

Homework 1

Student Number:

Name:

Problem 1. (30 points)

Doc 1 new home sales top forecasts

Doc 2 home sales rise in july

Doc 3 increase in home sales in july

Doc 4 july new home sales rise

Consider the documents above,

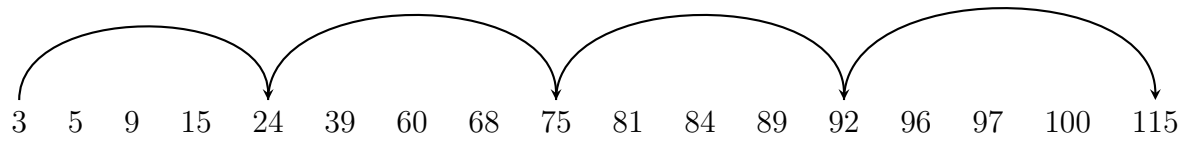
- a. Draw the term-document incidence matrix for this document collection.
- b. Draw the inverted index representation for this collection.
- c. For the document collection, what are the returned results for these queries:
 - i july AND rise
 - ii (NOT increase) AND (home OR sale)

Problem 2. (30 points) The following pairs of words are stemmed to the same form by the Porter stemmer. Which pairs would you argue shouldn't be conflated. Give your reasoning.

- a. abandon/abandonment
- b. absorbency/absorbent
- c. marketing/markets
- d. university/universe
- e. volume/volumes

Problem 3. (20 points) Write a query using Westlaw syntax which would find any of the words *professor*, *teacher*, or *lecturer* in the same sentence as a form of the verb *explain*.

Problem 4. (30 points) Consider a postings intersection between this postings list, with skip pointers:



and the following intermediate result postings list (which hence has no skip pointers):

3 5 89 95 97 99 100 101

Trace through the postings intersection algorithm(pdf of lecture 1, under section Skip Pointers)

- How often is a skip pointer followed?
- How many postings comparisons will be made by this algorithm while intersecting the two lists?
- How many postings comparisons would be made if the postings lists are intersected without the use of skip pointers?