

# Assignment 4 Solution

## Problem 1

For all keywords in the query: only Doc 3 is relevant for the query, and this doc contains all keywords. Thus,  $N = 3$ ,  $S = 1$ , and  $s = 1$ .

Obama:

$$df\_Obama = 3$$

$$c\_Obama = \log\left(\frac{(1+\frac{1}{2})/(1-1+\frac{1}{2})}{(3-1+\frac{1}{2})/((3-3)-(1-1)+\frac{1}{2})}\right) = -0.22$$

health:

$$df\_health = 2$$

$$c\_health = \log\left(\frac{(1+\frac{1}{2})/(1-1+\frac{1}{2})}{(2-1+\frac{1}{2})/((3-2)-(1-1)+\frac{1}{2})}\right) = 0.48$$

plan:

$$df\_plan = 2$$

$$c\_plan = c\_health = 0.48$$

Thus, the RSV scores for the 3 documents are:

$$RSV\_doc1 = -0.22 + 0.48 = 0.26$$

$$RSV\_doc2 = -0.22 + 0.48 = 0.26$$

$$RSV\_doc3 = -0.22 + 0.48 + 0.48 = 0.74$$

## Problem 2

### Parts 1 and 2

Priors:

$$P(e) = 1/9$$

$$P(\text{not } e) = 8/9$$

Conditional probabilities:

$$P(b|e) = 1/6$$

$$P(n|e) = 1/6$$

$$P(o|e) = 2/6$$

$$P(u|e) = 1/6$$

$$P(z|e) = 1/6$$

$$P(b|\text{not } e) = 1/6$$

$$P(n|\text{not } e) = 0/6, \text{ smoothed to: } 0.01$$

$$P(o|\text{not } e) = 1/6$$

$$P(u|\text{not } e) = 2/6$$

$$P(z|\text{not } e) = 2/6$$

Document scores:

$$\text{Score}(e|\text{zoo}) = 1/9 * 1/6 * 2/6 * 2/6$$

$$\text{Score}(\text{not } e|\text{zoo}) = 8/9 * 2/6 * 1/6 * 1/6$$

The document is classified as not English.

### Part 3

The vocabulary in this case is: o1, z2, b3, u1, u3, z1, o2, o3, b1, u2, n3.

Conditional probabilities:

$$P(z1|e) = 1/6$$

$$P(o2|e) = 1/6$$

$$P(o3|e) = 1/6$$

$$P(z1|\text{not } e) = 0/6, \text{ smoothed to: } 0.01$$

$$P(o2|\text{not } e) = 0/6, \text{ smoothed to: } 0.01$$

$$P(o3|\text{not } e) = 0/6, \text{ smoothed to: } 0.01$$

Document scores:

$$\text{Score}(e|\text{zoo}) = 1/9 * 1/6 * 1/6 * 1/6 = 0.00051$$

$$\text{Score}(\text{not } e|\text{zoo}) = 8/9 * 0.01 * 0.01 * 0.01 = 0.00000088$$

The document is classified as English.

### Problem 3

(best score per query in bold)

Query	Doc 1	Doc 2	Doc 3	Doc 4
click	0.46875	<b>0.71875</b>	0.21875	0.34375
shears	0.125	0.0625	0.0625	<b>0.1875</b>
click shears	0.0585	0.0449	0.0136	<b>0.0644</b>

$$P(\text{click}|\text{doc1}) = (4/8 + 7/16) / 2 = 0.46875$$

$$P(\text{click}|\text{doc2}) = (2/2 + 7/16) / 2 = 0.71875$$

$$P(\text{click}|\text{doc3}) = (0/2 + 7/16) / 2 = 0.21875$$

$$P(\text{click}|\text{doc4}) = (1/4 + 7/16) / 2 = 0.34375$$

$$P(\text{shears}|\text{doc1}) = (1/8 + 2/16) / 2 = 0.125$$

$$P(\text{shears}|\text{doc2}) = (0/2 + 2/16) / 2 = 0.0625$$

$$P(\text{shears}|\text{doc3}) = (0/2 + 2/16) / 2 = 0.0625$$

$$P(\text{shears}|\text{doc4}) = (1/4 + 2/16) / 2 = 0.1875$$

$$P(\text{click shears}|\text{doc1}) = 0.46875 * 0.125 = 0.0585$$

$$P(\text{click shears}|\text{doc2}) = 0.71875 * 0.0625 = 0.0449$$

$$P(\text{click shears}|\text{doc3}) = 0.21875 * 0.0625 = 0.0136$$

$$P(\text{click shears}|\text{doc4}) = 0.34375 * 0.1875 = 0.0644$$