

START OF QUIZ

Student ID:

**95174918, Maurin-
Jones, Kai**

Question 1

Topic: Lecture 5

Source: Lecture 5

Explain the logic behind the IDF part of TF-IDF (ie, why does it give higher weights to more "interesting" words?). (1)

Question 2

Topic: Lecture 6

Source: Lecture 6

Why don't we just use k-means to cluster document-vectors (sparse or dense)? (1)

Question 3

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 4

Topic: Lecture 5

Source: Lecture 5

Why can we be confident that a low-rank approximation of a matrix contains the most important information in a document? (1)

Question 5

Topic: Lecture 8

Source: Lecture 8

What do we mean by interpolation? (1)

Question 6

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 7

Topic: Lecture 6

Source: Lecture 6

In some ways, we could consider Theta distributions themselves to be an embedding of a topic. Explain, and explain how we might be able to leverage that. (2)

Question 8

Topic: Lecture 8

Source: Lecture 8

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

Question 9

Topic: Coding

Source: Coding

Write a function that returns the most likely n documents given a term-document matrix, a smoothing parameter, and a query. (3)

END OF QUIZ