Perform the boolean query Vikings AND Spam given the following inverted index:

- Vikings: 2, 3, 5, 7, 11, 13, 17, 23, 29, 31, 37, 41, 47, 53, 59, 61

• Spam: 5, 13, 17, 23, 29, 31, 43, 47, 59

Present the result of the query [1 point] and the number of comparisons with [2 points], and without [2 points] a *skip list*. Assume the skip list has frequency $\lfloor \sqrt{|P|} \rfloor$, where |P| is the cardinality of a posting P and |.| is the floor function. (Example: $P = \{1, 2, 3\}, |P| = 3, \sqrt{|P|} \approx 1.73, \lfloor \sqrt{|P|} \rfloor = 1$)

Results of the query: 5, 13, 17, 25, 29, 31, 47, 59 Number of comparisons.

$$(37, 43), (41, 43), (47, 43), (47, 47), (53, 59), (51, 59)$$

$$(2,5), (3,5), (5,5), (7,13), (11,13), (13,13), (17,17),$$

Compute an unbiased estimate of a text retrieval system's precision, recall, and the F_1 measure on the first five results [3 points], and the precision at 40% recall [2 points] given the following lists of results for queries q_1 , and q_2 , where R is a relevant result, and N is a non-relevant result:

- Results for q_1 : RNNRRNR (10 relevant results for q_1 exist in the collection.)
- Results for q_2 : NRNRRRN (5 relevant results for q_2 exist in the collection.)

First five results:
$$Pq_{1} @ 5 = \frac{3}{5}, Pq_{2} @ 5$$

$$- Results for q_{1}: RNNRR
$$- Results for q_{2}: NRNRR$$

$$Macro-averaging P_{1}R$$

$$P@ 5 = (P_{1} @ 5 + P_{1} @ 5)/2 = \frac{3}{5} = 0.6$$

$$R@ 5 = (R_{1} @ 5 + R_{1} @ 5)/2 = \frac{3}{20} = 0.45$$

$$R@ 5 = \frac{3+3}{5+5} = \frac{3}{5} = 0.6$$

$$R@ 5 = \frac{3+3}{5+5} = \frac{3}{5} = 0.6$$

$$R@ 5 = \frac{3+3}{5+5} = \frac{6}{15} = \frac{2}{5} = 0.4$$$$

$$P_{41} @ 5 = \frac{3}{5}$$
 $P_{42} @ 5 = \frac{3}{5}$
 $K_{41} @ 5 = \frac{3}{10}$
 $K_{42} @ 5 = \frac{3}{5}$

Micro-averaging P, R

Pa
$$5 = \frac{3+3}{5+5} = \frac{5}{5} = 0.6$$

Ra $5 = \frac{3+3}{10+5} = \frac{6}{15} = \frac{2}{5} = 0.4$

Macro-averaging
$$F_1$$

$$F_1 = \frac{2 \cdot 1005}{1005} \cdot 1005$$

$$= \frac{0.54}{1.05} = 0.5143$$

Micro-averaging
$$F_1$$

 F_1 $a 5 = 2 \cdot R_4$ $a 5 \cdot P_4$ $a 5$
 R_4 $a 5 + P_4$ $a 5$
 $= \frac{0.36}{0.4} = 0.4$

Macro-averaging
$$F_1$$
 | Micro-averaging F_1 | Macro-averaging F_2 | Micro-averaging F_3 | $F_4 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 + P_0 \cdot 5}$ | $F_4 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 + P_0 \cdot 5}$ | $F_4 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 + P_0 \cdot 5}$ | $F_4 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 + P_0 \cdot 5}$ | $F_5 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 + P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | $F_6 = \frac{2 \cdot R_0 \cdot 5 \cdot P_0 \cdot 5}{R_0 \cdot 5 \cdot P_0 \cdot 5}$ | F_6

$$F_{192} = \frac{2 \cdot R_{91} \cdot \Omega_{5} \cdot P_{91} \cdot \Omega_{5}}{R_{91} \cdot \Omega_{5} + P_{91} \cdot \Omega_{5}}$$

$$= \frac{0.72}{1.2} = 0.6$$

$$\frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{1} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{1} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{2} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{3} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{4} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{5} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{1} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{2} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{3} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{4} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{5} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{5} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{7} = \frac{1}{1,2} = \frac{1}{1,2}$$

$$F_{7$$

192 R9, @5+P9, @5

Results with 40% recall:

$$P_{q_1} = \frac{4}{7}$$

$$R_{q_2} = \frac{1}{2}$$

$$R_{q_2} = \frac{1}{2}$$

$$R_{q_2} = \frac{2}{5}$$

$$Micro-averaging P$$

$$P = \frac{4+1}{5} = \frac{5}{5} = 0.5$$

Macro-averaging P

$$P = (Pq_1 + Pq_2)/2 = \frac{15}{28} = 0.5357$$

Define the two assumptions the *Naive Bayes* classifier makes [2 points]. Explain the advantage of computing a product of probability estimates as a sum in the logarithmic space [1 point]. Given an observation x, and the classes c_1 , and c_2 , is the knowledge of $P(\mathbf{x} \mid c_1)P(c_1) > P(\mathbf{x} \mid c_2)P(c_2)$ sufficient to decide whether $P(c_1 \mid \mathbf{x}) > P(c_2 \mid \mathbf{x})$? Why or why not? [2 points]

Given the following list of observations, use the Naive Bayes classifier to decide whether to play golf when it is sunny, hot, windy, and the humidity is normal. [5 points]

Outlook	Temperature	Humidity	Windy	Play golf
Sunny	Mild	High	False	Yes
Sunny	Mild	Normal	False	Yes
Overcast	Hot	High	False	Yes
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	True	No
Rainy	Hot	High	False	No
Sunny	Cool	Normal	True	No
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Hot	High	True	No
Rainy	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	True	Yes

 $P(x_{i+1}, x_{i+2}, ..., x_{h}, c_{b}) = P(x_{i} | c_{b})$ Numerical stability of multiplying small real numbers.

By Bayes theorem, $P(c_1|\vec{x}) = \frac{P(\vec{x}|c_1)P(c_1)}{P(\vec{x})}$ and $P(\vec{x}|c_2)P(c_2)$ $P(c_2|\vec{x}) = \frac{1}{P(\vec{x})}$ · Clearly, P(x/4) P(c1) > P(xlez) P(cz) implies P(c1/x) > P(cz/x), and vise versa. P(Yes | Shing, Hot, Normal, True) = P(Shing | Yes). P(Hot | Yes).

• P(Normal | Yes) • P(True | Yes) = $\frac{3}{4} \cdot \frac{2}{9} \cdot \frac{6}{9} \cdot \frac{3}{9} \cdot \frac{9}{19} = \frac{972}{91859}$ P(Wolsung, Hot, Wormal, True) = P(Sunny IWO) . P(Hot INO). · P(Normal | No) · P(True | No) = $\frac{2}{5}$ · $\frac{2}{5}$ · $\frac{1}{5}$ · $\frac{3}{5}$ · $\frac{5}{14}$ = $\frac{60}{17500}$ P(Tes | Sunny, Hot, Wormal, True) > P(Nol Sunny, Hot, Normal, True). Therefore, play 20lf.

Given a directed graph G that represents three Web pages $V(G) = \{a, b, c\}$, and the links $E(G) = \{a, b, c\}$ $\{(b,a),(c,a),(c,b),(b,c)\}$ between these three pages, draw G [1 point] and produce the adjacency matrix (also known as the link matrix) A [1 point], and the Markov transition matrix P [2 points].

Describe the intuition behind the PageRank algorithm [1 point]. Compute the PageRank of the pages a, b, and c using a single iteration of the PageRank algorithm [2 points].

Describe what we mean, when we call a page a hub, or an authority [1 point]. Compute the hub, and authority scores of the pages a, b, and c [2 points].

$$A = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

$$P = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix} \circ \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/2 & 1/2 & 1/2 \end{bmatrix} \cdot (1-d) + d \cdot \frac{1}{3}$$

Figure: graph a

where . O. is the Hadamard product.

The Page Rand algorithm compuses the probability that a hypo-He ticul vandon surfer will end up at a given web page. $\vec{x}_0 = (100)$ $\vec{x}_1 = \vec{x}_0 \cdot P = [100]$ $\begin{bmatrix} 1/3 & 1/3 \\ ... & 1/3 \end{bmatrix}$

$$\frac{1}{1/3} = \frac{1}{1/3} \cdot P = [100] \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ & & & & \\ & & & \\ & & & & \\ &$$

A hub is a web page pointing to many authorities. An authority is a web page that many hubs point to $AA^{T} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix}$ 2 $\frac{1}{h_{1}} = (AA^{T}h_{0}) = [033]^{T}/3 = [0117]^{T}$ $a_1 = (A^T A \vec{a}) = [422])/4 = [41/21/2]^T$