

K, L are input parameters

→ Input: → $D = \{x_i\}_{i=1}^N$

→ Preprocess D s.t.

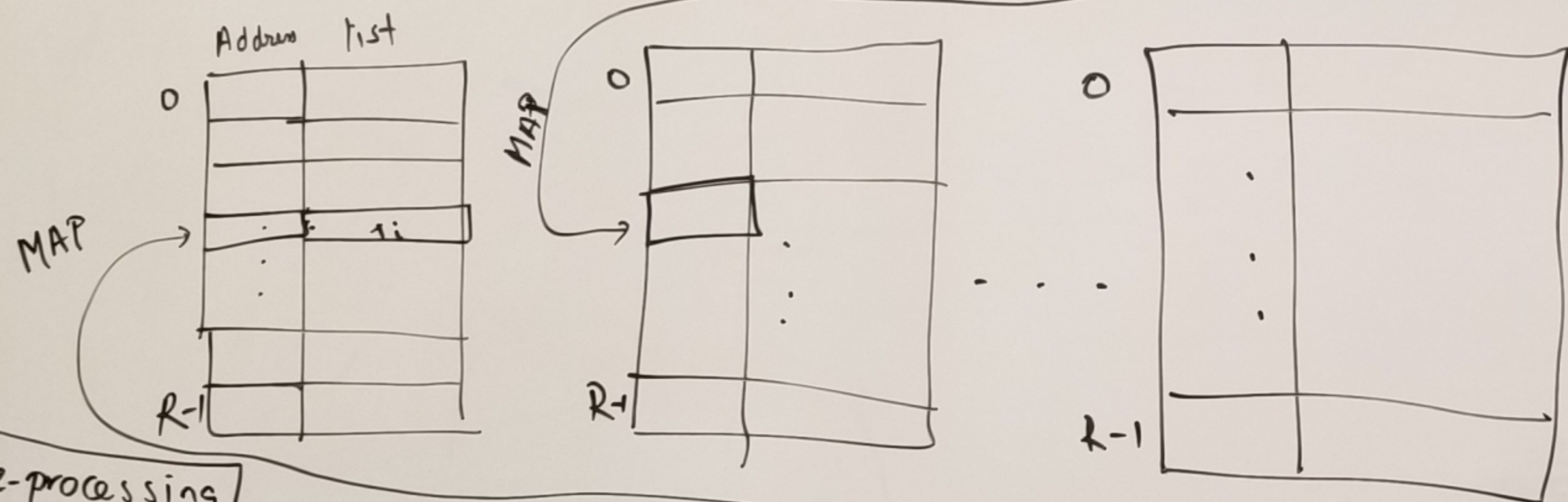
- Given a query q find $x_i \in D$
s.t. $\text{sim}(q, x_i)$ is high (very high)

- $\ll O(N)$ time.

→

- Initialize $K \times L$ independent hash functions. \rightarrow LSM [minwise hashing]

- Create L hash tables. $[0 - (R-1)]$



MAP(I_1, I_2, \dots, I_k)

return

$$(a_1 I_1 + a_2 I_2 + \dots + a_k I_k + a_{k+1}) \bmod P \bmod R$$

a_i 's are chosen randomly and fixed

Pre-processing

- for each $x_i \in D$.

Compute $K \times L$ Hash of x_i
for $j = 1$ to L

MAP $[h_{j1}, h_{j2}, \dots, h_{jk}]$ in hash table j
 \downarrow
 $[0 - (R-1)]$

HASH TABLE 1: $[h_{11}(x_i), h_{12}(x_i), \dots, h_{1k}(x_i)]$
 HASH TABLE 2: $[h_{21}(x_i), h_{22}(x_i), \dots, h_{2k}(x_i)]$... $h_{Lk}(x_i)$

Query: \rightarrow

Given $\rightarrow q$, $K \times L$ hash functions, L Hash tables from pre-processing.

- Compute $K \times L$ ^{LSM} hashes of query q .

Returned Set = ϕ

for $j = 1$ to L

Returned Set = Returned Set $\overset{\text{Union}}{=} H_j [\text{MAP} (h_{j1}(q), h_{j2}(q), \dots, h_{jk}(q))]$

Hash table j
 \downarrow

\rightarrow return Returned Set. (Candidates to search)

\rightarrow return top- k $x_b \in$ Returned Set with highest sim(q, x_b)

EXIT

→ Given q , the probability of $x_j \in$ Returned Set

$$\Rightarrow 1 - (1 - \text{Collprob}(q, x_j))^L$$

$$\text{Collprob}(q, x_j) = \Pr(h(x_j) = h(q)) = \text{Jaccard Similarity if } h \text{ is minwise hashing.}$$

→ Monotonic in $\text{Collprob}(q, x_j)$

$\Rightarrow x_j$ is returned with higher probability than x_i iff.

$$\text{Sim}(x_j, q) > \text{Sim}(x_i, q)$$

→ Given q , the probability of $x_j \in$ Returned Set

$$\Rightarrow 1 - \left(1 - \underbrace{\text{Collprob}(q, x_j)}^k\right)^L \left[\text{Collprob}(q, x_j)^k + (1 - \text{Collprob}(q, x_j)^k) \times \frac{1}{R} \right]$$

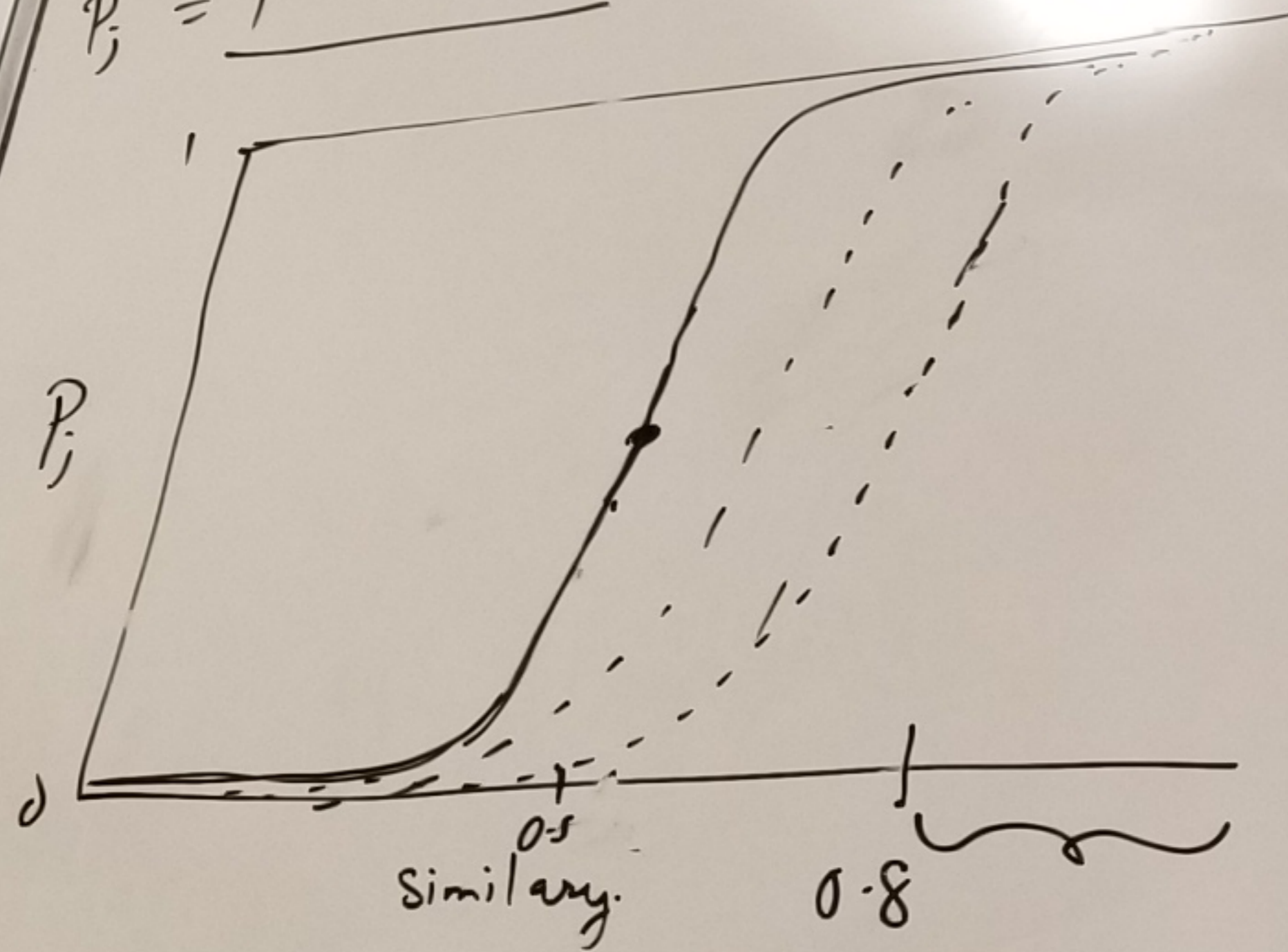
$\text{Collprob}(q, x_j) = \Pr(h(x_j) = h(q)) = \text{Jaccard Similarity}$ if h is minwise hashing.

→ Monotonic in $\text{Collprob}(q, x_j)$

$\Rightarrow x_j$ is returned with higher probability than x_i iff.

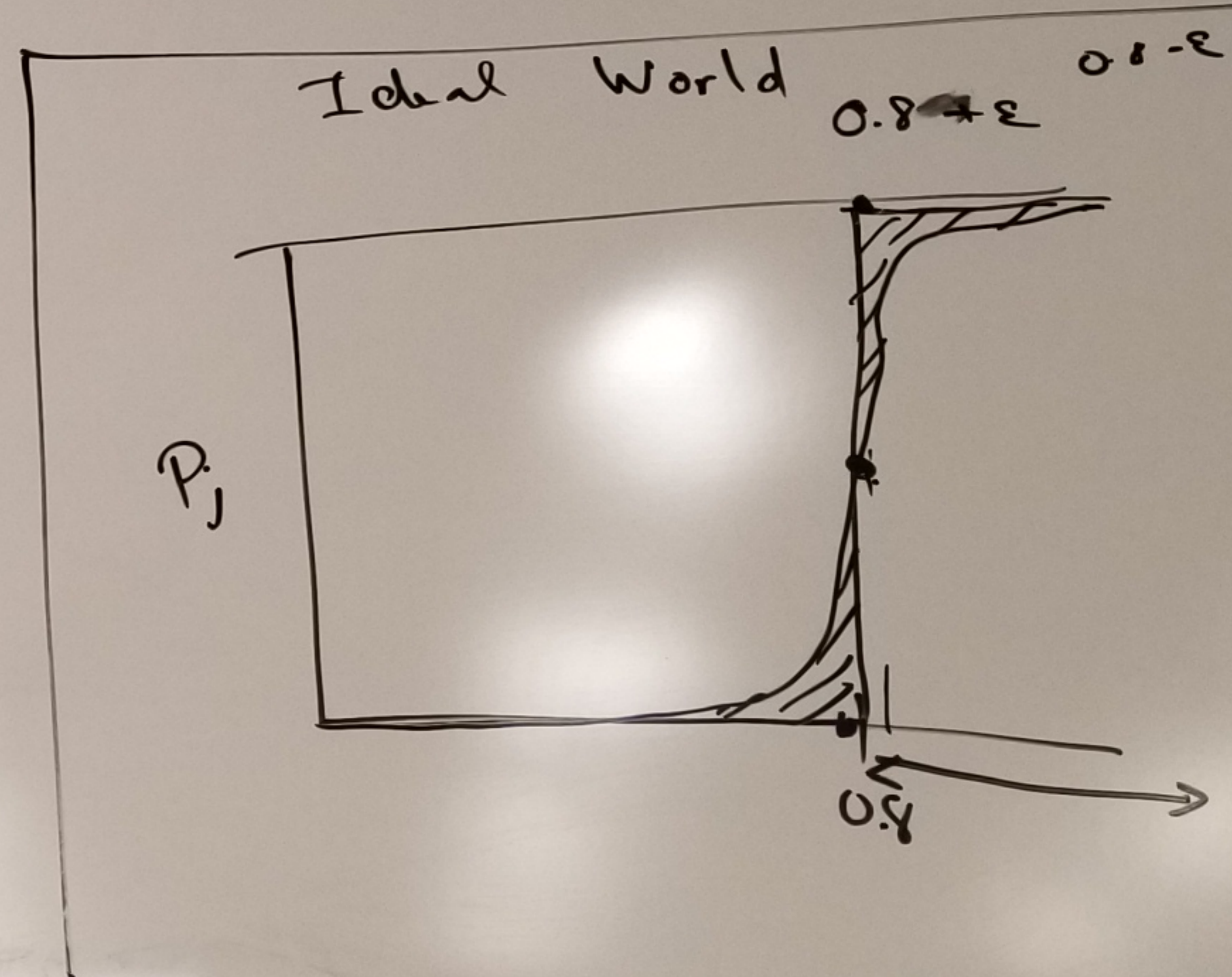
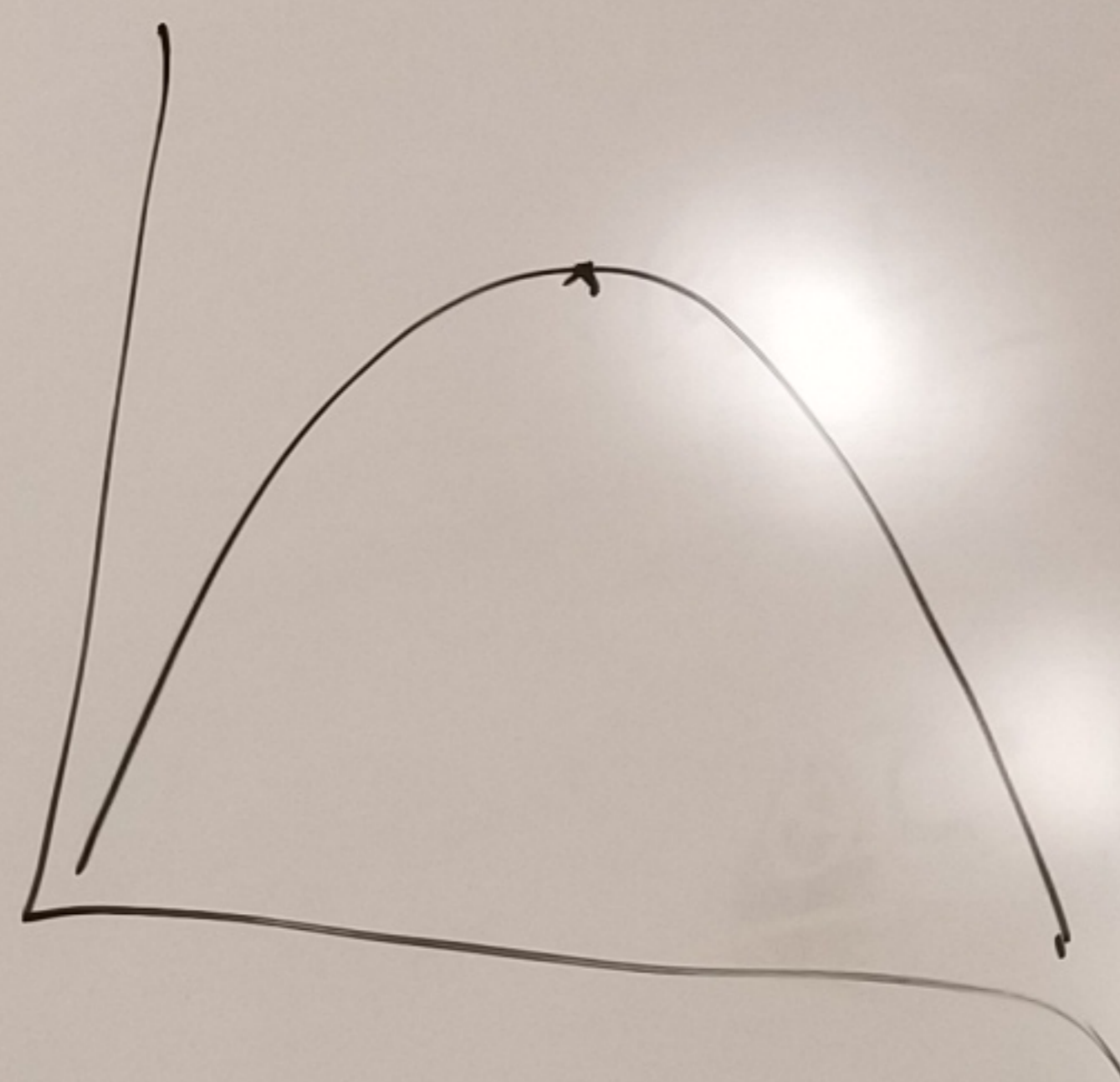
$$\text{Sim}(x_j, q) > \text{Sim}(x_i, q)$$

$$P_i = 1 - (1 - c_j^k)^L$$



c_j	P_i
.2	.006
.3	0.047
0.5	0.18
0.6	0.8
0.7	.975
.8	.9996

$L = 20$
 $K = 5$



- Initialize $K \times L$

→ SM P_i

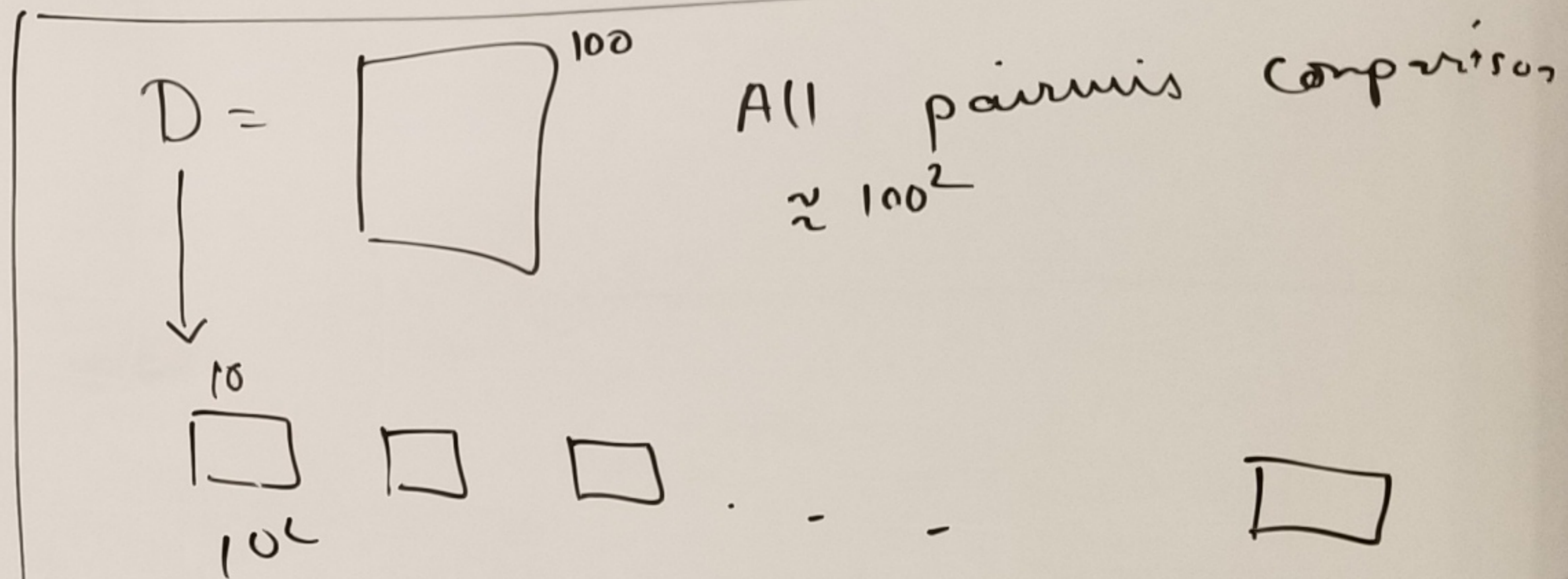
- Reducing all pairwise comparisons.
New-Duplicate Detection, Entity Resolution
- Estimate Number of Unique Entries in a Given Record.

$$D_1 = \{x_i\}_{i=1}^N ; D_2 = \{y_i\}_{i=1}^N \dots$$

↓
text

- 350,000 (190,000)

~ 63 billion



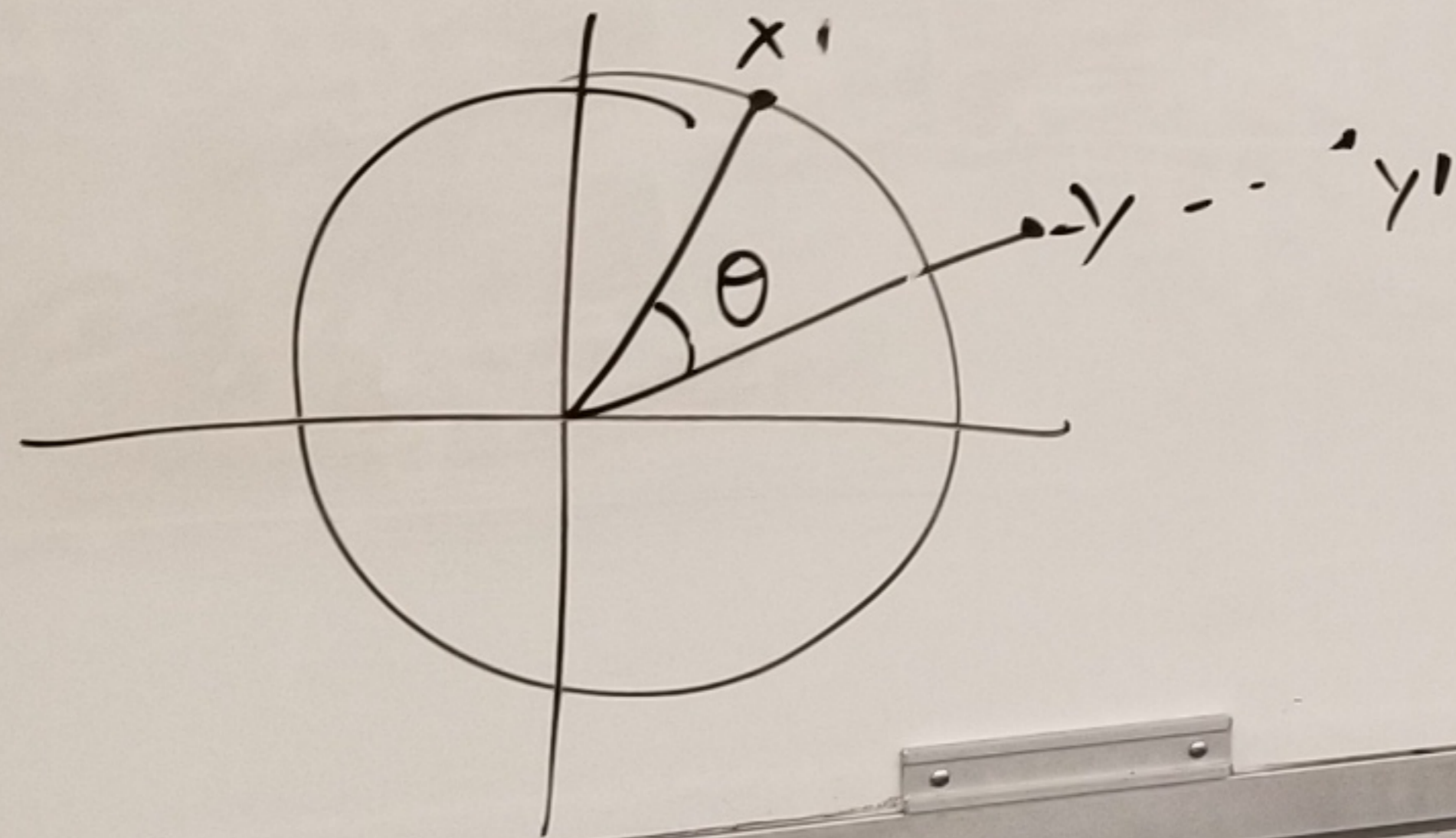
$$\rightarrow \Pr(h(x) = h(y)) = f(\text{sim}(x, y))$$

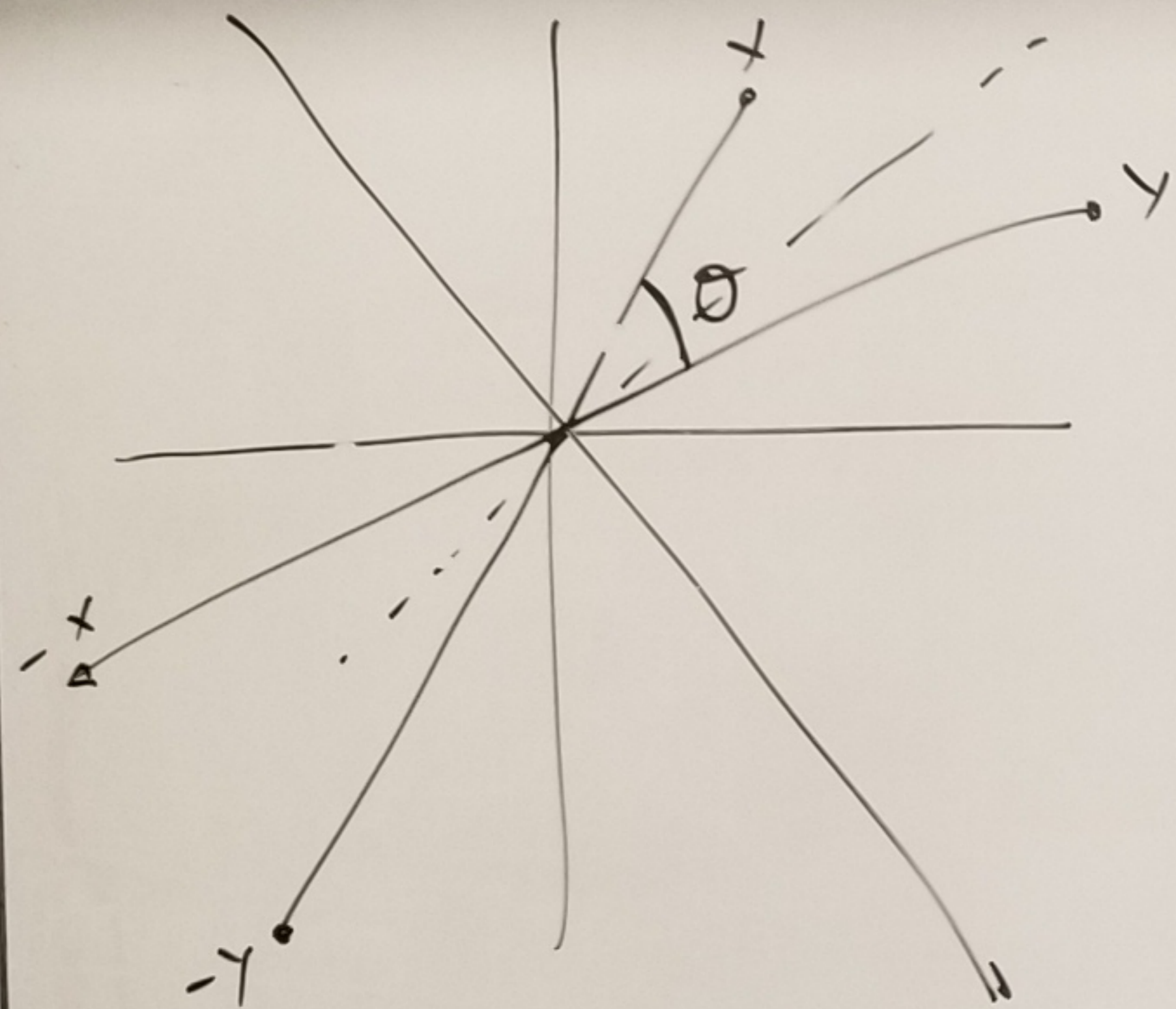
Jaccard \rightarrow Minwise Hashing.

Cosine Similarity between $x, y \in \mathbb{R}^D$

$$C(x, y) = \frac{x^T y}{\|x\|_2 \|y\|_2} \in [-1, 1] = \theta$$

\downarrow
vectors





Hash(x , seed)

generate $r_i \sim N(0, 1)$ using seed.

$$r_1, r_2, \dots, r_D = r$$

$$x \in \mathbb{R}^D$$

$$\text{Sign}(x^T r)$$

$$\Pr(h(x) \neq h(y)) = 1 - \frac{\theta}{\pi} = 1 - \frac{\cos^{-1}\left(\frac{x^T y}{\|x\| \|y\|}\right)}{\pi}$$