

COMPUTER SCIENCE TRIPOS Part IA – 2026 – Paper 3

9 Machine Learning and Real-world Data (sht25)

You want to classify sentence pairs of sentences $S1$ and $S2$ as to whether their meaning is identical or not. You have collected some sentence pairs and labels assigned by humans, as below.

Sentence pair $\langle S1, S2 \rangle$		Class
1) Mini-poodle feeds man delicious acorns.	Man feeds dog tasty acorns.	DIFFERENT
2) Elsa didn't serve chicken today.	Chicken wasn't cooked by Elsa.	SAME
3) Dog was bitten by squirrel.	Squirrel attacked tiny dog violently.	SAME
4) This school happily admits 6-year-olds.	This school admits no 6-year-olds.	DIFFERENT

- (a) As features, you want to use the tokens contained in the sentence pair. Give the formula for the Naive Bayes Classifier, defining all parameters in the formula, and explain how this formula can be applied to the situation above. [3 marks]
- (b) What is the main disadvantage of using tokens as features here? [3 marks]
- (c) In an alternative approach to feature extraction, we derive features for each sentence $S1$ and $S2$, expressing how long it is (F1, F2), how many negation tokens such as *not*, *none*, *nobody*... it contains (F3, F4), and whether it contains passive voice (F6, F7). For each sentence pair $\langle S1, S2 \rangle$ we report the number of shared tokens between $S1$ and $S2$ (F7), and whether the sentences have the same main verb (F8) and the same subject (F9), as shown in the table below.

	F1 S1 length	F2 S2 length	F3 S1 neg.	F4 S2 neg.	F5 S1 passive	F6 S2 passive	F7 overlap	F8 same verb	F9 same subj
a	5	5	1	1	N	Y	2	N	N
b	5	5	0	1	N	N	4	Y	Y
c	5	5	0	0	N	N	3	Y	N
X	?	?	1	2	N	N	5	N	Y

- (i) Feature sets a,b,c were derived from 3 of the sentence pairs 1)–4) above, but the correct assignment got lost. Assign each feature set to its sentence pair. [2 marks]
- (ii) For the sentence pair from the table above that remains non-aligned, provide the values for features F1–F9. [2 marks]
- (iii) Write a sentence pair that corresponds to feature set X. You can freely decide how long the sentences are. The class should be SAME. [3 marks]
- (d) State the central assumption behind the NB classifier, and whether you think it holds in the scenario in Part (c). If so, explain why. If not, provide a counter example. [3 marks]
- (e) Give two ways how the classifier in Part (c) might be improved. [4 marks]