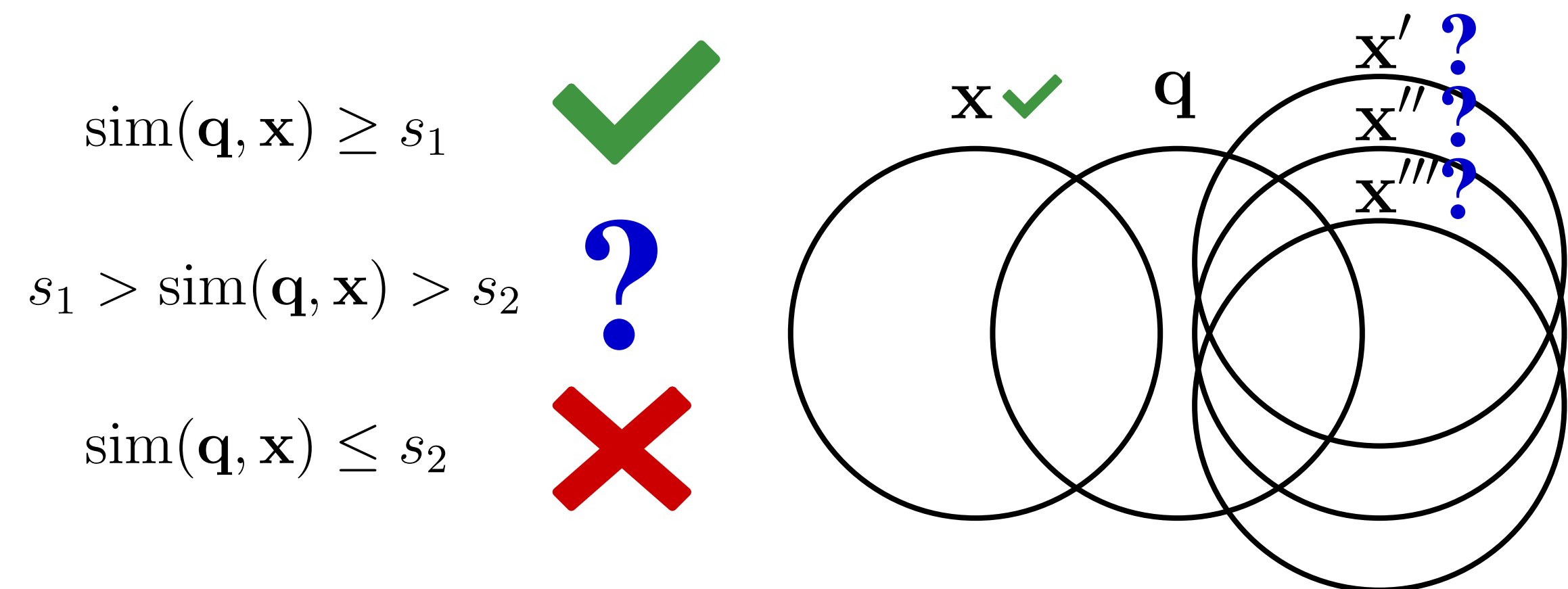


## Problem

- Set similarity search:
  - Preprocess a collection  $P$  of  $n$  subsets of  $[d] = \{1, \dots, d\}$ .
  - Given a query set  $q \subseteq [d]$  return  $x \in P$  such that  $\text{sim}(q, x) \geq s_1$ .
- Simplification: All sets have size  $t$ .
  - $\text{sim}(x, y) = |x \cap y|/t$ .
- Approximation: Allow returning  $x \in P$  with  $s_1 > \text{sim}(q, x) > s_2$ .



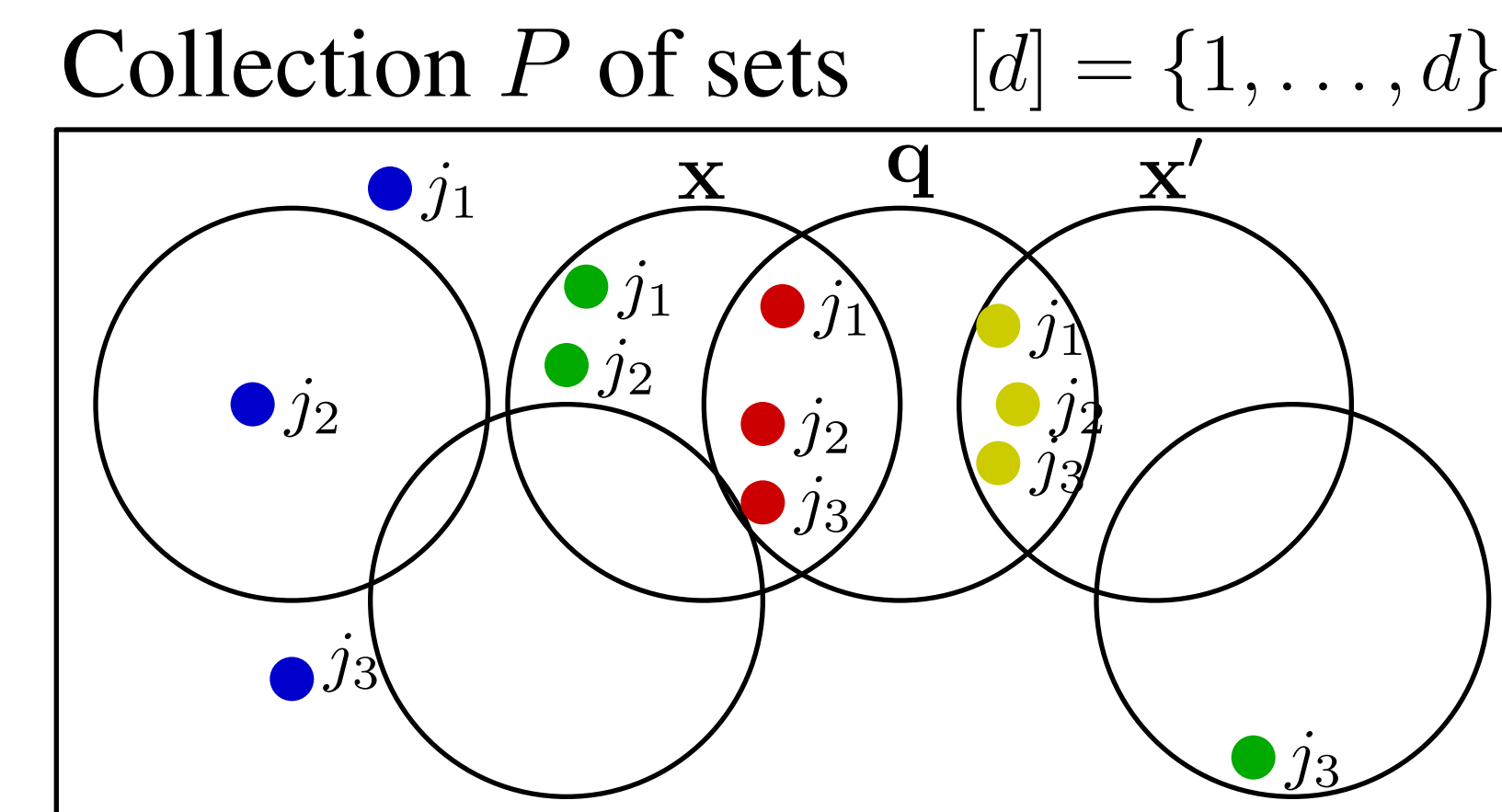
## Contribution

Data structure solution	
Space: $n^{1+\rho} + tn$	Query time: $tn^\rho$
Upper bound $\rho = \frac{\log(1/s_1)}{\log(1/s_2)}$	Lower bound* $\rho \geq \frac{\log(1/s_1)}{\log(1/s_2)} - o_d(1)$

- Example:  $s_1 = 0.5, s_2 = 0.25$  gives  $\rho = 0.5$ 
  - MinHash  $\rho = 0.56$
- Solution for Braun-Blanquet similarity:

$$\text{sim}(x, y) = \frac{|x \cap y|}{\max(|x|, |y|)}$$

## Random buckets



- Collection  $B$  of random buckets  $(j_1, \dots, j_k) \in [d]^k$ .
- $B(x) = \{(j_1, \dots, j_k) \in B \mid j_1 \in x \wedge \dots \wedge j_k \in x\}$

$$\text{Query time: } \underbrace{|B(q)|}_{\text{Lookups}} + \underbrace{\sum_{x \in P: \text{sim}(q, x) \leq s_2} |B(q) \cap B(x)|}_{\text{Collisions}} + \underbrace{|B|}_{\text{Buckets}}$$



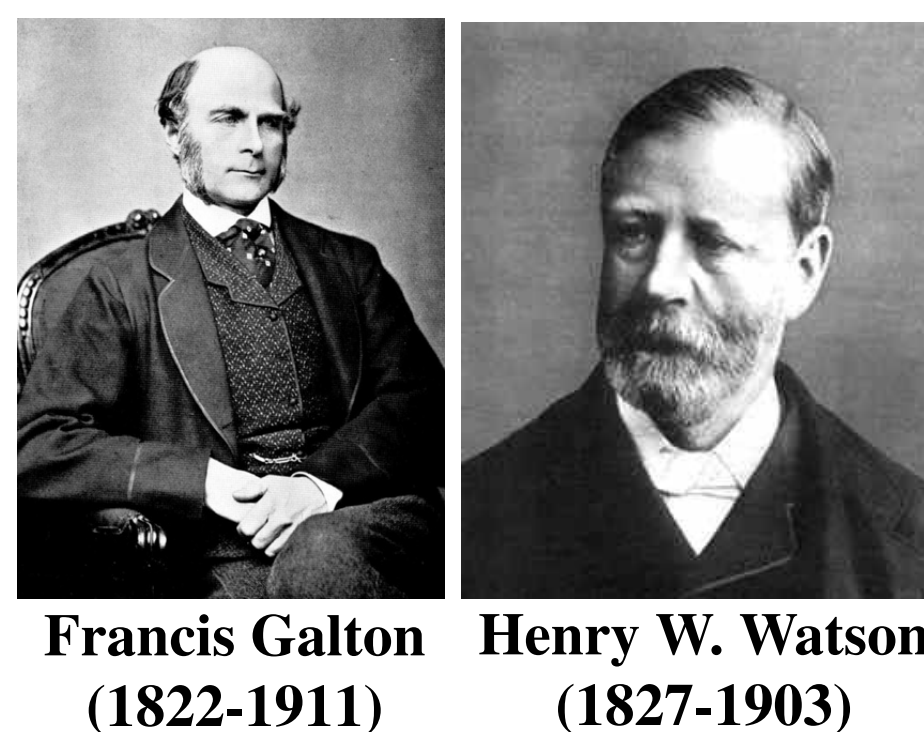
$|B| \gg n \gg |B(q)|$

## Analysis

- Correctness:  $q, x$  with  $\text{sim}(q, x) \geq s_1$ .
  - Pr. collision  $(|q \cap x|/d)^k \geq (s_1 t/d)^k$ .
  - Use  $b = \lceil (s_1 t/d)^{-k} \rceil$  buckets.
- Lookups:
  - $E[|B(q)|] = b(t/d)^k = (1/s_1)^k$ .
- Comparisons:  $q, x'$  with  $\text{sim}(q, x') \leq s_2$ .
  - $E[|B(q) \cap M(x')|] = b(|q \cap x'|/d)^k \leq b(s_2 t/d)^k = (s_2/s_1)^k$ .
- Balance:  $k = \log(n)/\log(1/s_2)$ 
  - Space:  $n^{1+\rho} + tn$ .
  - Query time:  $tn^\rho$

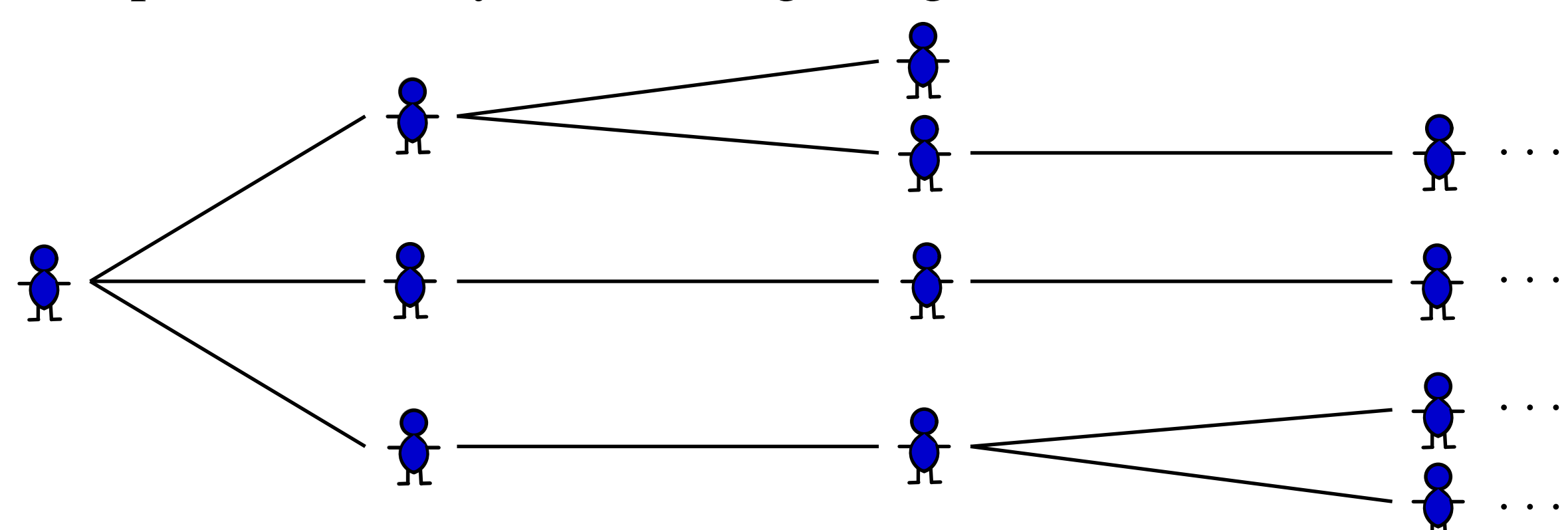
## Branching processes

- Study of population growth in the 1870's.
  - Will my surname survive  $k$  generations?
- Offspring  $X$  i.i.d. across generations.
  - Example:  $X \sim \text{bin}(3, 1/3)$



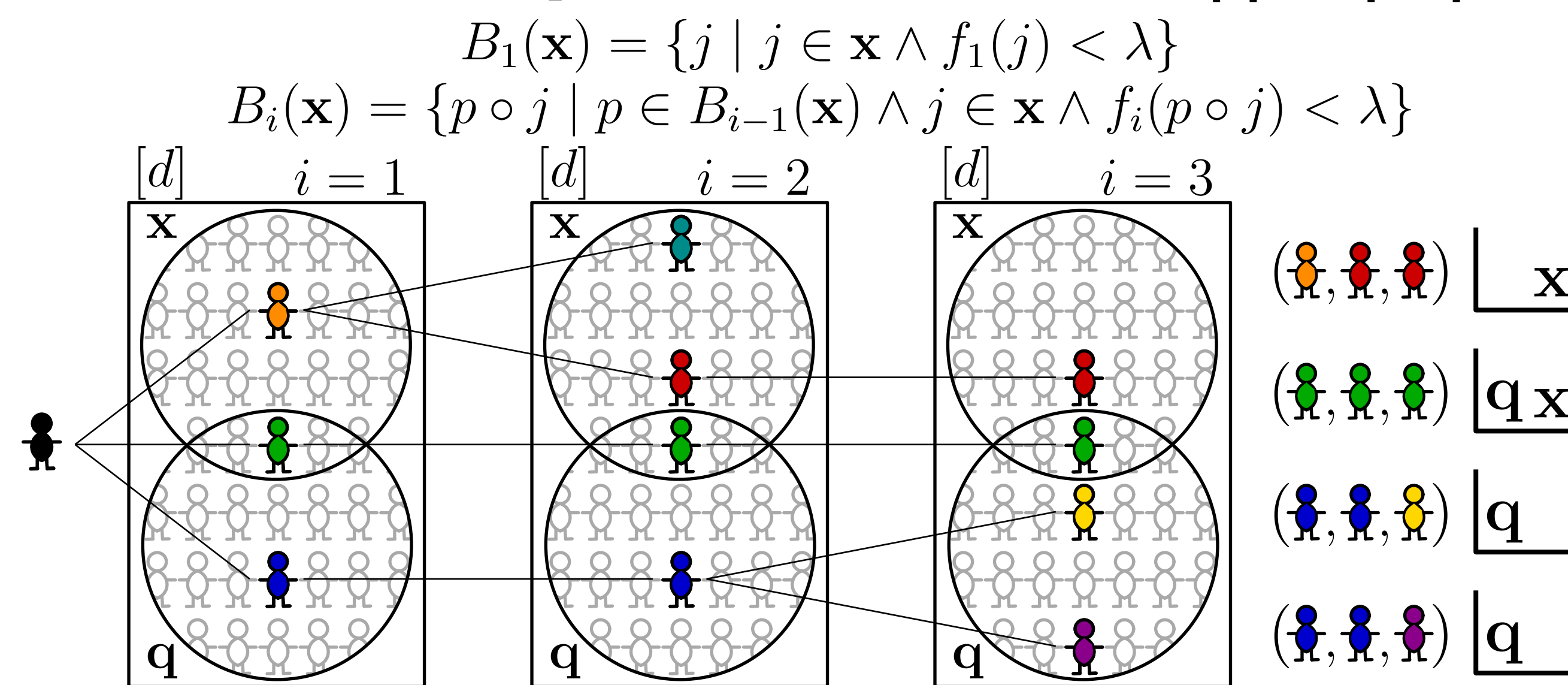
**Lemma** [Agresti 1974] Let  $\mathbb{E}[X] = 1$ , then the probability of surviving  $k$  generations is at least  $1/(k\text{Var}[X] + 1)$ .

- Example: Probability of surviving 100 generations is at least 0.01.



## Chosen Path

- Buckets  $B(x)$  as paths  $p = (j_1, \dots, j_k)$  in a branching process.
- For  $i = 1, \dots, k$  sample random hash functions  $f_i: [d]^i \rightarrow [0, 1]$ .



- Correctness: Offspring  $X = |B_1(q) \cap B_1(x)| \sim \text{bin}(|x \cap q|, \lambda)$ .
  - Set  $\lambda = 1/s_1 t$ : For  $\text{sim}(q, x) \geq s_1$  we have  $\mathbb{E}[X] \geq 1$ .

Compute  $B(x)$  in time  $O(|B(x)|)!!!$

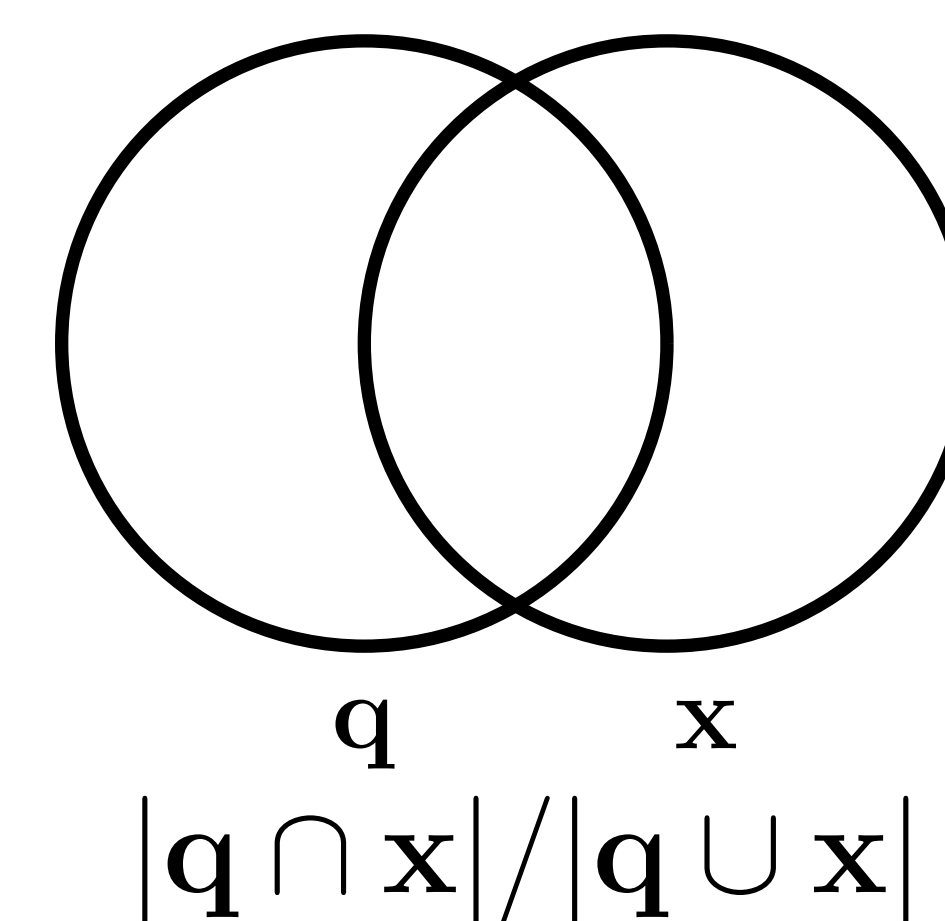
## Beyond MinHash?

- Probability  $p(\text{sim}(q, x))$  of sets colliding in bucket.
  - Performance  $\rho = \log(1/p(s_1))/\log(1/p(s_2))$ .

### MinHash [B '97]

$$h(q) = \arg \min_{j \in q} f(j)$$

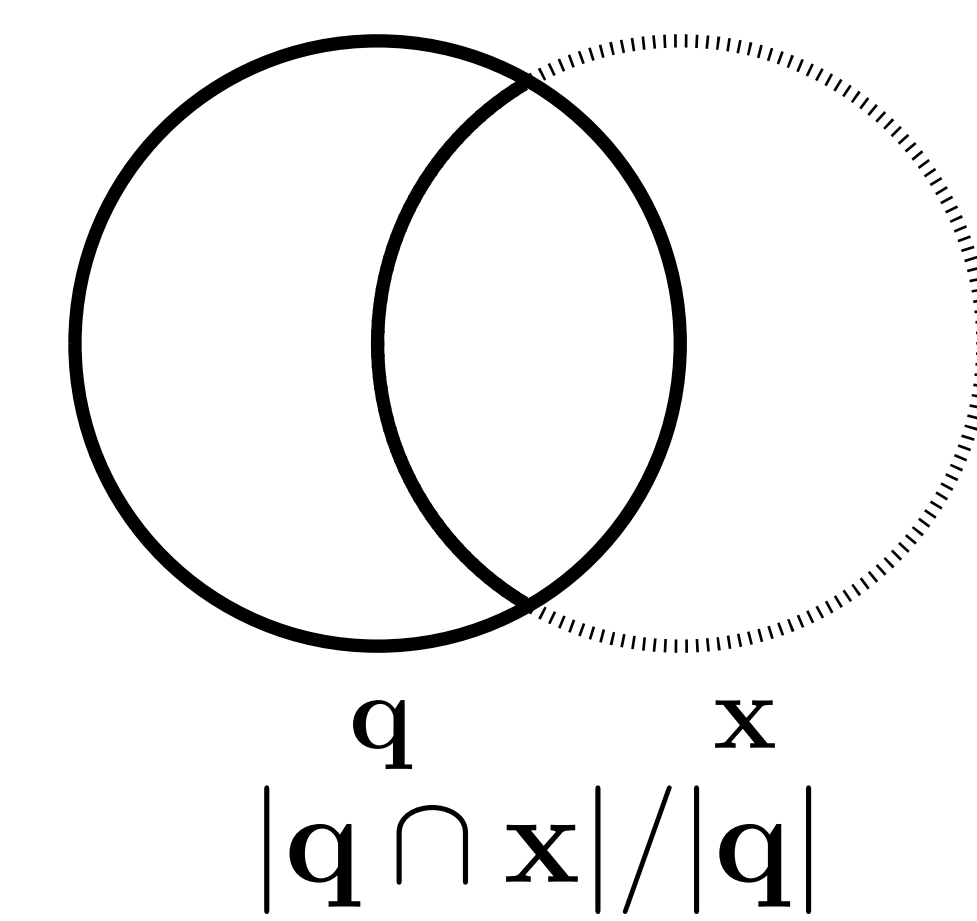
Random  $f: [d] \rightarrow [0, 1]$ .



### Chosen Path

$$B(q) = \{j \in q \mid f(j) < \lambda\}$$

Random  $f: [d] \rightarrow [0, 1]$ .



## Previous work

Reference	$\rho$
This paper	$\log(1/s_1)/\log(1/s_2)$
MinHash [B'97 + IM'98]	$\log \frac{s_1}{2-s_1} / \log \frac{s_2}{2-s_2}$
Angular LSH [AILRS'15]	$\frac{1-s_1}{1+s_1} / \frac{1-s_2}{1+s_2}$
Data-dep. LSH [AR'15]	$\frac{1-s_1}{1+s_1-2s_2}$

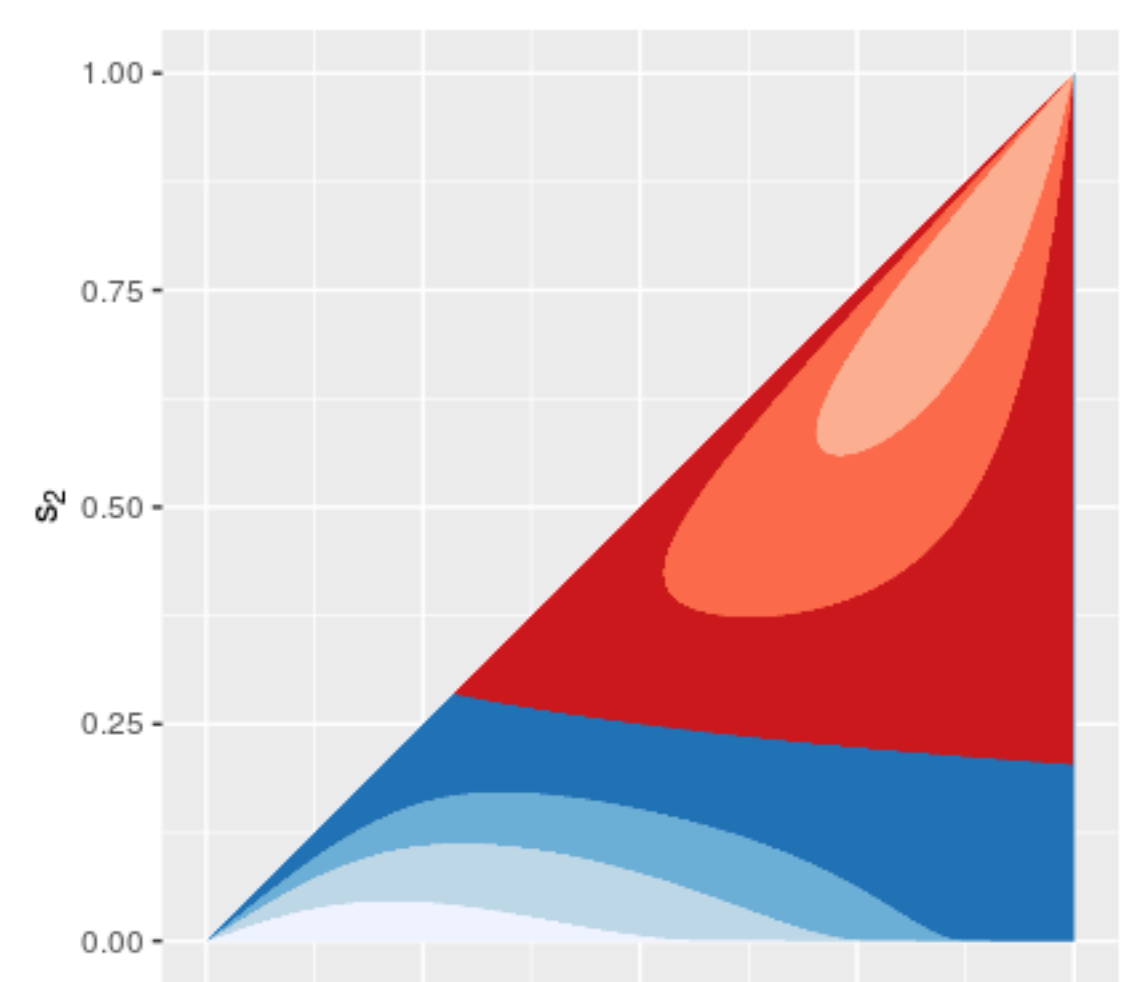


Fig:  $\rho_{\text{new}} - \rho_{\text{data}}$  for  $0 < s_2 < s_1 < 1$

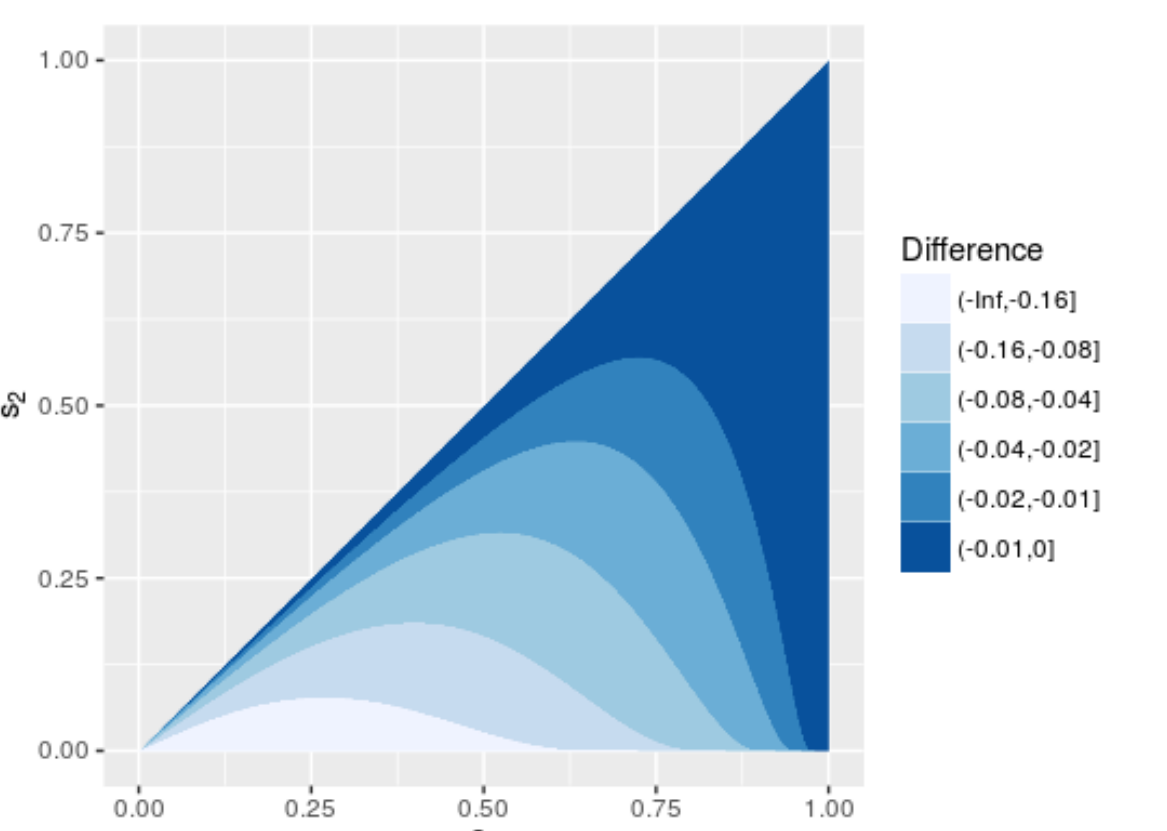


Fig:  $\rho_{\text{new}} - \rho_{\text{angular}}$  for  $0 < s_2 < s_1 < 1$

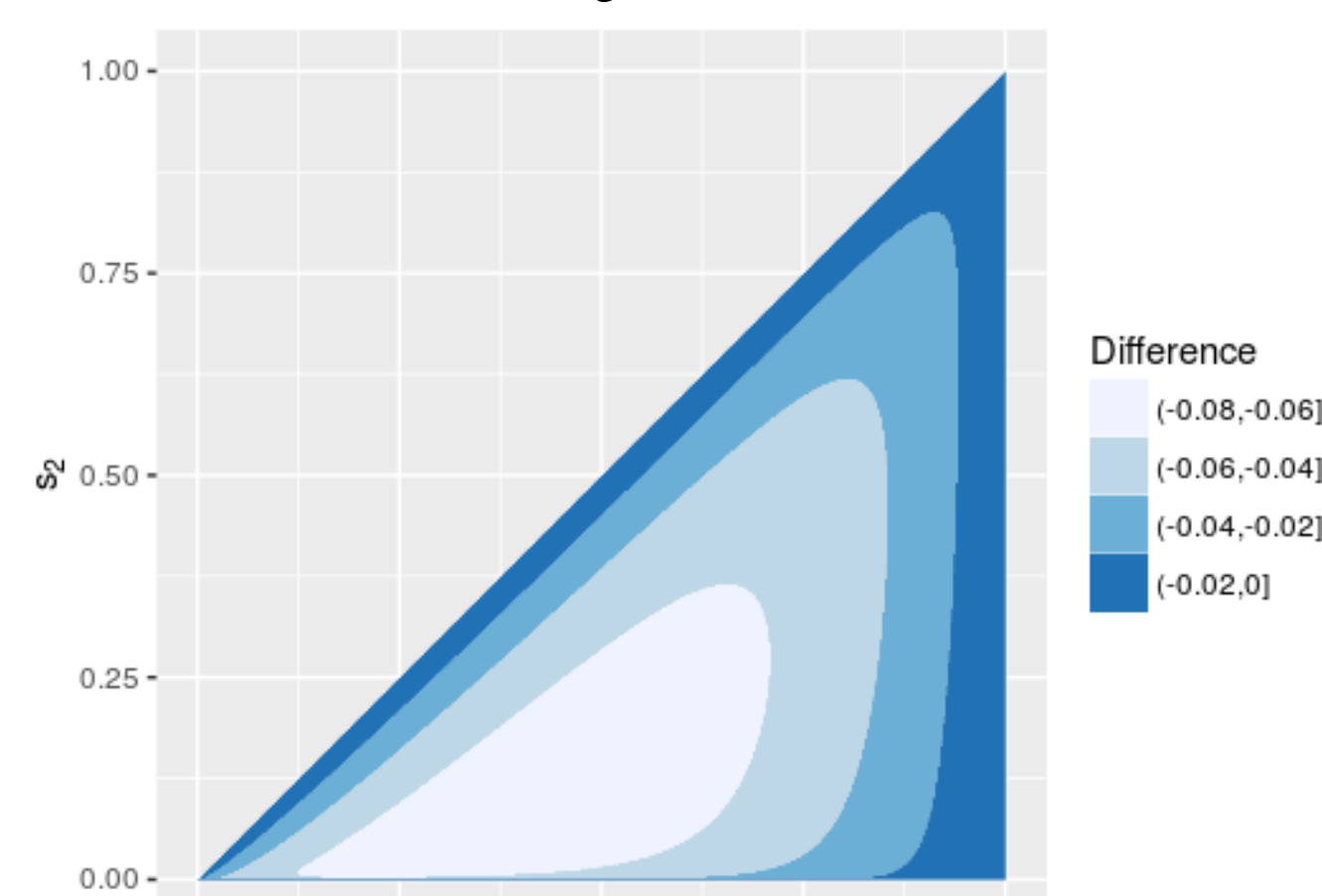


Fig:  $\rho_{\text{new}} - \rho_{\text{minhash}}$  for  $0 < s_2 < s_1 < 1$

## Lower bound

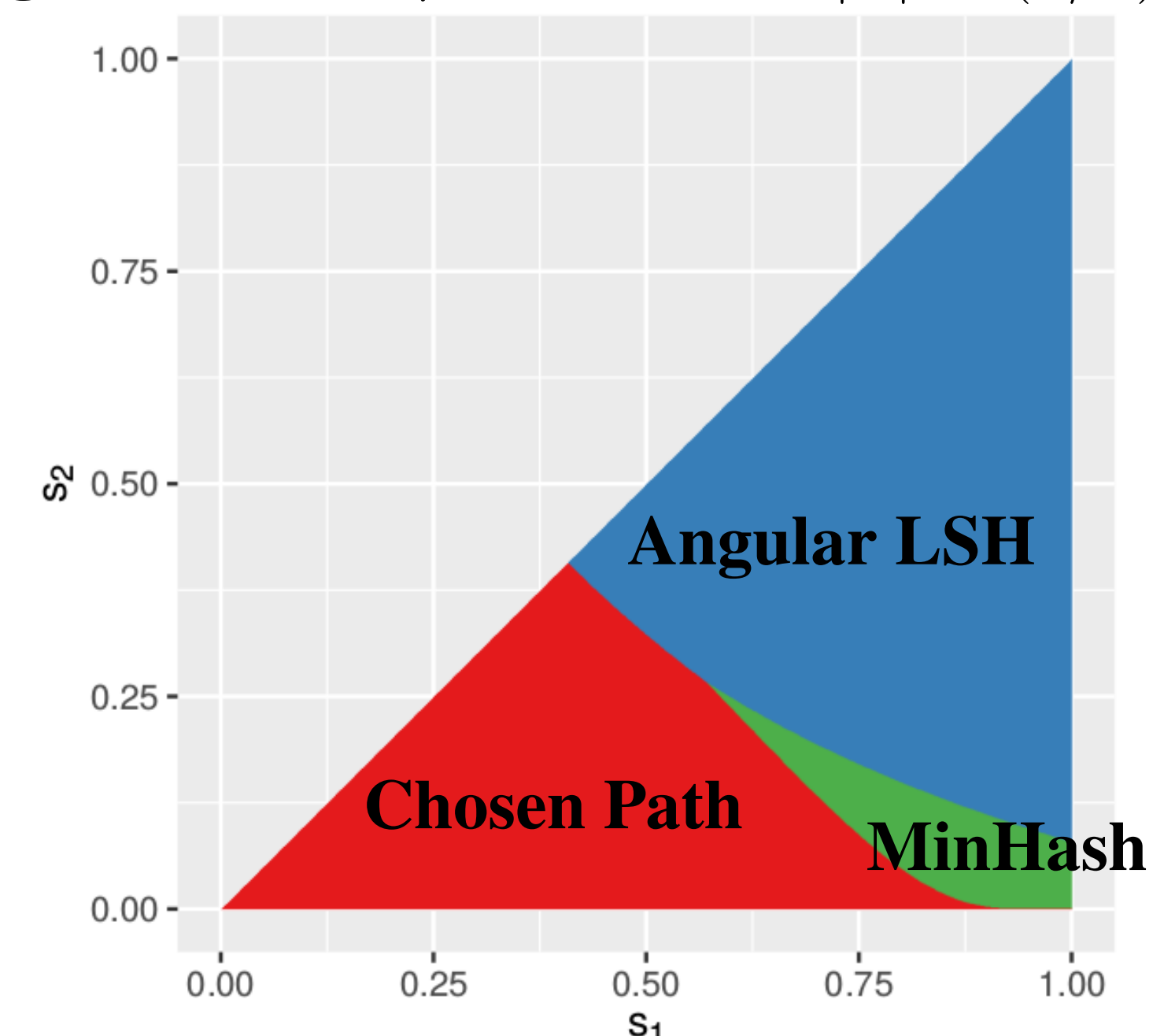
- $(r, cr)$ -near neighbor problem in Hamming space:
  - Preprocess a collection  $P$  of  $n$  points in  $\{0, 1\}^D$ .
  - Given a query point  $q \in \{0, 1\}^D$ , if there exists  $x \in P$  with  $\|q - x\|_1 \leq r$  return  $x' \in P$  with  $\|q - x'\|_1 < cr$ .
- [OWZ'11] lower bound:  $\rho \geq 1/c$  for space  $n^{1+\rho}$ , time  $n^\rho$  solution.
- Suppose a solution for the  $t$ -regular  $(s_1, s_2)$ -similarity problem with:
  - $\rho < \log(1/s_1)/\log(1/s_2)$
- $(r, cr)$ -nn. in Hamming space as set similarity search:
  - $\tilde{q}, \tilde{x} \in \{0, 1\}^D$  with  $\|\tilde{q} - \tilde{x}\|_1 = \Delta$ .
  - Define  $x = \{2j - \tilde{x}_j \mid j \in [D]\}$ .
  - Tensor:  $x^{\otimes \tau} = \{(j_1, \dots, j_\tau) \mid j_i \in x\}$ .
  - $|q^{\oplus \tau}| = |x^{\oplus \tau}| = D^\tau$  and  $\text{sim}(q^{\oplus \tau}, x^{\oplus \tau}) = (1 - \Delta/D)^\tau$ .
- Choose  $D, r, cr, \tau$  s.t.  $(1 - r/D)^\tau \approx s_1$  and  $(1 - r/D)^\tau \approx s_2$ .

$$\rho < \frac{\log(1/s_1)}{\log(1/s_2)} \approx \frac{\log(1/(1 - r/D))}{\log(1/(1 - cr/D))} \leq 1/c$$

## Sets of different sizes

- General set similarity search:
  - Query sets have size  $|q| = t$ .
  - Data sets have size  $|x| = t'$ .
  - $\text{sim}(q, x) = |q \cap x|/\min(|q|, |x|)$ .

Figure: Lowest  $\rho$ -value when  $|q| = (1/2)|x|$ .



- Data-dependent:
  - $(s_1, s_2, P) \rightarrow \mathcal{M}$ .
- Data-independent:
  - $(s_1, s_2, |P|) \rightarrow \mathcal{M}$ .