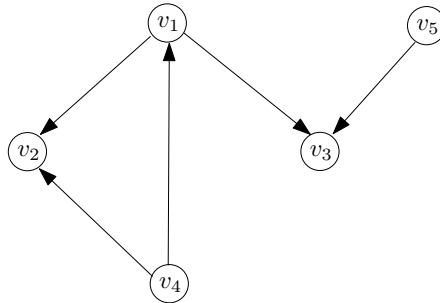


## WST540: Exercise List 3

**Problem 1.** Consider the following graph:



Each node represents a website and each edge represents a hyperlink. Suppose that a web crawler initially knows only  $v_4$ . Using the BFS algorithm discussed in Lecture 4, which websites will be discovered by the crawler when it finishes?

**Problem 2.** Suppose that we have the following document collection:

document ID	content
1	the old night keeper keeps the keep in the town
2	in the big old gown in the big old house
3	the house in the town had the big old keep
4	where the old night keeper never did sleep
5	the night keeper keeps the keep in the night
6	and keeps in the dark and sleeps in the light

Also, consider that our dictionary has words {big, dark, gown, house, keep, light, night, old, sleep, town}. Assume that, after stemming, we have the following equivalence:

$$\begin{aligned} \text{keep} &= \text{keeper, keeps} \\ \text{sleeps} &= \text{sleep} \end{aligned}$$

Give the document-level inverted lists of all the words in the dictionary. Each entry of an inverted list should have the format (doc id, term frequency).

**Problem 3.** Give the Elias' gamma and delta codes of 23.

**Problem 4.** Consider the following inverted list, where each entry is in the format of (doc id, term freq):

$$(1, 1), (4, 1), (5, 2)$$

Give the bit sequence that compresses the above list based on the following ideas:

- For the  $i$ -th ( $i \geq 2$ ) pair, represent its id by storing in Elias' gamma code the difference from the id of the  $(i - 1)$ -th pair.
- Store each term-frequency value in Elias' gamma code.

**Problem 5.** Give the word-level inverted lists for Problem 3. Each entry of an inverted list should have the format (doc id, term frequency, position 1, position 2, ...).