

Question 2. LSH.

ci) document A ("the sky is blue the sun is bright")

B ("the sun in the sky is bright")

a) compute 2-shingles for A and B

b) compute their Jaccard similarity based on their 2-shingles.

Answer: a) $k=2$, set of 2-shingles

2-shingles for A: $\{(\text{the, sky}), (\text{sky, is}), (\text{is, blue}), (\text{blue, the}),$
 $(\text{the, sun}), (\text{sun, is}), (\text{is, bright})\}$

2-shingles for B: $\{(\text{the, sun}), (\text{sun, in}), (\text{in, the}), (\text{the, sky}),$
 $(\text{sky, is}), (\text{is, bright})\}$

the Jaccard similarity is $\text{sim}(A, B) = |A \cap B| / |A \cup B|$
 $= 4/8$

(ii)

$$h_1(n) = 5n - 1 \pmod{11}$$

$$h_2(n) = 2n + 1 \pmod{11}$$