

Quiz 5: Similar Items

Name: GRADER ID: _____

Consider the following characteristic matrix of two sets: S1 and S2.

| Row # | S1 | S2 | H1 | H2 |
|-------|----|----|----|----|
| 0 | 1 | 0 | 1 | 2 |
| 1 | 1 | 1 | 2 | 3 |
| 2 | 1 | 1 | 3 | 4 |
| 3 | 0 | 0 | 4 | 5 |
| 4 | 1 | 1 | 5 | 6 |
| 5 | 1 | 0 | 6 | 0 |
| 6 | 0 | 1 | 0 | 1 |

1. [6 points] What are the two minhash values of S1 and S2 based on the permutation using $h_1(x) = (x + 1) \bmod 7$ (3pts) and $h_2(x) = (x + 2) \bmod 7$ (3pts). You need to calculate the minhash values with only one scan on the entire table.

1) mod 7 (3pts) and $h_2(x) = (x + 2) \text{ mod } 7$ (3pts). You need to calculate the value of h_2 for each element in the table after one scan on the entire table.

Row 0:

| | s_1 | s_2 |
|-------|----------|----------|
| h_1 | ∞ | ∞ |
| h_2 | ∞ | ∞ |

Row 1:

| | s_1 | s_2 |
|-------|-------|----------|
| h_1 | 1 | ∞ |
| h_2 | 2 | ∞ |

Row 2:

| | s_1 | s_2 |
|-------|-------|-------|
| h_1 | 1 | 2 |
| h_2 | 2 | 3 |

Row 3:

| | s_1 | s_2 |
|-------|-------|-------|
| h_1 | 1 | 2 |
| h_2 | 2 | 3 |

Row 4:

| | s_1 | s_2 |
|-------|-------|-------|
| h_1 | 1 | 2 |
| h_2 | 2 | 3 |

Row 5:

| | s_1 | s_2 |
|-------|-------|-------|
| h_1 | 1 | 2 |
| h_2 | 0 | 3 |

Row 6:

| | s_1 | s_2 |
|-------|-------|-------|
| h_1 | 1 | 0 |
| h_2 | 0 | 1 |

2. [4 points] Construct a signature for S1 and S2 based on the minhash values obtained from h1(x) and h2(x) above. Estimate the Jaccard similarity of S1 and S2 using the signature. What is the actual Jaccard similarity of S1 and S2 (2pt)? Is the estimate close to the actual Jaccard similarity? If not, suggest a way to improve the estimate (2pts).

Actual Jaccard similarity of s_1 and $s_2 = 3/6 = 1/2$

[1 point]

Taccard Similarity using Signature - 0/2
[Estimated]

[1 point]

NO, the estimate is not close to the actual Jaccard similarity. To improve, increasing the number of permutations or hash functions will help

[2 points]