

CSE 586 Final Exam

Olcay Taner YILDIZ

I. QUESTION (15 POINTS)

Write the algorithm that handles the boolean query $x \Delta y$ (symmetric difference) using inverted index.

II. QUESTION (16 POINTS)

Let say there are N documents in a collection. Given that the terms x and y occur in n_x and n_y distinct documents respectively, how many results will be returned from these queries (minimum and maximum)?

- x and y
- x or y
- x and not y
- x or not y

III. QUESTION (15 POINTS)

Let say we have a collection of M documents with N terms. What will be the time complexity of

- Learning the mean vectors in the Rocchio algorithm with K classes?
- Testing a new document in the 1-NN algorithm?
- Finding which cluster the document belongs to if we have used K -means clustering?

IV. QUESTION (12 POINTS)

Suppose that the task is to cluster the following eight points (with (x, y) representing location) into three clusters. $A_1(2, 10)$, $A_2(2, 5)$, $A_3(8,4)$, $A_4(5, 8)$, $A_5(7,5)$, $A_6(6,4)$, $A_7(1,2)$, $A_8(4,9)$. The distance function is Euclidian distance. Suppose initially we assign A_1 , A_4 , and A_7 as the center of each cluster, respectively. Run the k-means clustering one iteration on these data.

V. QUESTION (14 POINTS)

docID	words in document	in Turkey?
1	Turkish Ankara Turkish	yes
2	Turkish Turkish Adana	yes
3	Turkish Turkish Adana Istanbul	yes
4	Athens Greece	no
5	Athens Greece Turkish	no
6	Athens Greece Selanik	no

Based on the data given above,

- Estimate a multinomial Naive Bayes classifier
- Apply classifier to the document
Turkish Turkish Turkish Athens Greece

VI. QUESTION (14 POINTS)

Based on the data given above

- Implement Rocchio classifier, that is find the mean vector of each class.
- Apply classifier to the document
Turkish Turkish Turkish Athens Greece

VII. QUESTION (14 POINTS)

Based on the data given above, estimate the class of the document

- Turkish Turkish Turkish Athens Greece
for
- $k = 1$
 - $k = 3$