

# Locality Sensitive Hashing Scheme Based on $p$ -Stable Distributions

Mayur Datar (Stanford)

Nicole Immorlica (MIT)

Piotr Indyk (MIT)

Vahab Mirrokni (MIT)

## (Streaming) Massive Data Sets $\Rightarrow$ High Dimensional Vectors

- Massive data sets visualized as high dimensional vectors
- E.g. Number of IP-packets sent to address  $i$  from IP address  $j$

$$\mathbf{v}^j = \{v_1^j, v_2^j, \dots, v_i^j, \dots, v_N^j\}$$

Dimensionality =  $2^{32}$

- E.g. Number of phone calls made from telephone number  $j$  to telephone number  $k$

$$\mathbf{v}^j = \{v_1^j, v_2^j, \dots, v_k^j, \dots, v_{N'}^j\}$$

Dimensionality =  $10^9$

## Update Model

- Vectors constantly updated as per *cash register model*
- Update element  $(i, a)$  for vector  $v$  changes it as follows:

$$v = \{v_1, v_2, \dots, (v_i + a), \dots, v_N\}$$

- Numerous high dimensional vectors  
E.g. one vector per (millions) telephone customers,  
one vector per (millions) IP-address etc.  
Rows of a huge matrix

## $l_p$ Norms

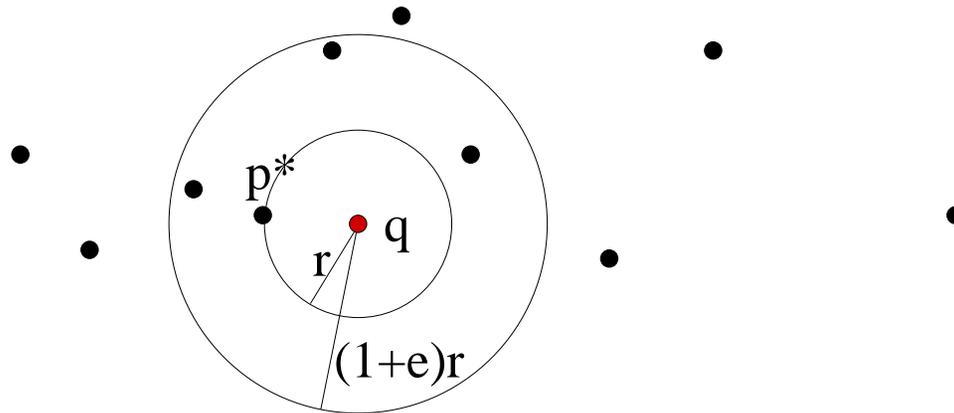
- $l_p(\mathbf{v}) = (\sum_{i=1}^N |v_i|^p)^{1/p}$   
E.g.  $l_1$  norm (Manhattan),  $l_2$  norm (Euclidean)
- $l_p$  norms usually computed over vector differences  
E.g.  $l_1(\mathbf{v}^j - \mathbf{v}^k)$ ,  $l_2(\mathbf{v}^j - \mathbf{v}^k)$ ,  $l_{0.005}(\mathbf{v}^j - \mathbf{v}^k)$  etc.
- What do  $l_p$  norms capture?
  - $l_1$  norm applied to telephone vectors: symmetric (multi) set difference between two customers
  - $l_p$  norms for small values of  $p$  (0.005): capture Hamming norms, distinct values [CDIM'02]

## Proximity Queries

- **Nearest Neighbor:** Given a query  $q$  find the closest (smallest  $l_p$  norm) point  $p$
- **Near Neighbor:** Given a query  $q$  and distance  $R$  find all (or most) points  $p$  s.t.  $l_p(p - q) \leq R$
- Applications: Classification, fraud detection etc.  
E.g. find cell phone customers whose calling pattern is similar to that of XYZ (UBL)

## Approximate Nearest Neighbor

- Curse of dimensionality
- Error parameter  $\epsilon$ : Find any point that is within  $(1 + \epsilon)$  times the distance from true nearest neighbor



## Approximate Near Neighbor $((R, \epsilon)$ -PLEB)

- $B(c, R)$  denotes a *ball* of radius  $R$  centered at  $c$
- Given: radius  $R$ , error parameter  $\epsilon$  and query point  $q$ :
  - if there exists data point  $p$  s.t.  $q \in B(p, R)$ , return YES and a point (or all points)  $p'$  s.t.  $q \in B(p', (1 + \epsilon)R)$ ,
  - if  $q \notin B(p, R)$  for all data points  $p$ , return NO,
  - if closest data point to  $q$  is at distance between  $R$  and  $R(1 + \epsilon)$  then return YES or NO

## Approximate Near Neighbor

- Useful problem formulation in itself
- Approximate *nearest* neighbor can be reduced to approximate near neighbor (binary search on  $R$ )
- Henceforth, we will concentrate on solving approximate near neighbor

## Our contribution

- Data structure for the approximate near neighbor problem  $((R, \epsilon)$ -PLEB)
- Small query time, **update time** and easy to implement
- works for  $l_p$  norms, for  $0 < p \leq 2$ . In particular  $0 < p < 1$
- Earlier result ([IM'98]) worked for  $l_1$ ,  $l_2$  and Hamming norm.
- Our technique improves the query time for  $l_2$  norm

## Locality Sensitive Hashing (LSH) ([IM'98])

- Intuition: if two points are close (less than dist  $r_1$ ) they hash to same bucket with prob at least  $p_1$ . Else, if they are far (more than dist  $r_2 > r_1$ ) they hash to same bucket with prob no more than  $p_2 < p_1$
- Formally: A family  $\mathcal{H} = \{h : S \rightarrow U\}$  is called  $(r_1, r_2, p_1, p_2)$ -sensitive for distance function  $D$  if for any  $\mathbf{v}, \mathbf{q} \in S$ 
  - if  $\mathbf{v} \in B(\mathbf{q}, r_1)$  then  $\Pr_{\mathcal{H}}[h(\mathbf{q}) = h(\mathbf{v})] \geq p_1$ ,
  - if  $\mathbf{v} \notin B(\mathbf{q}, r_2)$  then  $\Pr_{\mathcal{H}}[h(\mathbf{q}) = h(\mathbf{v})] \leq p_2$ .
  - $r_1 < r_2, p_1 > p_2$

## Using LSH to solve $(R, \epsilon)$ -PLEB ([IM'98])

- Let  $c = 1 + \epsilon$

**Theorem.** *Suppose there is a  $(R, cR, p_1, p_2)$ -sensitive family  $\mathcal{H}$  for a distance measure  $D$ . Then there exists an algorithm for  $(R, c)$ -PLEB under measure  $D$  which uses  $O(dn + n^{1+\rho})$  space, with query time dominated by  $O(n^\rho)$  distance computations, and  $O(n^\rho \log_{1/p_2} n)$  evaluations of hash functions from  $\mathcal{H}$ , where  $\rho = \frac{\ln 1/p_1}{\ln 1/p_2}$*

- Bottom-line: Design LSH scheme with small  $\rho$  for  $l_p$  norms

## Recap

- Proximity problems reduced to designing LSH schemes
- Design LSH schemes for  $l_p$  norms with small  $\rho$ , update time etc.
- A family  $\mathcal{H} = \{h : S \rightarrow U\}$  is called  $(r_1, r_2, p_1, p_2)$ -sensitive for distance function  $D$  if for any  $\mathbf{v}, \mathbf{q} \in S$ 
  - if  $\mathbf{v} \in B(\mathbf{q}, r_1)$  then  $\Pr_{\mathcal{H}}[h(\mathbf{q}) = h(\mathbf{v})] \geq p_1$ ,
  - if  $\mathbf{v} \notin B(\mathbf{q}, r_2)$  then  $\Pr_{\mathcal{H}}[h(\mathbf{q}) = h(\mathbf{v})] \leq p_2$
- $r_1 = R = 1$ ,  $r_2 = R(1 + \epsilon) = 1 + \epsilon = c$

## $p$ -Stable distributions

- **$p$ -stable distribution ( $p \geq 0$ ):** A distribution  $\mathcal{D}$  over  $\Re$  s.t.
  - $n$  real numbers  $v_1 \dots v_n$ ,
  - i.i.d. variables  $X_1 \dots X_n$  with distribution  $\mathcal{D}$ ,
  - r.v.  $\sum_i v_i X_i$  has the same distribution as the variable  $(\sum_i |v_i|^p)^{1/p} X = l_p(\mathbf{v})X$ , where  $X$  is a r.v. with distribution  $\mathcal{D}$
- E.g.  $p$ -Stable distr for  $p = 1$  is Cauchy distr, for  $p = 2$  is Gaussian distr
- for  $0 < p < 2$  there is a way to sample from a  $p$ -stable distribution given two uniform r.v.'s over  $[0, 1]$  [NoI]

## How are $p$ -Stable distributions useful?

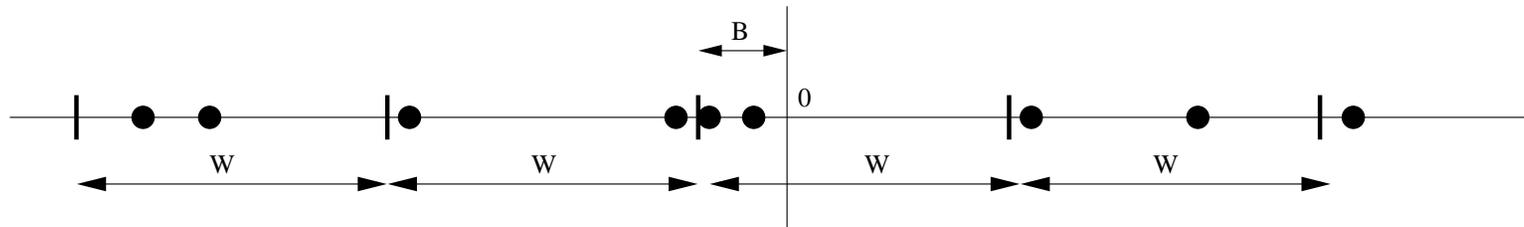
- Consider a vector  $\mathbf{X} = \{X_1, X_2, \dots, X_N\}$ , where each  $X_i$  is drawn from a  $p$ -Stable distr
- For any pair of vectors  $\mathbf{a}, \mathbf{b}$   $\mathbf{a} \cdot \mathbf{X} - \mathbf{b} \cdot \mathbf{X} = (\mathbf{a} - \mathbf{b}) \cdot \mathbf{X}$  (by linearity)
- Thus  $\mathbf{a} \cdot \mathbf{X} - \mathbf{b} \cdot \mathbf{X}$  is distributed as  $(l_p(\mathbf{a} - \mathbf{b}))X'$  where  $X'$  is a  $p$ -Stable distr r.v.
- Using multiple independent  $\mathbf{X}$ 's we can use  $\mathbf{a} \cdot \mathbf{X} - \mathbf{b} \cdot \mathbf{X}$  to estimate  $l_p(\mathbf{a} - \mathbf{b})$  [Ind'01]

## How are $p$ -Stable distributions useful?



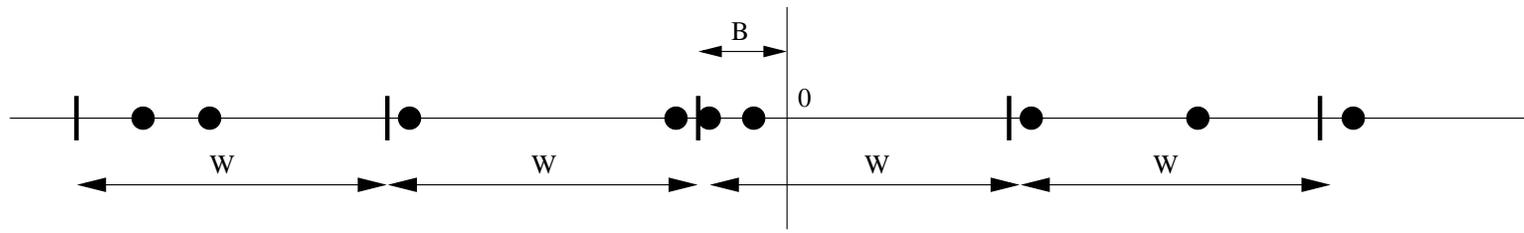
- For a vector  $\mathbf{a}$ , the dot product  $\mathbf{a} \cdot \mathbf{X}$  projects it onto the real line
- For any pair of vectors  $\mathbf{a}, \mathbf{b}$  these projections are “close” (w.h.p.) if  $l_p(\mathbf{a} - \mathbf{b})$  is “small” and “far” otherwise
- Divide the real line into segments of width  $w$
- Each segment defines a hash bucket, i.e. vectors that project onto the same segment belong to the same bucket

## Hashing (formal) definition



- Consider  $h_{\mathbf{a},b} \in \mathcal{H}^w$ ,  $h_{\mathbf{a},b}(\mathbf{v}) : \mathcal{R}^d \rightarrow \mathcal{N}$
- $\mathbf{a}$  is a  $d$  dimensional random vector whose each entry is drawn from a  $p$ -stable distr
- $b$  is a random real number chosen uniformly from  $[0, w]$  (random shift)
- $h_{\mathbf{a},b}(\mathbf{v}) = \lfloor \frac{\mathbf{a} \cdot \mathbf{v} + b}{w} \rfloor$

## Collision probabilities



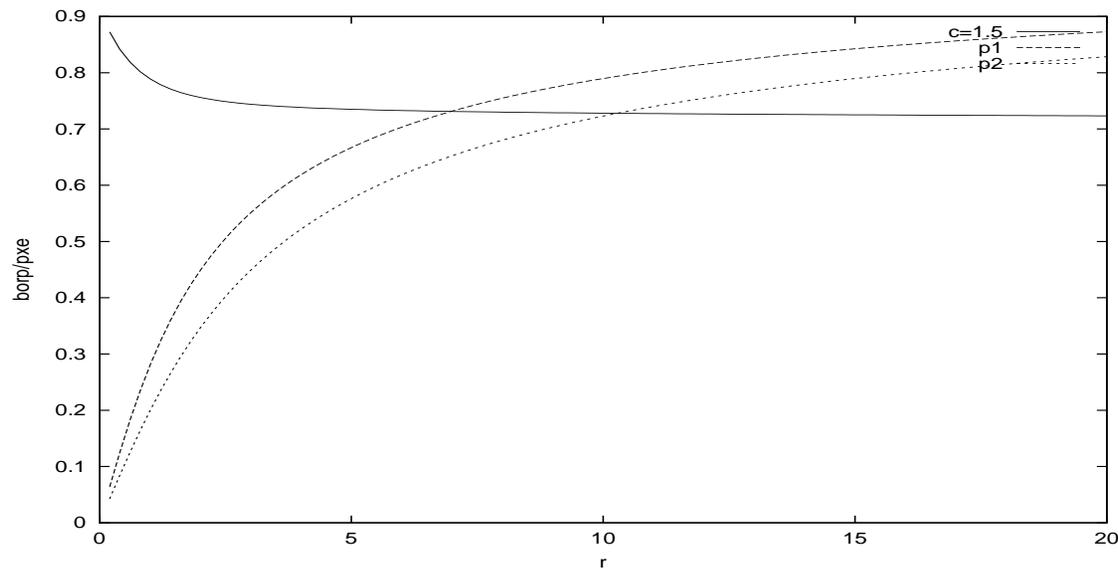
- Consider two vectors  $v_1, v_2$  and let  $\ell = l_p(v_1, v_2)$
- Let  $Y$  denote the distance between their projections onto the random vector  $a$  ( $Y$  is distributed as  $\ell X$  where  $X$  is a  $p$ -stable distr r.v.)
- if  $Y > w$ ,  $v_1, v_2$  will not collide
- if  $Y \leq w$ ,  $v_1, v_2$  will collide with probability equal to  $(1 - (Y/w))$  (random shift  $b$ )

## Collision probabilities

- $f_p(t)$ : p.d.f. of the absolute value of a  $p$ -stable distribution
- $\ell = l_p(\mathbf{v}_1, \mathbf{v}_2)$
- $\ell \leq 1$ ,  $p_1 = Pr[h_{\mathbf{a},b}(\mathbf{v}_1) = h_{\mathbf{a},b}(\mathbf{v}_2)] \geq \int_0^w f_p(t)(1 - \frac{t}{w})dt$
- $\ell > 1 + \epsilon = c$ ,  $p_2 = Pr[h_{\mathbf{a},b}(\mathbf{v}_1) = h_{\mathbf{a},b}(\mathbf{v}_2)] \leq \int_0^w \frac{1}{c} f_p(\frac{t}{c})(1 - \frac{t}{w})dt$
- $\mathcal{H}^w$  hash family is  $(r_1, r_2, p_1, p_2)$ -sensitive for  $r_1 = 1$ ,  $r_2 = c$  and  $p_1, p_2$  given as above

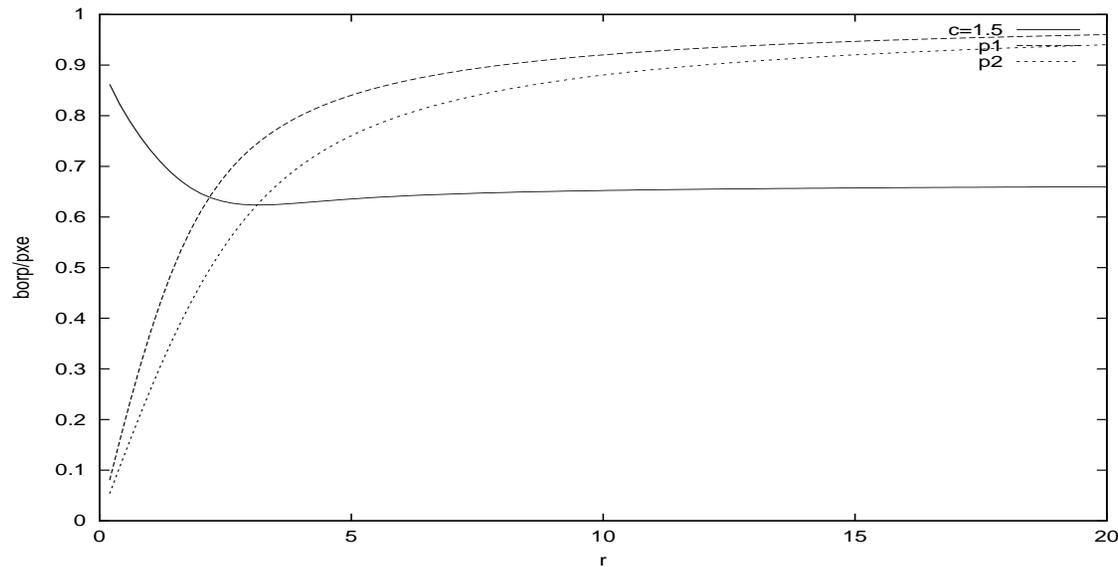
## Special cases

- $p = 1$  (Cauchy distr):  $f_p(t) = \frac{2}{\pi} \frac{1}{1+t^2}$
- $p_2 = 2 \frac{\tan^{-1}(w/c)}{\pi} - \frac{1}{\pi(w/c)} \ln(1 + (w/c)^2)$
- $p_1$  obtained by substituting  $c = 1$



## Special cases

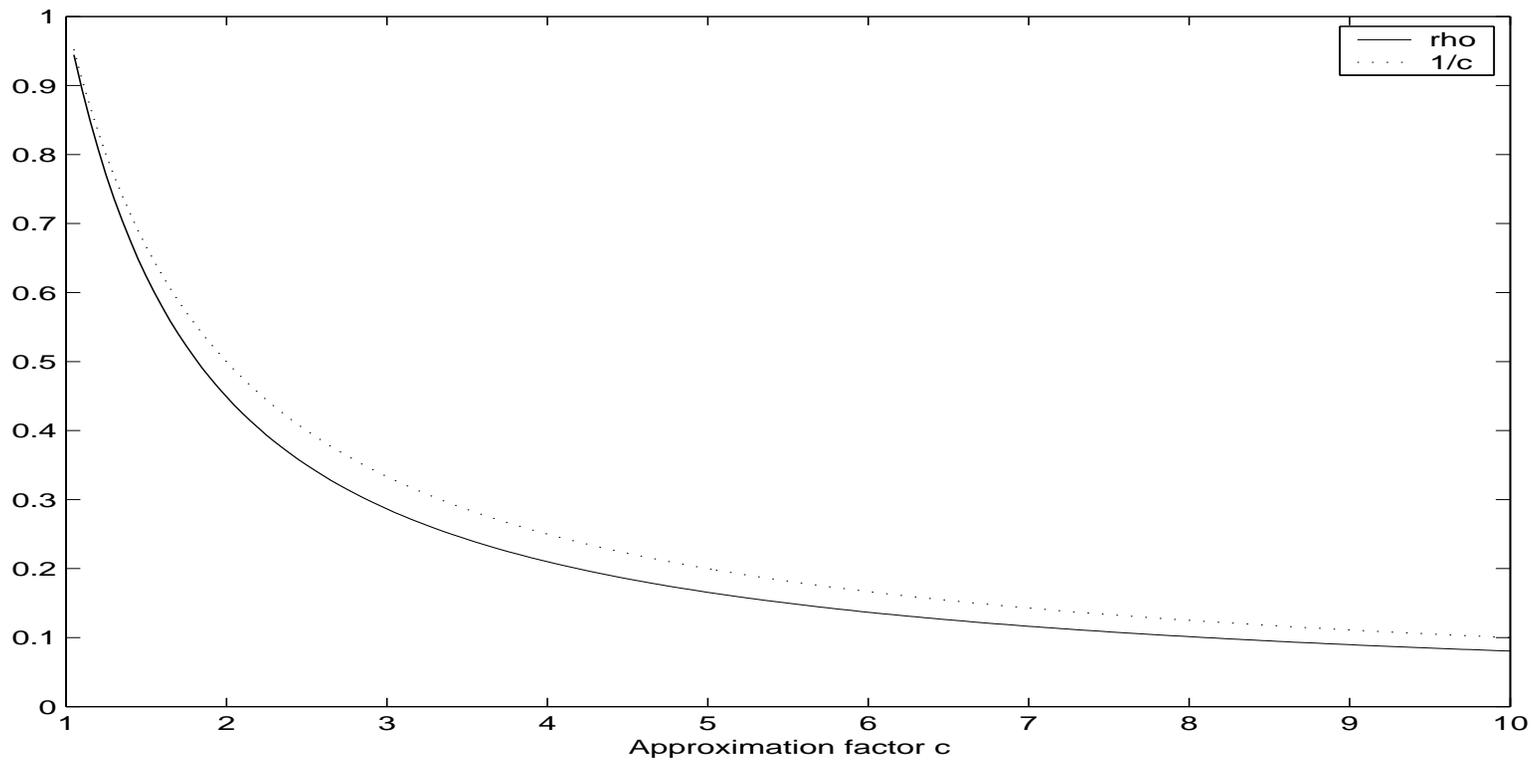
- $p = 2$  (Gaussian distr):  $f_p(t) = \frac{2}{\sqrt{2\pi}} e^{-t^2/2}$
- $p_2 = 1 - 2\text{norm}(-w/c) - \frac{2}{\sqrt{2\pi}w/c} (1 - e^{-(w^2/2c^2)})$
- $p_1$  obtained by substituting  $c = 1$



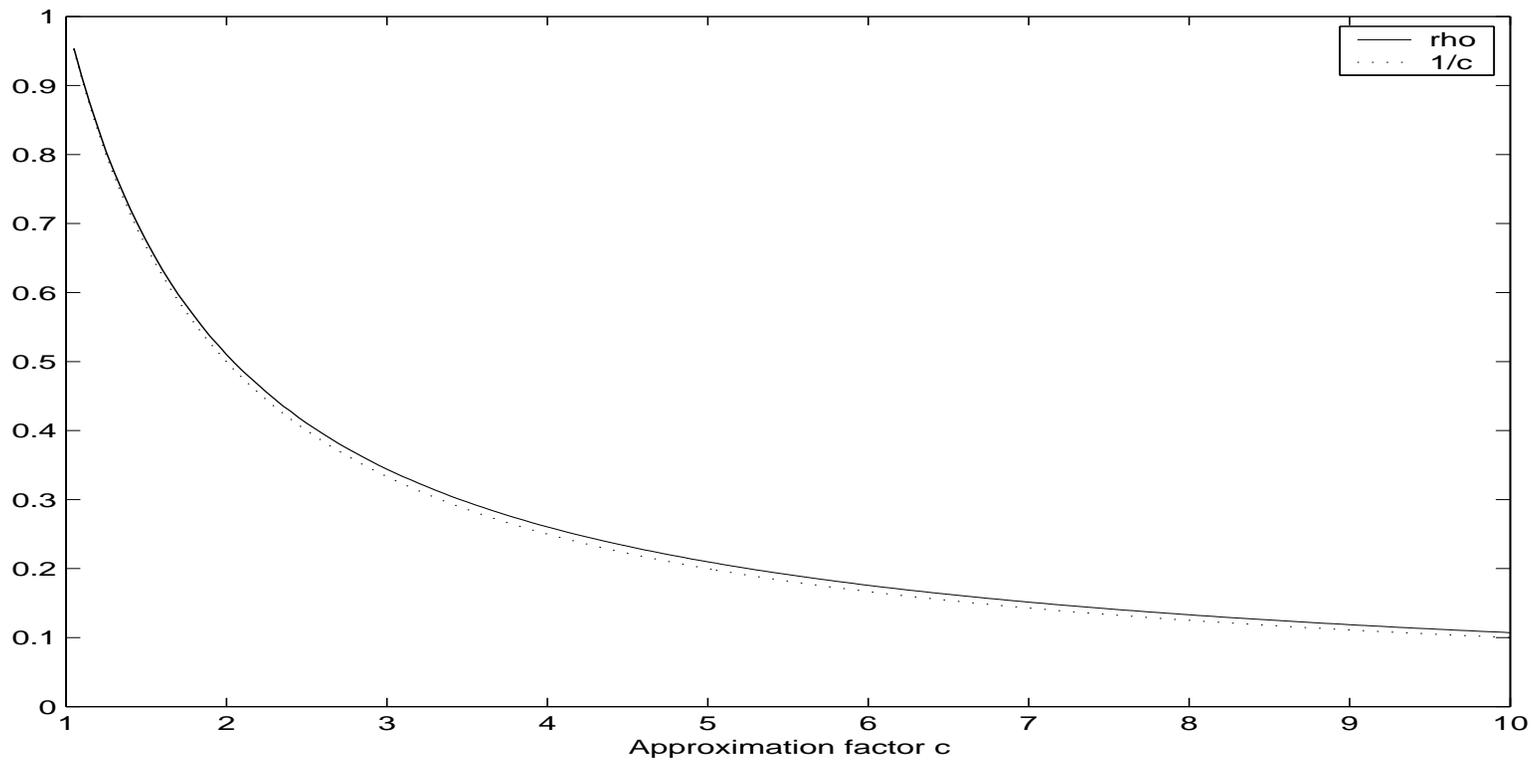
## Comparison with previous scheme

- Previous hashing scheme for  $p = 1, 2$  achieved  $\rho = 1/c$
- Based on reduction to hamming distance
- New scheme achieves smaller  $\rho$  (than  $1/c$ ) for  $p = 2$
- Large constants and log factors for  $p = 2$  in query time besides  $n^\rho$
- Achieves  $\rho = 1/c$  for  $p = 1$

## $\rho$ for $p = 2$



## $\rho$ for $p = 1$



## General case

- what about general case, i.e.  $p \neq 1, 2$ ?

**Theorem.** For any  $p \in (0, 2]$  there is a  $(r_1, r_2, p_1, p_2)$ -sensitive family  $\mathcal{H}^w$  for  $l_p^d$  such that for any  $\gamma > 0$ ,

$$\rho = \frac{\ln 1/p_1}{\ln 1/p_2} \leq (1 + \gamma) \cdot \max\left(\frac{1}{c^p}, \frac{1}{c}\right).$$

- Achieves  $\frac{1}{c^p}$  for  $p < 1$

## Conclusions

- New LSH scheme for  $0 < p \leq 2$ . First one for  $0 < p < 1$
- Easy to implement (experiments in progress)
- Easy to update hash value in cash register model
- Improves running time for  $p = 2$  over previous scheme