

Part I : Text Processing

[25 marks in total]

1. **Regular Expressions:** Given the following “regular expression”:

 $\backslash s*f?re\{1,3\}\backslash w$

For each of the following strings, indicate (with “yes” or “no” in your script book) whether the regular expression will match. [2 marks]

- (a) freedom Yes
- (b) a really good book Yes
- (c) sfrew Yes
- (d) \tf\tfr\tfre\tfree Yes

2. **Approximate String Search:** By referring to the definitions from this subject:

- (a) Explain the problem that we wish to solve in “approximate string search”, including one typical strategy and any necessary data structure(s). [2 marks]
- (b) Explain why “approximate string search” is a “Knowledge Technology” (or “knowledge task”). [2 marks]

(a) the problem is given a string, how to find a approximately match string from the dictionary, for example, to solve the problem of spelling correction and computational genomics. Strategy: N-grams, Patricia structures?

(b) The solution of it is ill-defined and ambiguous since it is not an exact string match method, the result may be different, and the context is critical for example, the exact word a user want is different in different context no single solution provided and computer media between users and data

continued ...

3

(a) it $\rightarrow A(1)$ happens $\rightarrow A(1)$ over $\rightarrow A(3) \rightarrow B(1) \rightarrow C(4)$ and $\rightarrow A(2) \rightarrow B(2)$ again $\rightarrow A(1)$ under $\rightarrow B(1)$ in $\rightarrow B(1)$ out $\rightarrow B(1)$

$$\begin{bmatrix} 1, 1, 3, 2, 1, 0, 0, v \\ 0, 0, 1, 2, 0, 1, 1, 1 \end{bmatrix}$$

$$\begin{bmatrix} 0, 0, 4, 0, 0, 0, 0, 0 \\ 1, 1, 1, 1, 1 \end{bmatrix}$$

$$\begin{bmatrix} 0, 0, 4, 0, 0, 0, 0, 0 \\ 1, 1, 1, 1, 1 \end{bmatrix}$$

$$\frac{1+1+9+4+1}{4}$$

$$\frac{9+4+36}{4}$$

(b) ① Fetch the inverted list for each query

Over $\rightarrow A(3) \rightarrow B(1) \rightarrow C(4)$ and $\rightarrow A(2) \rightarrow B(2)$ out $\rightarrow B(1)$

$$4+4$$

② the conjunction may be : over AND and AND out
so the intersection of these three list is {B1}

$$(c) W_1 = [1, 1, 3, 2, 1, 0, 0, 1]$$

$$W_2 = [0, 0, 1, 2, 1, 1, 1, 1]$$

$$W_3 = [0, 0, 4, 0, 0, 0, 0, 0]$$

$$W_q = [0, 0, 1, \frac{3}{2}, 3, 0, 0, 0]$$

$$\cos(1, q) = \frac{3+3+3}{4 \cdot \frac{7}{2}} = \frac{9}{4 \cdot \frac{7}{2}}$$

$$\cos(2, q) = \frac{1+3+3}{3 \cdot \frac{7}{2}} = \frac{7 \cdot 4}{4 \cdot \frac{7}{2}}$$

$$\cos(3, q) = \frac{4}{4 \cdot \frac{5}{2}}$$

 $\therefore \text{rank: doc2} > \text{doc1} > \text{doc3}$

4. (a) $W_{q,t} \times W_{d,t}$ is stored in the Accumulator, which can provide more weight to document where the query terms appear many time and more weight for terms appear in little specific document. With accumulator, we can provide ranking more efficiently since we reduce the number of document we need to consider each step. The document with a higher result generated by its accumulator is thought to have a higher rank.

(b) cost time and space, not efficient.

(c) The "Limitting" approach

① create an empty set of Accumulator and set a limit L

② For each query term t, order by decreasing $W_{q,t}$
fetch the inverted list of tfor each $(d, f_{t,i})$ if $|A| < L$, generate A_d if d has no accumulator

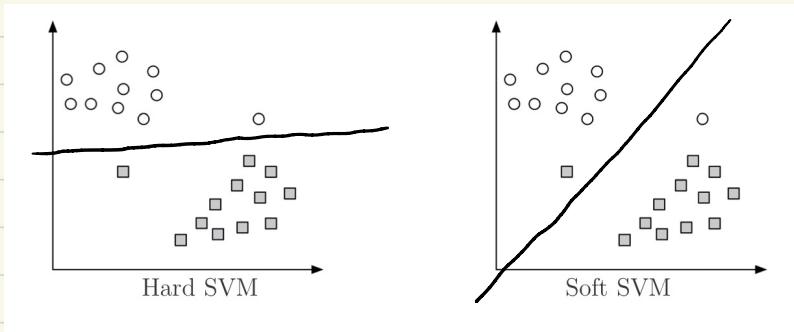
$$A \leftarrow A + W_{q,t} \times W_{d,t}$$

$$A_d \leftarrow A_d / W_d$$

③ return the r documents with the greatest A_d (order by decreasing A_d)

5

(a)



(b) Yes, because soft SVM relax the notion of linear separability. A small number of points are allowed to be in the wrong side of the line; Soft Margin SVMs are less susceptible to noise because they introduce a slack variable that allows misclassification.

(c) ① linear separable

② Only contain two classes, A binary classification problem

(d) main parameters in the "primal":

$$\frac{b}{\|w\|} \text{ and } w$$

parameter in "dual":

$$f(x) = \sum \alpha_i y_i x_i^T x + b$$

Each non-zero α_i indicates that corresponding x_i is a support vector

b ① consider 1 element: (support $\geq \frac{1}{3}$)

null \rightarrow morning

evening
coffee
tea
pastry

$\emptyset \rightarrow$ hot

② consider 2 elements: (Support $\geq \frac{1}{3}$)

(morning, pastry) (coffee, pastry),

③ consider 3 elements: (support $\geq \frac{1}{3}$)
none

④ confidence $\geq \frac{3}{4}$:

(morning \rightarrow pastry)

(coffee \rightarrow pastry)

7.0

	P ₁	P ₂	P ₃	P ₄ , P ₅
P ₁	0	1	5	10
P ₂	1	0	3.5	8
P ₃	5	3.5	0	4
P ₄ , P ₅	10	8	4	0

②

	P ₁ , P ₂	P ₃	P ₄ , P ₅
P ₁ , P ₂	0	5	10
P ₃	5	0	4
P ₄ , P ₅	10	4	0

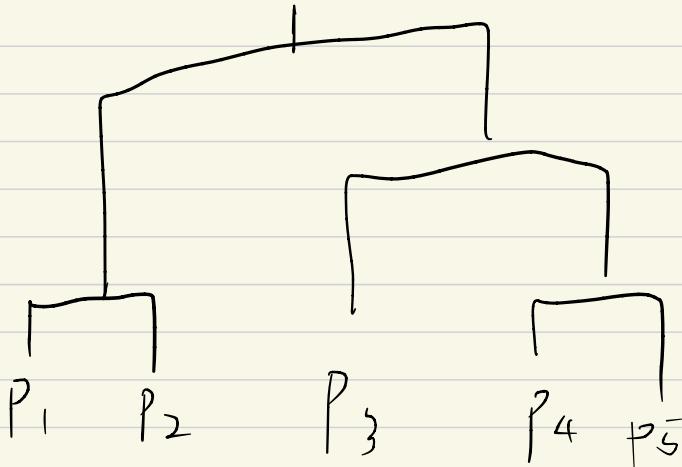
③

	P ₁ , P ₂	P ₃ , P ₄ , P ₅
P ₁ , P ₂	0	10
P ₃ , P ₄ , P ₅	10	0

④

	P ₁ , P ₂ , P ₃ , P ₄ , P ₅
P ₁	0
P ₂	0
P ₃	0
P ₄	0
P ₅	0

Complete link



conditional

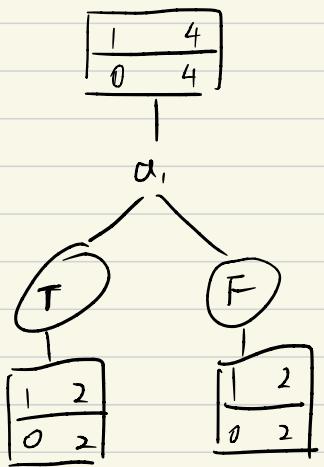
8. The features (attributes) in the data set is independent to each other.

It is necessary because $P(x_1, x_2, \dots, x_n | c) = \prod P(x_i | c)$ only if $(x_i | c)$ is conditional independent to each other.

Because, it is hard to ensure features in the real world are conditional independent to each other, they can be correlative to each other.

$$Q1. H(\text{LABEL}) = -\left(\frac{1}{2} \log_2 \frac{1}{2}\right) \times 2 = 1$$

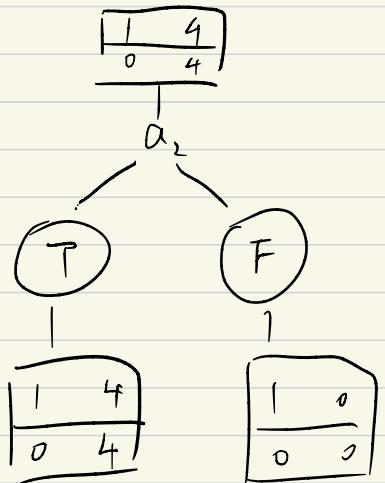
a_1 :



$$-\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right) \cdot \frac{1}{8} \cdot 2 = 1$$

$$IG(a_1 | R) = 1 - 1 = 0$$

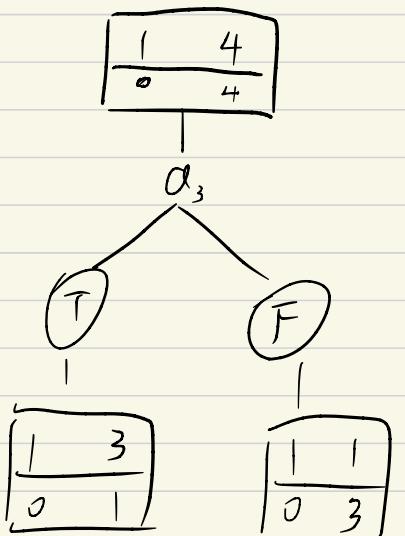
a_2 :



$$-\left(\frac{1}{2} \log_2 \frac{1}{2}\right) \times 2 = 1$$

$$IG(a_2 | \text{LABEL}) = 1 - 1 = 0$$

a_3



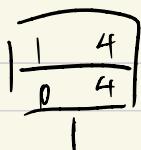
$$-\frac{1}{2} \left(\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right) + \left[-\frac{1}{2} \left(\frac{3}{7} \log_2 \frac{3}{7} + \frac{1}{7} \log_2 \frac{1}{7} \right) \right]$$

$$= 0.81$$

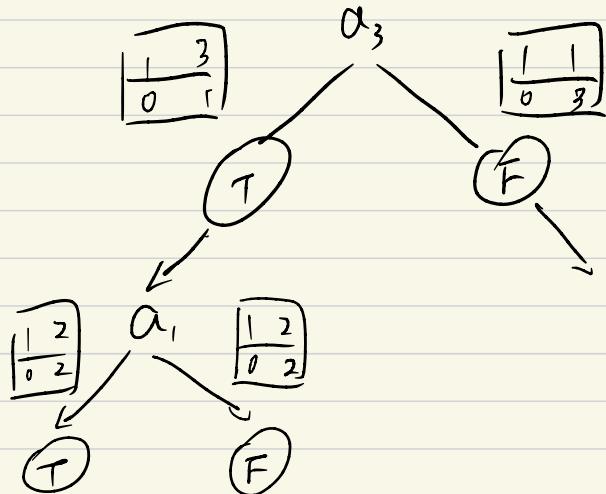
$$IG(a_3 | \text{LABEL}) = 1 - 0.81 = 0.19$$

So a_3 is the root.

(2)



when a_1, a_2, a_3 can not give any information
 $a_1 a_2 a_3 = TTF$, Label can be 1 or 0.



- 10 ① initialize a weight
② Forward pass, calculate the hidden H_i and output.
③ Calculate the error between o_k and t_n , update the weight in a specific learn rate, a small learning weight would lead to a local minimal, Large learning weight lead to over shooting, until network has converged.

11. (a) SVM is appropriate, SVM can handle a huge attribute set in an efficiently. efficiently to find the weight of attribut-

(b) K - Nearest Neighbour is inappropriate for many attributes which would make similarity become useless, it is not right to assume each attribute to have the same weight, calculating similarity in a high dimension is computational heavy.

(1) not appropriate, two many attributes, it would need a lot of time to make the net converged

(d) Not appropriate, Too many attributes, probabilities based on numeric values will be essentially random. It is over sensitive to redundant and irrelevant attributes.

1b) Yes, there must be some book highly correlative to subject, by using this book to predict, for example, in DE tree, through those books with high IG can provide good performance

$\rho_{He} - R$

We may use the keyword or tag of it to represent it, but sometime user may search it based on their content.

Goal:

① Relevance

② Novelty

③ serendipitous ,

④ increase recommendation diversity .