد بالای همسایگی با افزایش شعاع



$d=1$



$d=2$



$d=3$

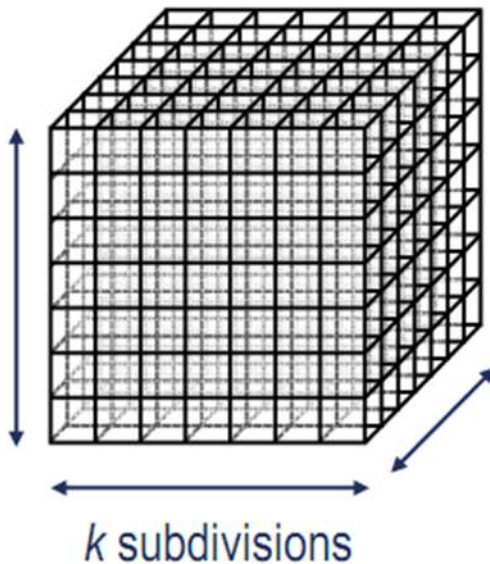
$$\text{vol}(r) \in \Theta(r^d)$$

$d \rightarrow \infty$

توری کمکی نمی کند.

3

- تکنیک های تقسیم معمولی کار نمی کنند.
- توری های $k \times k$ دارای k^d خانه هستند.
- توری های adaptive هم جواب نمی دهند.



□ دورترین همسایه \sim نزدیک ترین همسایه

□ لم JL

JL-Lemma: [Dasgupta et al. 99]

- Point set P in \mathbb{R}^d , $n := \#P$
- There is $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$, $k \in O(\varepsilon^{-2} \ln n)$
($k \geq 4(\varepsilon^2/2 - \varepsilon^3/3)^{-1} \ln n$)
- ...that preserves all inter-point distances up to a factor of $(1 + \varepsilon)$

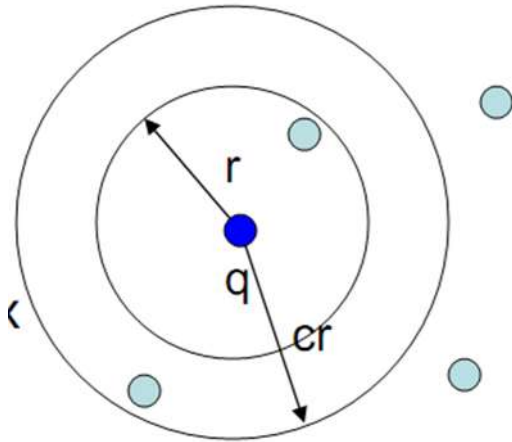
کاربرد لم JL

- می توانیم به صورت تصادفی جهت هایی را انتخاب کنیم و روی آنها تصویر عمودی نقاط را به دست بیاوریم.
- با احتمال بیشتر از $1-1/n$ جواب می دهد.
- لم JL به ما می گوید فاصله ی زوج نقاط در مجموعه کوچکی از نقاط (نمایی بر حسب تعداد ابعاد) می تواند در ابعاد پایین حفظ شود.
- لم JL نمی گوید که خود نقاط هم حفظ می شوند، فقط فاصله ی زوج نقاط حفظ می شود.

تقریب همسایه نزدیک

6

- با ضریب تقریب c همسایگی به شعاع r حول یک نقطه
- ساختمان داده ای بدهید که برای هر نقطه کوثری q
- اگر نقطه ای مثل p هست که $\|p-q\| \leq r$ باشد، نقطه p' را برگرداند که $\|p'-q\| \leq c.r$



کاهش از مساله های دیگر

- نزدیک ترین همسایه C -تقریبی به همسایه نزدیک C -تقریبی
 - با اضافه شدن یک \log
- با پیدا کردن همه ی همسایه های نزدیک می توان نزدیک ترین همسایه را دقیق حل کرد.
- کاربردها:
 - درخت پوشای کمینه C -تقریبی
 - خوشه بندی
 - ...

نتایج همسایه نزدیک تقریبی

8

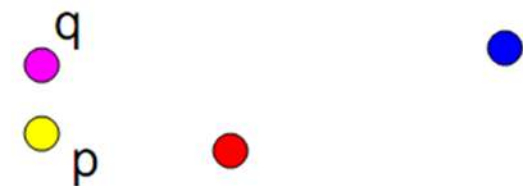
	Query time	Space used	Preprocessing time
Vornoi	$O(2^d \log n)$	$O(n^{d/2})$	$O(n^{d/2})$
Kd-tree	$O(2^d \log n)$	$O(n)$	$O(n \log n)$
LSH	$O(n^\rho \log n)$	$O(n^{1+\rho})$	$O(n^{1+\rho} \log n)$

کارهای گذشته

Space	Time	Comment	Norm	Ref
$dn+n^{4/\varepsilon^2}$	$d * \log n / \varepsilon^2$ or 1	$c=1+ \varepsilon$	Hamm, l_2	[KOR'98, IM'98]
$n^{\Omega(1/\varepsilon^2)}$	$O(1)$			[AIP'0?]
$\rightarrow dn+n^{1+\rho(c)}$	$d n^{\rho(c)}$	$\rho(c)=1/c$	Hamm, l_2	[IM'98], [Cha'02]
\rightarrow		$\rho(c)<1/c$	l_2	[DIIM'04]
$dn * \log s$	$dn^{\sigma(c)}$	$\sigma(c)=O(\log c/c)$	Hamm, l_2	[Ind'01]
		$\sigma(c)=O(1/c)$	l_2	[Pan'06]
$\rightarrow dn+n^{1+\rho(c)}$	$d n^{\rho(c)}$	$\rho(c)=1/c^2 + o(1)$	l_2	[Al'06]
$dn * \log s$	$dn^{\sigma(c)}$	$\sigma(c)=O(1/c^2)$	l_2	[Al'06]

هشینگ حساس به محل (LSH)

10

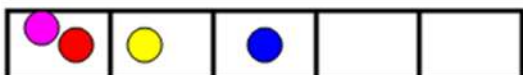
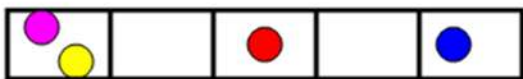


□ ایده: توابع هش $g: R^d \rightarrow U$ را برای هر زوج نقاط p و q تعریف می کنیم:

□ اگر $\|p-q\| \leq r$ آنگاه $\Pr(g(p)=g(q))$ خیلی کوچک نباشد.

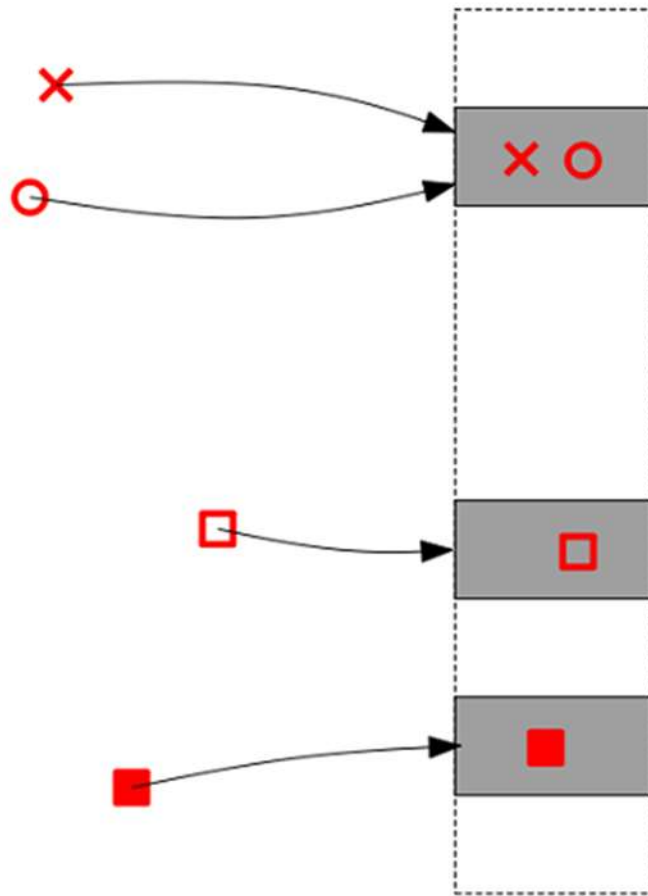
□ اگر $\|p-q\| > cr$ آنگاه $\Pr(g(p)=g(q))$ کوچک باشد.

□ در این صورت با هش کردن مساله را می توان حل کرد.



تصویر کلی عملکرد LSH

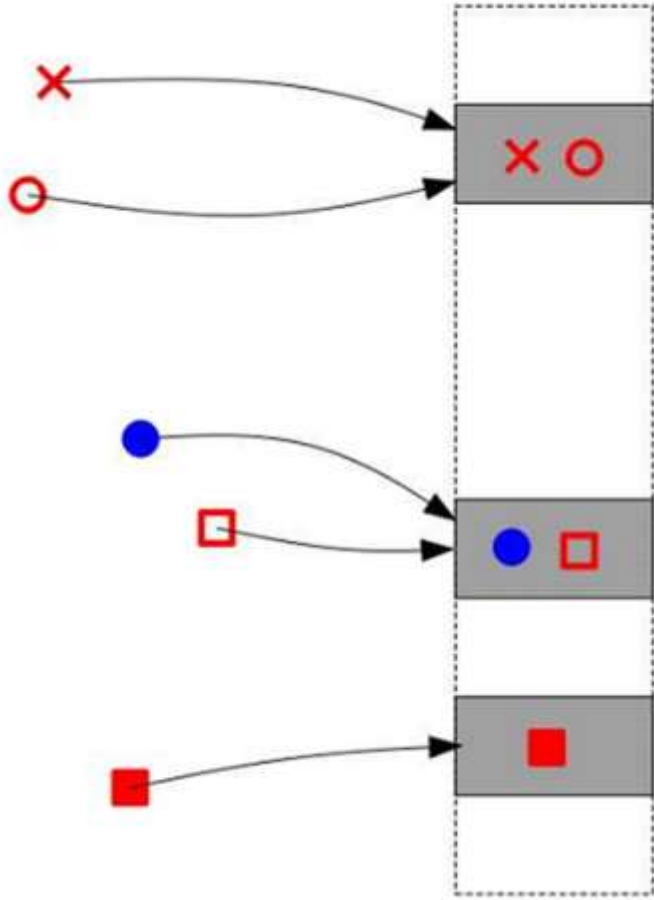
11



- پیش پردازش
- خانواده ای از توابع هش که نقاط نزدیک به یک سطل و نقاط دور به سطل های مجزا بروند.
- یک تابع تصادفی انتخاب کن و نقاط ورودی را هش کن.
- فقط سطل های ناتهی را نگهداری کن.

تصویر کلی عملکرد LSH

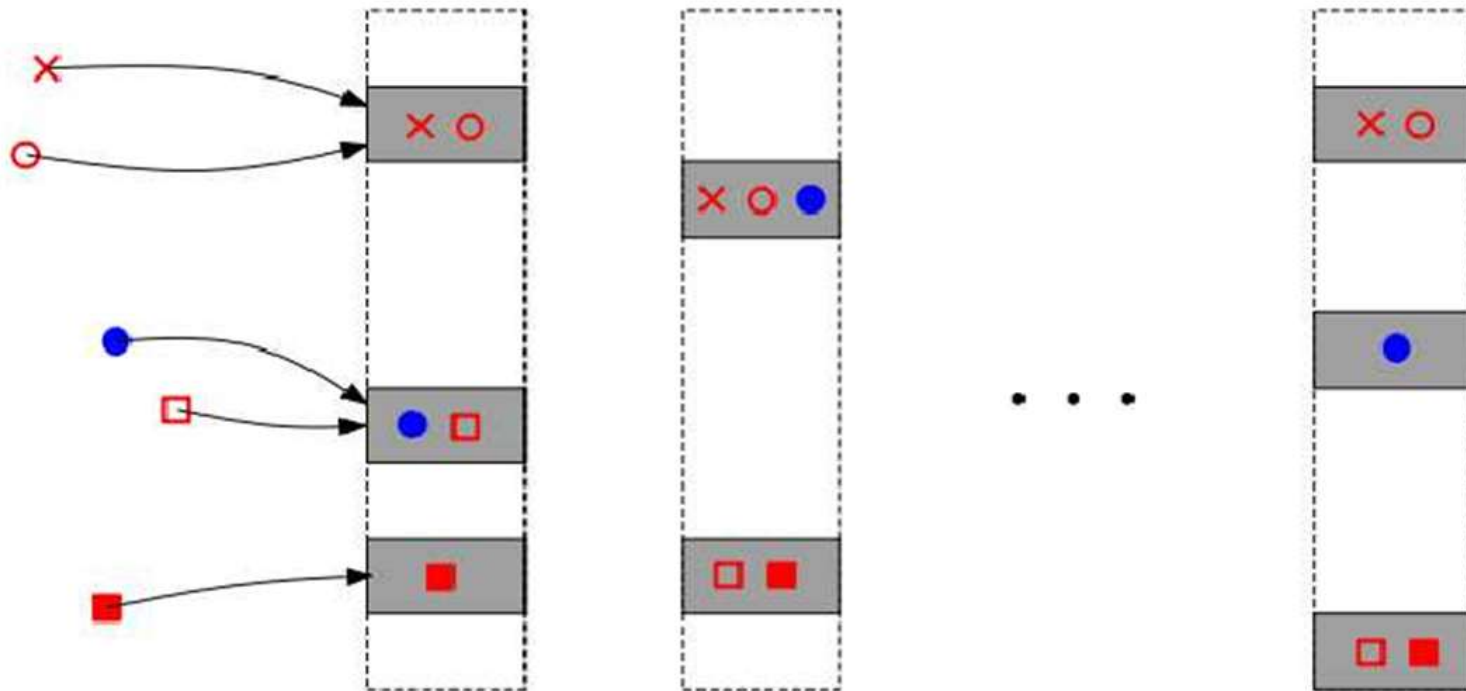
12



- کوئری
- سطل متناظر هر نقطه q را برای همسایه تقریبی نزدیکش بگرد.
- ایراد: ممکن است سطل q خالی باشد.

تصویر کلی عملکرد LSH

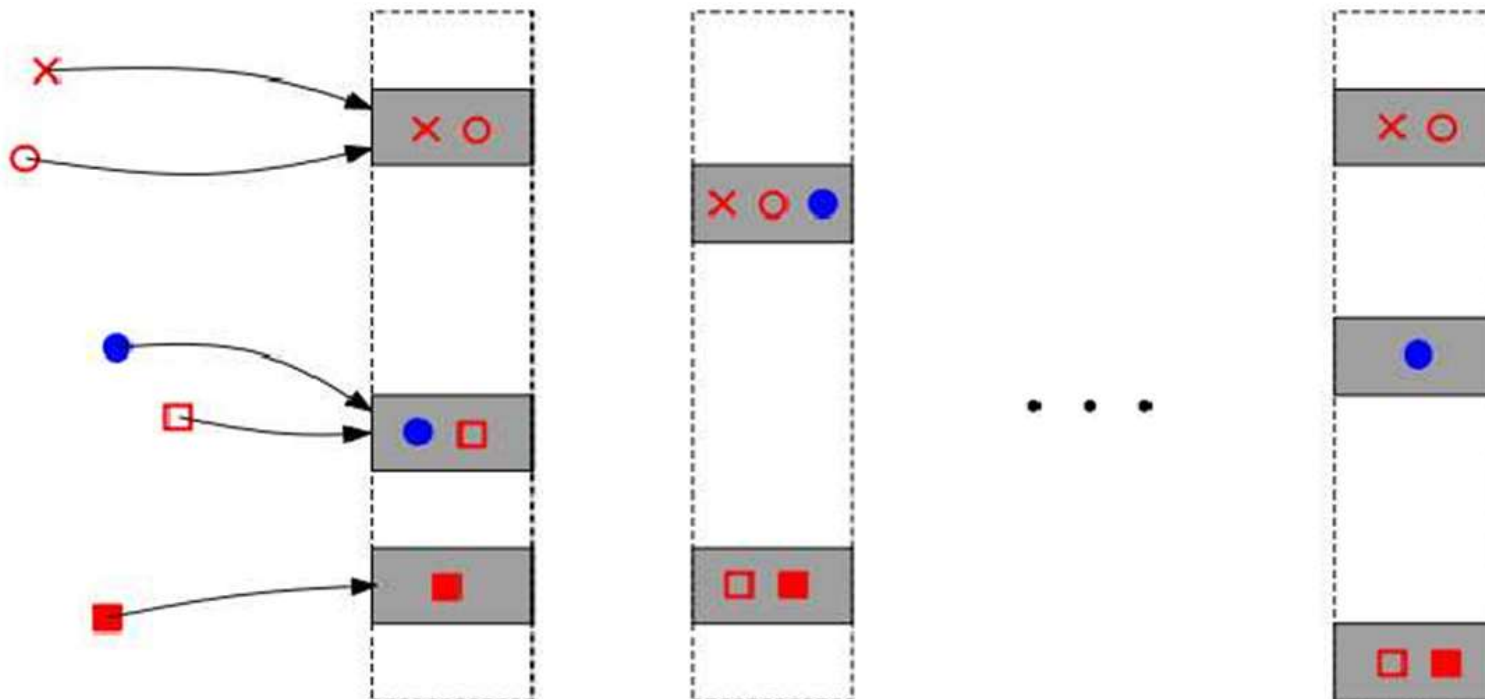
- راه حل: از چند جدول هش استفاده کن.
- اگر هر همسایگی نزدیک تقریبی پیدا شد، کار تمام است.



تصویر کلی عملکرد LSH

14

- ایراد: دقت کافی ندارد، چون کاندیداهای زیادی پیدا می شوند.
- راه حل: بعد از رسیدن به احتمال کوچک مشخصی متوقف شود.



تصویر کلی عملکرد LSH

□ پیدا کردن تابع هش

If $u \in B(q, r)$ then $\Pr[h(u) = h(q)] \geq \alpha$

If $u \notin B(q, R)$ then $\Pr[h(u) = h(q)] \leq \beta$

$r < R, \alpha \gg \beta$

□ که در آن

□ h به صورت تصادفی از یک خانواده تابع هش انتخاب شده است.

□ R به صورت روبرو به دست آمده است:

$$R = r(1 + \varepsilon)$$

- A family H of functions $h: \mathbb{R}^d \rightarrow U$ is called (P_1, P_2, r, cr) -sensitive, if for any p, q :
 - if $\|p-q\| < r$ then $\Pr[h(p)=h(q)] > P_1$
 - if $\|p-q\| > cr$ then $\Pr[h(p)=h(q)] < P_2$
- Example: Hamming distance
 - LSH functions: $h(p)=p_i$, i.e., the i -th bit of p
 - Probabilities: $\Pr[h(p)=h(q)] = 1-D(p,q)/d$

$p=10010010$

$q=11010110$

الگوریتم LSH

17

- توابعی به شکل $g(p) = \langle h_1(p), h_2(p), \dots, h_k(p) \rangle$
- پیش پردازش:
- g_1, \dots, g_L را انتخاب کنید.
- به ازای هر نقطه p ، آن را به سطلهای $g_1(p), \dots, g_L(p)$ هش کنید.
- کوئری:
- از سطلهای $g_1(p), g_2(p), \dots$ نقاط را بازیابی کنید تا
 - همه ی نقاط از L سطل بازیابی شده باشند.
 - تعداد نقاط بازیابی شده $2L$ شود.
- زمان کوئری به تعداد نقاط بازیابی شده بستگی دارد.
- زمان کل $O(dL)$

فاصله همینگ و مکعب d -بعدی یکه

- نقاط درون $H^d = \{0,1\}^d$ قرار می گیرند.
- هر نقطه یک رشته باینری است.
- فاصله همینگ (r) : تعداد مختصات متفاوت

$$\begin{array}{l} u : 0 \boxed{0} 1 0 \boxed{1} \boxed{1} 1 0 \boxed{1} \\ v : 0 \boxed{1} 1 0 \boxed{0} \boxed{0} 1 0 \boxed{0} \end{array}$$

- Define family F :

Given : Hypercube H^d , point $b = (b_1, \dots, b_d)$

$$h \in F : \left\{ h_i(b) = b_i \mid b = (b_1, \dots, b_d) \in H^d, \text{ for } i=1, \dots, d \right\}$$

$$\alpha = 1 - \frac{r}{d}, \quad \beta = 1 - \frac{r(1 + \varepsilon)}{d}$$

- Intuition: compare a random coordinate
- Called: $(r, r(1 + \varepsilon), \alpha, \beta)$ -sensitive family

- **Define family G :**

Given : $b \in H^d, F$

$g \in G$:

$$\left\{ g: \{0,1\}^d \rightarrow \{0,1\}^k \mid g(b) = \left(h^1(b), \dots, h^k(b) \right), \text{ for } h^i \in F \right\}$$

$$\alpha' = \left(1 - \frac{r}{d} \right)^k = \alpha^k, \quad \beta' = \left(1 - \frac{r(1+\varepsilon)}{d} \right)^k = \beta^k$$

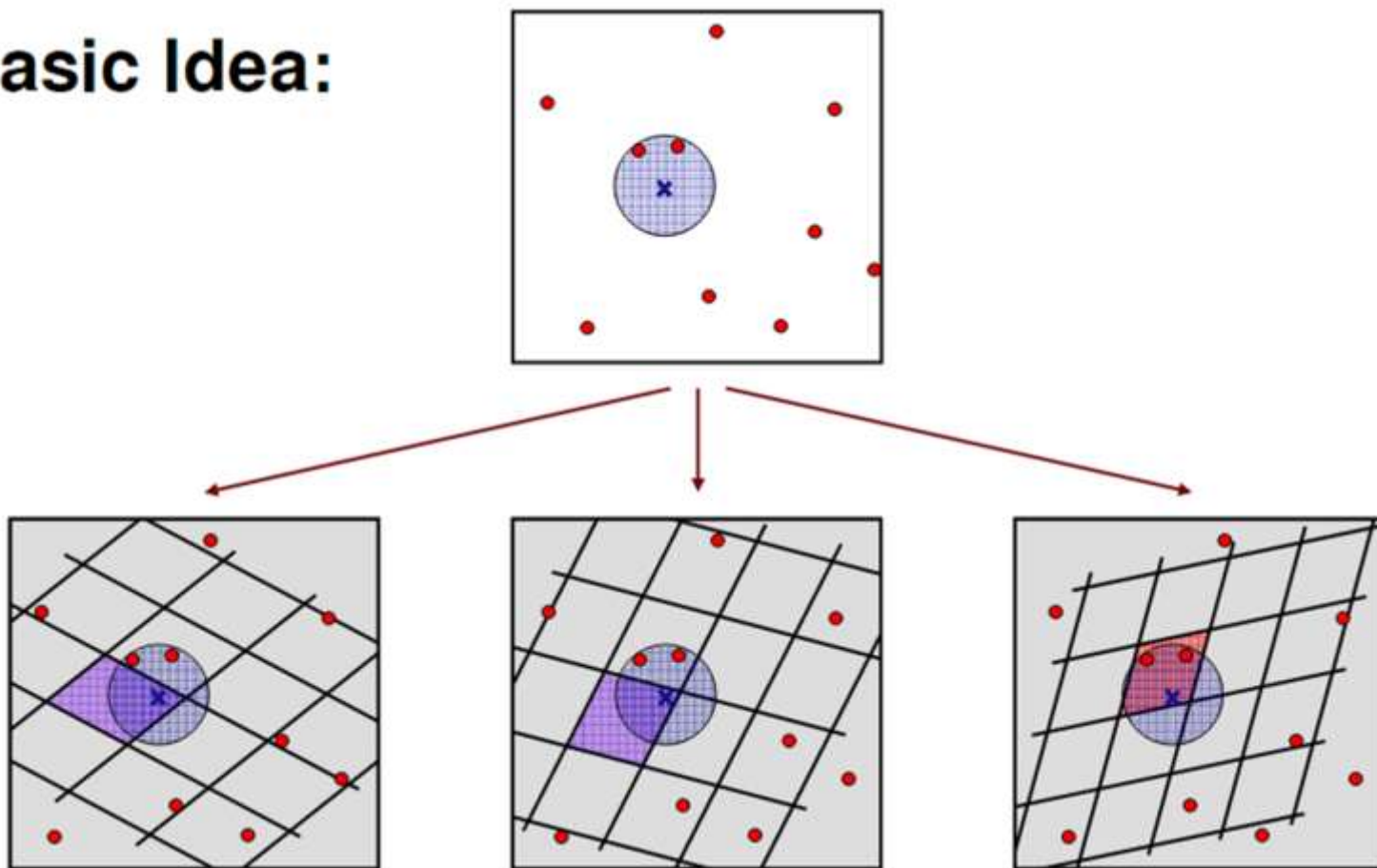
- Intuition: Compare k random coordinates
- Choose k later – logarithmic in n J-L lemma

- LSH solves c -approximate NN with:
 - Number of hash fun: $L=n^\rho$, $\rho=\log(1/P1)/\log(1/P2)$
 - E.g., for the Hamming distance we have $\rho=1/c$
 - Constant success probability per query q
- Questions:
 - Can we extend this beyond Hamming distance ?
 - Yes:
 - embed l_2 into l_1 (random projections)
 - l_1 into Hamming (discretization)
 - Can we reduce the exponent ρ ?

LSH با روش تصویر کردن تصادفی

22

Basic Idea:



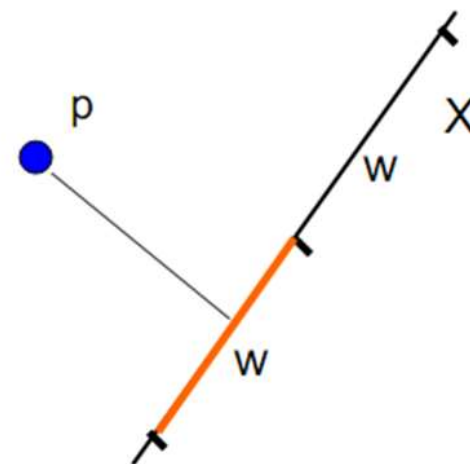
تصویر کردن تصادفی

- از چه توزیع هایی باید برای انتخاب بردارهای تصویر کردن استفاده کنیم؟
- اندازه سطل مناسب چیست؟
- حساسیت به محل
- چه تعداد خط در هر توری باشد؟
- چه تعداد خانه در کل باشد؟
- به حساسیت بستگی دارد.
- کارایی این روش چقدر است؟

LSH بر مبنای تصویر کردن

24

- Define $h_{X,b}(p) = \lfloor (p \cdot X + b) / w \rfloor$:
 - $w \approx r$
 - $X = (X_1 \dots X_d)$, where X_i is chosen from:
 - Gaussian distribution (for l_2 norm)
 - “s-stable” distribution* (for l_s norm)
 - b is a scalar
- Similar to the $l_2 \rightarrow l_1 \rightarrow$ Hamming route



* I.e., $p \cdot X$ has same distribution as $\|p\|_s Z$, where Z is s-stable

توزیع های p -stable

25

A prob. distribution D is called p -stable $:\Leftrightarrow$

- For any $v_1, \dots, v_n \in \mathbb{R}$
- And i.i.d. random variables $X_1, \dots, X_n \sim D$

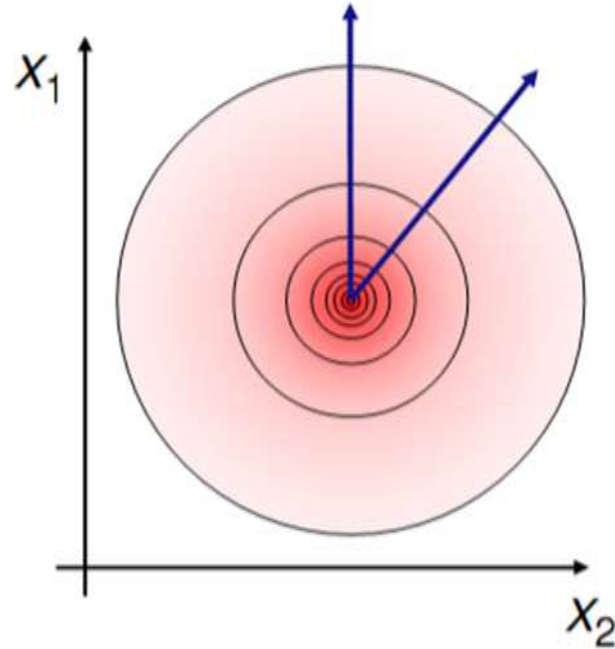
$\sum_i v_i X_i$ has the same distribution as $\left[\sum_i |v_i|^p \right]^{1/p} X$

where $X \sim D$

توزیع گاوسی

26

□ توزیع نرمال 2-stable است.



الگوریتم تصویر کردن

27

- Chose p according to metric of the space l_p
- Compute vector with entries according to a p -stable distribution
[for example: Gaussian noise entries]

- Each vector v_i yields a hash function h_i

- Compute: $h_i(x) = \left\lfloor \frac{\langle v_i, x \rangle + b}{r} \right\rfloor$
random value $\in [0 \dots r]$
bucket size

محاسبه حساسیت به محل برای تصویر کردن

28

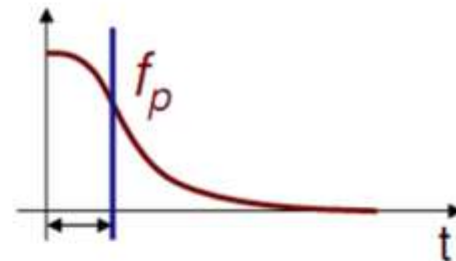
Computing the Locality “Sensitivity”

Distance $c = \|v_1 - v_2\|_p$

cX -distributed, X from p -stable distr.

$$\underbrace{\Pr(\text{collision})}_{=: p(c)} = \int_0^r \underbrace{\frac{1}{c}}_{\text{abs. density}} \underbrace{f_p\left(\frac{t}{c}\right)}_{\text{hit}} \underbrace{\left(1 - \frac{t}{r}\right)}_{\text{bucket}} dt$$

$$\left[h_i(x) = \left\lfloor \frac{\langle v_i, x \rangle + b}{r} \right\rfloor \right]$$



The constructed family of hash functions is

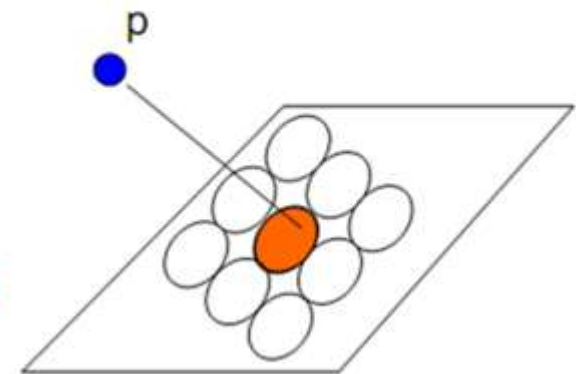
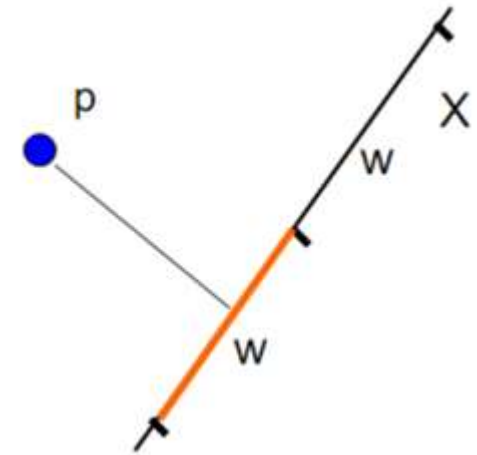
$(r_1, r_2, \alpha, \beta)$ -sensitive for

$$\alpha = p(1), \beta = p(c), r_2/r_1 = c$$

روش جدید LSH

29

- Instead of projecting onto \mathbb{R}^1 , project onto \mathbb{R}^t , for constant t
- Intervals \rightarrow lattice of balls
 - Can hit empty space, so hash until a ball is hit
- Analysis:
 - $\rho = 1/c^2 + O(\log t / t^{1/2})$
 - Time to hash is $t^{O(t)}$
 - Total query time: $dn^{1/c^2 + o(1)}$
- [Motwani-Naor-Panigrahy'06]: LSH in l_2 must have $\rho \geq 0.45/c^2$



MinHash

30

- =Min-wise independent permutations
- Let h be a hash function that maps the members of A and B to distinct integers, and for any set S define $\text{hmin}(S)$ to be the member x of S with the minimum value of $h(x)$.
- Then $\text{hmin}(A) = \text{hmin}(B)$ exactly when the minimum hash value of the union $A \cup B$ lies in the intersection $A \cap B$. Therefore,
- $\Pr[\text{hmin}(A) = \text{hmin}(B)] = J(A,B)$.

Min-Hashing

31

□ Consider

- $S_A, S_B \subset U$
- Pick – random permutation π of U
- Define $\alpha = \pi^{-1}(\min\{\pi(S_A)\})$ and $\beta = \pi^{-1}(\min\{\pi(S_B)\})$
- Meaning? – minimal element under permutation π

$$P[\alpha = \beta] = \frac{|S_A \cap S_B|}{|S_A \cup S_B|}$$

□ Lemma:

- Let $\delta = \min\{\pi(S_A \cup S_B)\}$
- Claim: $\alpha = \beta \Leftrightarrow \pi^{-1}(\delta) \in S_A \cap S_B$
- Clearly

$$P[\pi^{-1}(\delta) \in S_A \cap S_B] = \frac{|S_A \cap S_B|}{|S_A \cup S_B|}$$

Min-Wise Indep Permutations

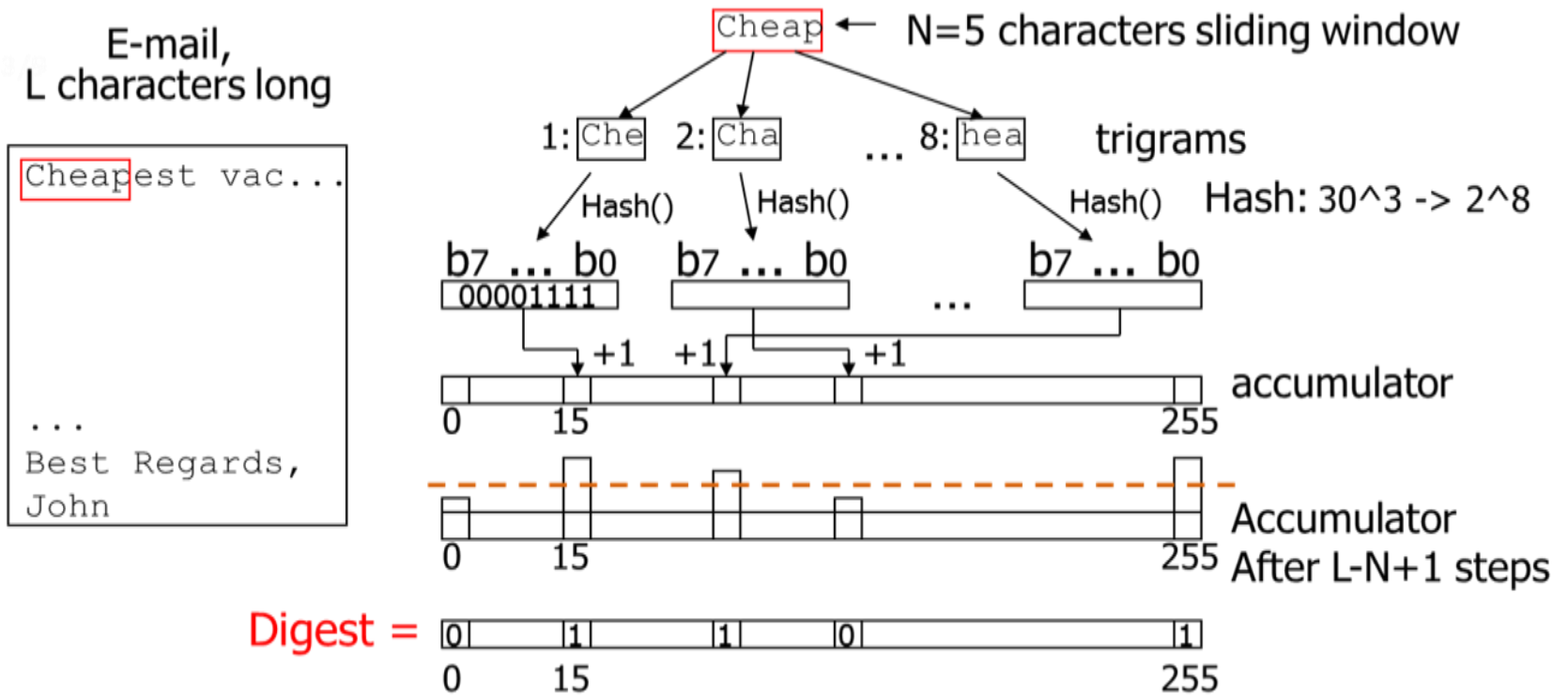
32

- Problem
 - ▣ Truly-random π over $U = [0 \dots N-1]$ is infeasible
 - ▣ But – do we really need true randomness?
- Solution
 - ▣ Poly-size family of permutations $F \subset S_N$ over U
 - ▣ Choosing/representing random $\pi \in F$ is easy
- Min-Wise Independence (MWI) Property:
 - ▣ For all sets $X \subset U$, for all $x \in F$,

$$P_{\pi \in F}[\min(\pi(X)) = x] = \frac{1}{|X|}$$

Nilsimsa similarity hashing

- Definition: Nilsimsa Compare Value (NCV) between two digests is equal to the number of bits at corresponding positions that are equal, minus 128.



Methods

34

- Bit sampling for Hamming distance
- Min-wise independent permutations
 - ▣ Jaccard index
- Nilsimsa Hash
- Random projection
 - ▣ Cosine distance
- Stable distributions

Conclusion

35

- Hamming: 1 dimension from d dimensions
- Jaccard: k dimensions from d dimensions
 - ▣ Binary values only
- Random Projection (Cosine): k dimensions (optional)
- Nilsimsa: k consecutive dimensions
 - ▣ Alphabet
- Stable Distributions: k dimensions