

MLFD 2022 Paper 3

7a) i) System A

		Actual	
		T	F
S	T	5	1
	F	9	119,985
			119,985

System B

		Actual	
		T	F
S	T	12	20
	F	2	119,966
			119,966

$$FN = 11$$

$$TN = 23995$$

$$FP = 21$$

$$TP = 17$$

$$\text{i)} A_A = \frac{119,985 + 5}{120,000} = 0.9999$$

$$A_B = \frac{119,966 + 12}{120,000} \approx 0.9998$$

$$\text{iii) Precision} = \frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

$$P_A = \frac{5}{5+1} \approx 0.83$$

$$P_B = \frac{12}{12+20} \approx 0.38$$

$$R_A = \frac{5}{5+9} \approx 0.36$$

$$R_B = \frac{12}{12+2} \approx 0.86$$

5) i) Null hypothesis: System A is no better than System B.

N is the total number of classified documents. ($120,000$)

P is the probability of system B beating system A ($P=0.5$).

Count number of times A beats B, B beats A, and they predict same.
Use the formula to find if probability is less than some α (usually 0.05 or 0.01).
If it is then reject null hypothesis.

$$\text{i)} A \text{ beats } B \text{ (plus)} = 1 + 20 = 21$$

$$B \text{ beats } A \text{ (minus)} = 8 + 0 = 8$$

$$\text{Praw (null)} = 120,000 - 29 = 119,971$$

ii) No since the difference in accuracy is very small and the sign test shows A doesn't beat B that much nor than B beats A.

iv) Comparing precision and recall:

$$P_A = 0.83 > P_B = 0.38$$

$$R_A = 0.36 < R_B = 0.86$$

so A has much higher precision but B has much higher recall.

iv) We can calculate the F-measure for each system as well:

$$F = \frac{2 \cdot P \cdot R}{P + R}, F_A \approx 0.50, F_B \approx 0.53$$

So B slightly better, but not much. The bigger difference depends if we want higher precision or recall specifically.

We can't do a sign test since we can't see for each document which system has higher precision or recall. These can only be calculated over all documents.

(a) Transition probs:

$$[L \rightarrow L: 0, L \rightarrow M: 1, L \rightarrow H: 0, \\ M \rightarrow L: 0, M \rightarrow M: \frac{1}{2}, M \rightarrow H: \frac{1}{2}, \\ H \rightarrow L: 0, H \rightarrow M: \frac{1}{3}, H \rightarrow H: \frac{2}{3}]$$

Emission probs

$$[(L, +): 1, (L, ++): 0, (L, +++) : 0, \\ (M, +): 0, (M, ++): 1, (M, +++) : 0, \\ (H, +): 0, (H, ++): \frac{2}{3}, (H, +++) : \frac{1}{3}]$$

b) 1. For a first order HMM the probability of a state given all previous states is equal to the probability of the state given the previous one state.

2. The observed feature depends only on the hidden state.

Assumption 1 isn't that appropriate as the infection number could depend on multiple previous states. For example if it's been H for 3 or 4 timesteps we might expect it to lower since most people will have developed immunity.

Assumption 2 is appropriate as it's likely the positivity rate depends on the number of currently infected people.

$$c) \quad \delta_L(t) = \max_{1 \leq i \leq 3} \left[a_{ij} \cdot \delta_i(t-1) \cdot b_j(0) \right]$$

$$L=1, M=2, N=3$$

$$t=8 \quad b_8(x)=$$

$$\delta_L(8) = \max \left\{ \begin{array}{l} 0 \\ \frac{1}{2} \times \frac{1}{2} \times 0 \\ \frac{1}{2} \times \frac{1}{3} \times 0 \end{array} \right\} = 0$$

$$\delta_M(8) = \frac{1}{2} \times 1 \times 0 = \frac{1}{2} \cdot 0$$

$$\delta_H(8) = \frac{1}{2} \times 1 \times \frac{1}{3} = \frac{1}{3}$$

so what is

$$\delta_L(a) = \max \left[\begin{array}{l} 0 \\ \frac{1}{2} \times \frac{1}{2} \times 0 \\ \frac{1}{2} \times \frac{1}{3} \times 0 \end{array} \right] = 0$$

$$\delta_M(a) = \max \left[\begin{array}{l} 1 \times 0 \times 1 \\ \frac{1}{2} \times \frac{1}{2} \times 1 \\ \frac{1}{3} \times \frac{1}{3} \times 1 \end{array} \right] = \frac{1}{4} \cdot \frac{1}{9}$$

$$\delta_H(a) = \max \left[\begin{array}{l} 0 \\ \frac{1}{2} \times \frac{1}{2} > \frac{2}{3} \\ \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} \end{array} \right] = \frac{1}{6} \cdot \frac{4}{27}$$

c) $\Sigma_{\text{L}}(10) = 0$

$$\Sigma_{\text{M}}(10) = \max \left[\begin{array}{l} 1 \times 0 \times 1 \\ \frac{1}{2} \times \frac{1}{9} \times 1 \\ \frac{2}{3} \times \frac{4}{27} \times 1 \end{array} \right] = \frac{8}{81} \approx \frac{1}{18}$$

$$\Sigma_{\text{H}}(10) = \max \left[\begin{array}{l} 0 \\ \frac{1}{2} \times \frac{1}{9} \times 1 \\ \frac{2}{3} \times \frac{4}{27} \times 1 \end{array} \right] = \frac{1}{9} \approx \frac{8}{81}$$

So At 10 q 8
H H H

- d)
1. The first markov assumption from (b) is not realistic for modelling infections
 2. It doesn't take into account lots of other important factors such as changes in behaviour, social distancing rule changes etc.

$$q) P_{NB}(c) = \operatorname{argmax}_{c \in \text{class}} [\log P(c) + \sum_{w \in c} \log P(w; |c|)]$$

$$P(c) = \frac{\text{count}(c)}{\text{count}(\text{all documents})}$$

$$P(w; |c|) = \frac{\text{count}(w; |c|)}{\text{count}(\text{all words in } c)}$$

- b) 4500 used for data training 90 for each brand.
 500 for testing, 10 for each brand.

This allows a high amount for training while still leaving enough for testing.

c) i) Attacked : Terrible, avoid, not

Safe : Awesome, Fantastic, Support

Terrible and, avoid, awesome, fantastic should generalise well as they are strong words with one main use.

Support and not probably won't generalise well as they aren't strong sentiment indicators and can be used for good or bad e.g. 'it's not great' or 'it's not bad at all'.

c) ii) Extract words following not as a group press with not to better get the context and sentiment.
(Maybe 3?)

Tag parts of speech as verbs, adjectives etc. to prevent confusion if a word has two meanings.

d) 1. Retrain often with new data to allow it to know any new words that are good indicators

2. Try to use post in context of surrounding posts to get better understanding.

3. let the classifier annotate new data to train itself.

e) Adv

- More fine-grained analysis for each post
- Might be quicker to spot individual posts that are starting campaigns against boards
- It's more flexible - can be adapted for different needs.

e) Disadv

- More computational complexity
- Isolating posts may lose contextual data
- Lose brand-level insights such as overall public perception