

Problem 62.

Nodes	m	n	o	p	q	r	s	t	u	v	w	x	y	z	remaining count
# Edges	0	0	2	0	2	3	2	2	2	2	1	2	1	2	14
/ = deleted-edge	/	0	2	0	1	2	2	2	2	2	1	1	1	2	13
(o)	/	/	1	0	0	2	2	2	1	2	1	1	1	2	12
	/	/	0	1	0	2	1	2	1	2	1	1	1	1	11
	/	/	1	1	0	1	0	2	1	1	1	1	1	1	10
	/	/	1	1	1	1	1	0	1	1	1	1	1	1	9
	/	/	1	1	1	1	0	1	1	1	1	1	1	1	8
	/	/	1	1	1	1	1	1	0	1	1	1	0	1	7
	/	/	1	1	1	1	1	1	1	0	1	1	1	0	6
	/	/	1	1	1	1	1	1	1	1	1	0	1	1	5
	/	/	1	1	1	1	1	1	1	1	1	0	1	1	4
	/	/	1	1	1	1	1	1	1	1	0	0	1	1	3
	/	/	1	1	1	1	1	1	1	1	0	0	1	0	2
	/	/	1	1	1	1	1	1	1	1	1	1	0	1	1
	/	/	1	1	1	1	1	1	1	1	1	1	1	0	0
	/	/	1	1	1	1	1	1	1	1	1	1	1	1	0

So, we can see that all the edges can be systematically removed and re-organized, without any cycling occurring.

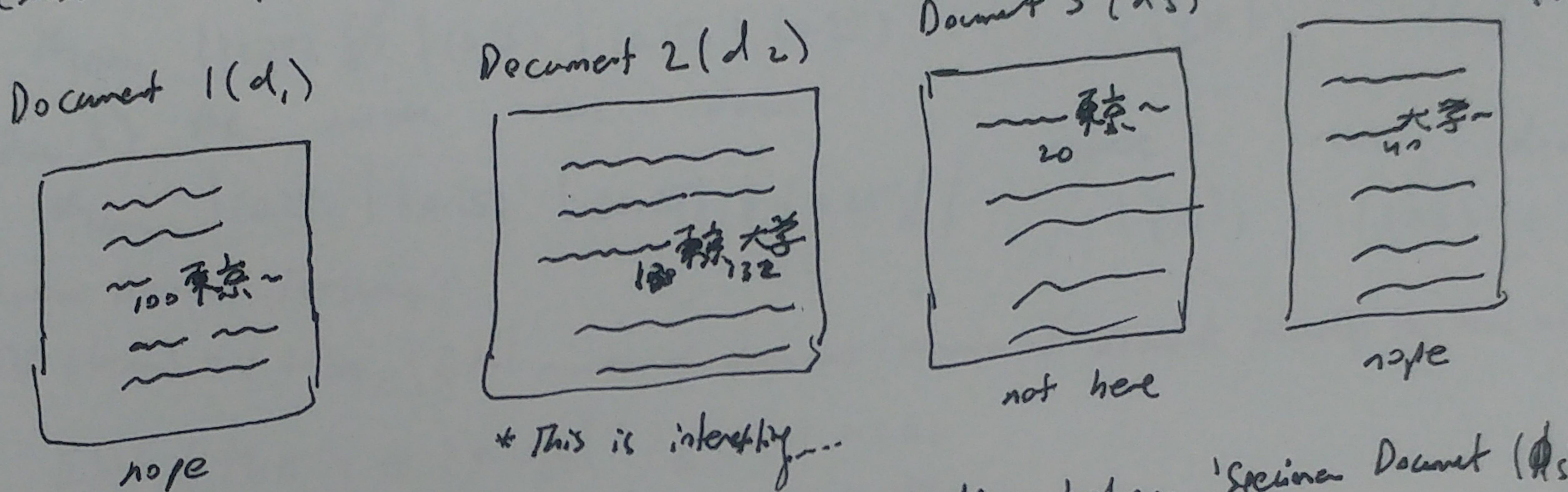
This figure will not encounter any cycling pattern.

Alex Taniguchi
08-174510
Database HW #6 11.1

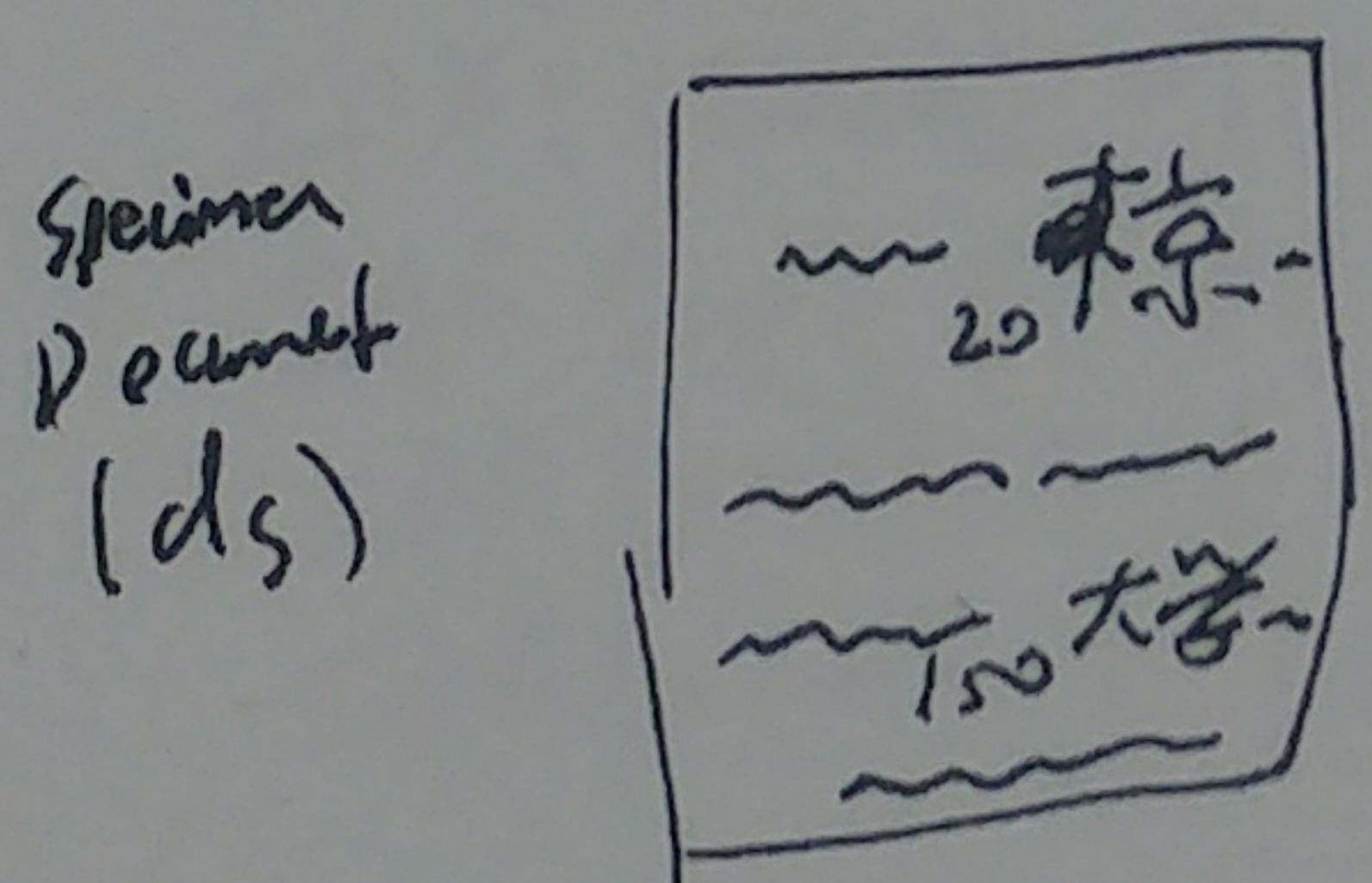
Problem 71. At its core, the problem is neatly sortable as an intersection problem, where we want the Boolean set of documents (and positions) for "Tokyo" (東京) AND "University" (大学). We would normally start with the term known to have less matches and then search within those smaller subset of matches for the term which had a larger subset of matches.

It is worth noting here that if the answer returned is $\{\emptyset\}$ (empty) then we care except from worries about position within the document, as having the term "東京大学" (Tokyo University) is impossible.

(Shall we try?)



Now, my case above cannot handle the difference between 'Specimen Document (d_3)' and Document 2 (d_2) above.



Even with 130-words between 'Tokyo' and 'University', the overlap is enough to differ my boolean algorithm from above. How do I ensure that I have '東京' and '大学' next to each other?

The easiest (yet $O(n)$ -time most expensive) solution is to make a positional index proximity-search algorithm with an extra-variable ' k '. While adjacency in English is '1 word' so $k=1$ in that case, adjacency in Japanese kanji by position is dependent on the length of the kanji. A very expensive algorithm could calculate reading the length of your term (here, $k=2$) for you.

In the case of a posting index structure of max-logarithm 'N', then individual search terms are solved in individual logarithmic time - $O(\log N)$, and the boolean search is $N^{\log N}$ scans so it would be $O(N \log N)$. The position-adjacency check adds another (in generalized form) - " $O((\log N)^2)$ " relative search, so the total time for my algorithm above is an admittedly dismal $O((N+1) \log N)$.

72 (continued). It is worth noting that we now have our numbers necessary for creating the vectors, that is:

$$d_1 = \langle 0.33, 0.16, 0.16, 0.16, 0.16 \rangle$$

$$d_2 = \langle 0.25, 0.25, 0.25, 0.25, 0.25 \rangle \text{ for TF - regular case - } ①$$

$$d_3 = \langle 0.25, 0.25, 0.25, 0.25, 0.25 \rangle$$

$$d_1 = \langle 0.25, 0.25, 0.25, 0.25, 0.25 \rangle \text{ filler 1's should be } 0 \text{ (?)}$$

$$d_2 = \langle 0.33, 0.33, 0.33, 0.33, 0.33 \rangle \text{ for TF - stop words removed case - } ②$$

$$d_3 = \langle 0.33, 0.33, 0.33, 0.33, 0.33 \rangle \text{ for TF - stop words removed case - } ③$$

$$d_1 = \langle 0.8833, 0.2343, 0.2343, 0.3498, 0.3498 \rangle$$

$$d_2 = \langle 0.1667, 0.5247, 0.5247, 0.5247, 0.5247 \rangle \text{ for TF-IDF - regular case - } ④$$

$$d_3 = \langle 0.1667, 0.3514, 0.3514, 0.5247, 0.5247 \rangle$$

$$d_1 = \langle 0.3514, 0.3514, 0.5247, 0.5247 \rangle$$

$$d_2 = \langle 0.6995, 0.6995, 0.6995 \rangle \text{ for TF-IDF - stop words removed case - } ⑤$$

$$d_3 = \langle 0.4685, 0.4685, 0.6995 \rangle \text{ for TF-IDF - stop words removed case - } ⑥$$

$$\text{in the } \cos(d_1, d_2, d_3) = \frac{d_1 \cdot d_2 \cdot d_3}{\|d_1\| \|d_2\| \|d_3\|}$$

$$\text{case } ① \quad \cos(d_1, d_2, d_3) = \frac{0.050625}{(0.4597)(0.5)(0.5)} \approx \frac{0.050625}{0.114918 \dots} \approx 0.4405 \rightarrow 44.05 \text{- cos similarity}$$

$$\text{case } ② \quad \cos(d_1, d_2, d_3) = \frac{0.081675}{(0.5)(0.57157 \dots)^2} \approx 0.5 \rightarrow 50.00 \text{- cos similarity}$$

$$\text{case } ③ \quad \cos(d_1, d_2, d_3) = \frac{0.191945777707}{(0.6824)(0.9239)(0.7717 \dots)} \approx 0.4105 \rightarrow 41.05 \text{- cos similarity}$$

$$\text{case } ④ \quad \cos(d_1, d_2, d_3) = \frac{0.477176616275}{(0.8931)(1.2116)(0.9635)} \approx 0.4577 \rightarrow 45.77 \text{- cos similarity}$$

There are a couple of interesting points, esp. where case ① (TF - regular) is a near-match to the TF-IDF score when reduced by stop words.

Considering that a targeted TF (remove stop words) has highest similarity, the IDF was not kind in having repeated words over short documents.

As this is a children's poem with repetition and short structure, the 50+ish similarity between all three lines is not a surprise. Indeed, with more time it would be fun to calculate similarities for other children's rhymes.