

**START OF QUIZ**

**Student ID:**

**37289428, Yun, Michelle**

## Question 1

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 2

Topic: Lecture 8

Source: Lecture 8

$P(d|q)$  is not what we are solving with the language model. Why is this not generally a problem? (1)

### Question 3

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 4

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 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 5

Source: Lecture 5

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

## Question 7

Topic: Lecture 6

Source: Lecture 6

Why can't we just run an HMM over documents to discover the latent states like we do for POS-tagging? (1)



## Question 8

Topic: Lecture 5

Source: Lecture 5

Why can we represent a rank- $m$  matrix as the sum of  $m$  rank-1 matrices \*or\* the product of an  $n \times m$  matrix and an  $m \times n$  matrix (ie, what is matrix multiplication doing that we can take advantage of)? Explain. (2)

## Question 9

Topic: Coding

Source: Coding

Write a short function that confirms that the sum of  $n$  rank-1 matrices is identical to the product of an  $n \times k$  matrix and a  $k \times n$  matrix. (3)

**END OF QUIZ**