Aragonda Gayathri

**Q1) 1) Test Document : "Predictable no-fun"**

Vocabulary Size $(v) = 20$

**Negative class**

Total tokens $= 14 \rightarrow$ Denominator $= 14 + 20 = 34$

Prior : $P(-) = 3/5 = 0.6$

Word likelihoods using Add-1 smothing

$P(\text{Predictable} | -) = (2+1)/34 = 3/34 \approx 0.0882$

$P(\text{no} | -) = (0+1)/34 = 1/34 \approx 0.0294$

$P(\text{fun} | -) = (0+1)/34 = 1/34 \approx 0.0294$

Score for Negative class $= 0.6 \times 0.0882 \times 0.0294 \times 0.0294$

$= 0.0000458$

**Positive class :**

Total tokens $= 12 \rightarrow$ Denominator $= 12 + 20 = 32$

Prior $P(+) = 2/5 = 0.4$

All words are unseen in the positive class

$P(w | +) = (0+1)/32 = 1/32 = 0.03125$

Score for positive class $= 0.4 \times 0.03125 \times 0.03125 \times$

$0.03125$

$= 0.0000122$

2) Compare the scores
- $P(-/d) \approx 0.000458$
- $P(+/d) \approx 0.0000122$

Since $P(-/d) > P(+/d)$, the system should assign the document to the : Negative class

Q2 1) → Representational harm occurs when AI systems or classifiers reinforce societal stereotypes or treat certain groups unfairly in their representation.
→ Kititchenko + Mohammad (2018) studied sentiment analysis models and found that these systems assigned more negative sentiment scores to names commonly associated with African American individuals - Even when the names where not in negative context
→ This is representational harm because the classifier was biased against names linked to a specific racial group, reflecting and amplifying social biases
2) → A key risk is the over censorship of marginalized communities
→ Dixon et ali (2018) found that phrases like "I am gay" or "I am muslim" were incorrectly flaged as toxic - not because of the meaning, but due to biases in the data.
→ This leads to silencing of valid identity expressions and can prevent marginalized voices from being heard online

3) → Although African American English (AAE) and Indian English legitimate varieties of English they differ in grammar, vocabulary and usage
→ classifiers often performs worse on these dialects because they are trained mostly on Standard American English, which dominates the trainning datasets. This lack of dialectal diversity results in poor accuracy and unfair outcomes for speakers of underrepresented English varieties

Q3 1)

Sentence $S_1$: $<s>$ I love NLP $</s>$
$P(I/<s>) = 2/3 = 0.6667$
$P(love|I) = 2/2 = 1.0$
$P(NLP|love) = 1/2 = 0.5$
$P(</s>/NLP) = 1/1 = 1.0$

Sentence $S_2$: $<s>$ I love deep learning $</s>$
$P(I/<s>) = 2/3 = 0.6667$
$P(love|I) = 2/2 = 1.0$
$P(deep|love) = 1/2 = 0.5$
$P(learning|deep) = 2/2 = 1.0$
$P(<s>/learning = 1/2 = 0.5$
$P(S_2) = 0.6667 \times 1.0 \times 0.5 \times 1.0 \times 0.5 = 0.1667$

Sentence 1 is more probable under the diagram model

2) P(noodle late) using MLE:

c(ate, noodle) = 0

Total after 'ate' = 12

P(noodle late) = 0/12 = 0.0

Assiging zero probability to an unseen biagram makes the entire sentence probability zero.

This causes issues when computing sentence probability or perprexity, leading to "unreliable" or answerable models

Vocabulary size $V = 10$

Total count after 'ate' = 12

P(noodle late) = $(0+1)/(12+10) = 1/22 = 0.0455$

Q4) 1) P(cats | I, like) = c(I like cats) / c(I like)

Since we don't have trigram count for I like cats, we back off to

bigram

P(cats | like) = c(like cats) / c(like)

c(like cats) = 2

c(like) = c(like cats) + c(like dogs) = 2 + 1 = 3

p(cats | I, like) = 2/3 = 0.6667

2) Try trigram first

P(dogs | you, like) = c(you like dogs) / c(you like)

But $c$(you like dogs) $= 0$ (unseen trigram)
So we back off to bigram
$p$(dogs like) $= c$(like dogs / $c$(like)
$c$(like dogs) $= 1$
$c$(like) $= 3$
$p$(dogs/ you, like) $= 1/3 = 0.3333$

3) → Back off is necessary because some trigrams are not observed in the training corpus. If we only use trigram probability, such cases would lead to a probability of zero for entire sentence.
→ Back off helps mitigate this by falling back to lower-order n-grams (bigrams or unigrams), enabling more to bust probability estimates and preventing zero probability issues

Q5) i) formulas
precision $= Tp/ (Tp + FP)$
Recall $= Tp/ (Tp + FN)$

class : cat
$Tp = 5$, $FP = 15$, $FN = 15$
precision $= 5/ (5 + 15) = 0.25$
Recall $= 5/ (5 + 15) = 0.25$

Class : Dog
$Tp = 20$, $FP = 25$, $FN = 25$

precision $= 20/(20 + 25) \approx 0.4444$
Recall $= 20(20 + 25) \approx 0.4444$

class: Rabbit
Tp $= 10$, Fp $= 15$, FN $= 15$
precision $= 10/(10 + 15) = 0.4$
Recall $= 10/(10 + 15) = 0.4$

2) 1) From per-class metrics
- precision values: Cat $= 0.25$, Dog $= 0.4444$, Rabit $= 0.4$

- Recall values: Cat $= 0.25$, Dog $= 0.4444$, Rabbit $= 0.4$

Macro Precision $= (0.25 + 0.4444 + 0.4)/3 = 1.0944/3 =$
$= 0.3648$

Macro recall $= (0.25 + 0.4444 + 0.4)/3 = 0.3648$

2) Total correct predictions (True positives):
$$5 + 20 + 10 = 35$$
Total predictions $= 90 \rightarrow$ Incorrect predictions
$$90 - 35 = 55$$
Micro precision $= 35/(35 + 55) = 35/90 = 0.3889$
Micro Recall $= 35/(35 + 55) = 35/90 = 0.3889$

3) Macro Averaging
- Averages precision and recall acrass all classes Equally

- usefull when each class should be treated with Equal importance especially in imbalanced datasets

Micro Averaging
- Computes Precision and recall by aggregating total TP, FP and FN across all classes.
- useful for overall performance measurment and reflects the influence of more frequent classes

Conclusion:
- use macro averaging to Evaluate per-class fairness
- use micro averaging to Evaluate overall system effectiveness.