points

Describe the SoundEx phonetic retrieval algorithm [2 points]. Give an example of two similarly-sounding words with the same SoundEx code [1 point]. Explain the weak point of SoundEx [1 point], and give an example of two different-sounding words with the same SoundEx codes [1 point].

Generally spenting. He Sound Ex algorithm naps similarly--sounding words together by removing consonants, and clustering vowels with similar English spelling together.

Juond and short have code 5630.

The major weak point of SoundEx is the fact that the initial letter of a word is retained.

Two different-sounding words with the same SoundEx codes are Grefan (S315), and Greven (S315).

## Consider the following XML document:

Menu item item stee note price little note price le l'alle note price l'Esser..." "Wish a..." 19.9

Draw the XML document as a graph [1 point], and count all *structural terms* in the document [2 points]. Compute the *structural similarity* between the structural term item/title#"That's not got much spam in it", and the queries //menu//price#30.9, and //item/title#"Lobster Thermidor" [2 points].

Structural terms: # "Egg,...", title# "Egg,...", item/title# "Egg,...", hote # "that's...", hote # "that's...", hote # "that's...", # 4.9,9 price # 9.9, item/price # 9.9, wenn/item/price # 9.9, # "Lobster..." # 106ster..." witle # "Lobster...", item/title # "Lobster...", # "Wish a...", nose # "Wish a...", item/nose # "Wish a...", mose # "Wish a...", # 40.9, price # 49.9, item/price # 49.9

6.4 = 24 structural terms.

/meny /price does not expand to item/title, the structural similarity is zero. //item/title is equivalent to item/title, so the sweetheral similarity is  $\frac{1+1}{1+1} = 1$ .

sheet

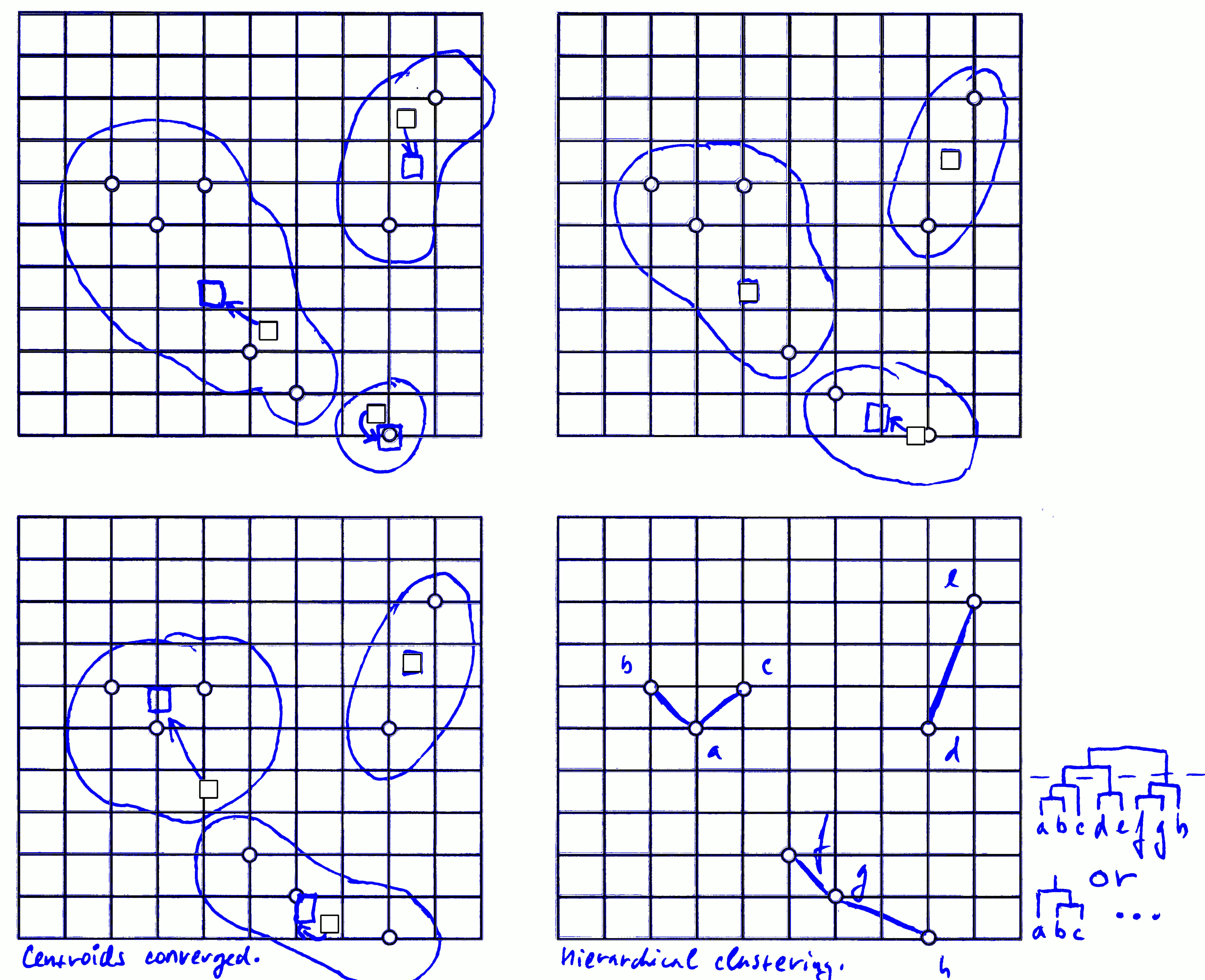
učo . . . . . . . . . . . . . point

Explain the following aspects of the K-means flat clustering algorithm [2 points]:

- 1. What do we need to know about our dataset before using the algorithm?
- 2. What is the input and the output of the algorithm?
- 3. What are the two steps that take place in every epoch?
- 4. How do we decide in which epoch to stop the algorithm?

1. We need to know the humber of elasses K and initial mean estimates (seeds). 2. The input are unclassified points and K seeds. S. Reassigning points, recomputing centroids. 4. Centroids converged.

Given the points O, and the seeds  $\square$ , run the K-means algorithm for three epochs. Draw the state of the algorithm at the beginning and after every epoch; no computation should be necessary. What is the output of the algorithm? [2 points]



Perform a hierarchical clustering of the above dataset into three classes using the single-link hierarchical agglomerative clustering algorithm, and draw the resulting dendrogram. [1 point] Is the output the same as the output of the K-means flat hierarchical clustering algorithm above? [1 point]

762 14 12

State the Zipf's law [1 point] and the Heaps' law [1 point] in its complete variant [1 point]. Let C denote the term-document matrix of a document collection that contains M unique terms. According to the Heaps' law, what is the number of documents N in the collection in relation to M? [3 points]

Zipfs law : cfi = c/i , where cfi is the collection frequency of the ith most frequent term, and e is a cohstant.

Heaps law & M= kTb, where Mis the vocabulary Size, Tis the collection size in tokens, and  $30 \le k \le 100$  and  $6 \approx 0,5$  are parameters.

the size of C is M x N < T, where N is He humber of documents in a collection. According to the Heaps! Care,  $T = \sqrt[4]{H/k} \propto M^2$ . There fore  $M \times N < M^2 \Rightarrow N < M$ .

You are the maintainer of a text retrieval system. Let  $E_1$  denote the complete set of documents in the index of your system and let  $E_2$  denote the complete set of documents in the index of a competing system. Suppose the indices of both systems are independent uniform random samples without replacement from the World Wide Web N. The size of  $E_1$  is  $|E_1| = 130$  trillion (130 · 10<sup>12</sup>) documents. You take a uniform random subsample of documents without replacement from  $E_1$  and you submit each document to the competing system. This gives you an unbiased estimate x = 0.2of the conditional probability  $P(d \in E_2 \mid d \in E_1), d \in N$ . You repeat the same procedure with  $E_2$ , obtaining an unbiased estimate y = 0.4 of the conditional probability  $P(d \in E_1 \mid d \in E_2), d \in N$ . Assume the estimates x, y are the true probabilities. What is the size  $|E_2|$  of the competing system's index? [4 points]

The grey parrot, native to equatorial Africa, is categorized as an endangered species by the International Union for Conservation of Nature (IUCN). Suppose you take a uniform random sample M without replacement of size |M| = 8000 from the grey parrot population N and mark the sampled animals. After returning the marked animals back into the population, you take a second independent uniform random sample T without replacement of the same size  $|T| = 8\,000$  from the population. The number of marked animals  $R = M \cap T$  in the second sample is |R| = 10. What is the size |N| of the population of the grey parrot? [4 points]

$$\frac{1}{1} + \frac{1}{1} + \frac{1$$