

Short Problem Set (Module 6)

1. [30%] For this problem we will use Cover Density Ranking for the following document which is two verses from an 18th century English nursery rhyme. The numeric superscripts indicate the word order in the document. For this problem we have this one document and a query consisting of the two words "**sing pie**".

sing¹ a² song³ of⁴ sixpence⁵
a⁶ pocket⁷ full⁸ of⁹ rye¹⁰
four¹¹ and¹² twenty¹³ blackbirds¹⁴
baked¹⁵ in¹⁶ a¹⁷ pie¹⁸

when¹⁹ the²⁰ pie²¹ was²² opened²³
the²⁴ birds²⁵ began²⁶ to²⁷ sing²⁸
wasn't²⁹ that³⁰ a³¹ dainty³² dish³³
to³⁴ set³⁵ before³⁶ the³⁷ king³⁸

- (a) Is (1,21) a *cover* for this query? Explain why or why not?
- (b) List the covers for this query.
- (c) Using a window size of $K=8$, calculate the similarity score for the document.

Cover Density Ranking is not discussed in the course text. However, there is an example in the lecture slides, and you can find the paper "*Relevance ranking for one to three term queries*" by Clarke et al., in the EReserves in Blackboard. That paper formally defines a cover and gives several examples that you may find useful.

2. [20%] In the statistical language model presented in the lecture and in Chapter 12 of the text we use linear interpolation (also called a "mixture model" or "Jelinek-Mercer smoothing") to make a probability estimate of a term. This estimate is based both on the term frequency in a document, and on the collection frequency of the term. See Equation 12.12 in IIR.

- (a) What is the purpose of the parameter λ ?
- (b) What would be the effect of setting λ to a value of 1?

3. [50%] Compute similarity scores for and rank documents D1 and D2 using query Q with a unigram statistical language model. Query Q contains the four words "aardvark bird dog dog". The corpus consists of only these eight documents and only these six indexing terms are found in the collection. The cells in the table below indicate the number of times a word appears in a document. You should use a mixture model with parameter $\lambda = 0.40$. Assume that the prior probability of relevance is equal for all documents.

Plainly show your work. It is fine to check your work with a program or spreadsheet, but I expect you to show how you derive probability estimates and to see the equations that you use to calculate document similarity.

	D1	D2	D3	D4	D5	D6	D7	D8
aardvark	5	2		1				2
bird	1	1	1	2	1	1	5	8
cat				2		3		
dog		1	3		2		2	
egret	1			1				
fish	3	2						

Report scores using scientific notation with three digits after the decimal point (e.g., 1.234×10^{-8}).