

**START OF QUIZ**

**Student ID:**

**68022458,Chan,Douglas**

## Question 1

Topic: Lecture 5

Source: Lecture 5

Why do we need methods like t-SNE? (1)

## Question 2

Topic: Lecture 6

Source: Lecture 6

In class, we talked about bookstores and streaming algorithms classifying books / movies. How can we tell that they don't use a topic modeling algorithm (or, if you think they do, what would be some clues)? (1)

### Question 3

Topic: Lecture 7

Source: Lecture 7

Explain why the cosine similarity between a document and query vector is roughly equivalent to adding up the TF-IDF scores of each word in the document that occurs in the query.  
(2)

## Question 4

Topic: Lecture 5

Source: Lecture 5

The Frobenius norm looks very similar to a distance metric we've already observed. Explain which one. (1)

## Question 5

Topic: Lecture 8

Source: Lecture 8

Why don't we use a higher-order language model to perform IR? (1)

## Question 6

Topic: Lecture 6

Source: Lecture 6

Why do we need a "human in the loop" for topic modeling? (1)

## Question 7

Topic: Lecture 7

Source: Lecture 7

Explain why boolean filtering is usually insufficient for retrieval, and why we normally need some way of scoring the documents. (2)



## Question 8

Topic: Lecture 8

Source: Lecture 8

In class, I mentioned that high  $k$  value for BM25 TF weighting rewards documents with many, many instances of a term in them. Explain why that's the case. (2)

## Question 9

Topic: Long

Source: Lecture 5

Imagine that we are working with a language other than English, such as Indonesian, with significant agglutinative morphology (words are inflected through the concatenation of affixes to a lemma). How do you think that this would impact our various vector space models? Which of them would be most affected, and which would be least affected? Explain. (3)

**END OF QUIZ**