

CS 585 Spring 2024 Written Assignment #03

Due: **Saturday, February 24, 2024 at 11:59 PM CST**

Points: **30**

Objectives:

1. (20 points) Demonstrate your understanding of the Naive Bayes classifier.
2. (10 points) Demonstrate your understanding of POS Tagging with Hidden Markov Model and Viterbi algorithm.

Problem 1 [20 pts]

Your task is to **manually** develop a Naive Bayes **SPAM** classifier / model. Training and test sets are provided below. Assume that no additional pre-processing is necessary:

Training set:	
Sample:	Label:
get cheap pills online	SPAM
online pharmacy pills fast	SPAM
online pills fast shipment	SPAM
my pills prescription	HAM
shipment not fast	HAM
order pills online	SPAM
fast pills delivery	SPAM
this online pharmacy	HAM
take your pills	HAM
get your pills cheap	SPAM
my online order	HAM

Test set:	
Sample:	Label:
order pills fast	SPAM
no prescription online pharmacy	SPAM
this order was online	HAM

A) [3 pt] Extract vocabulary from data (you can add / remove rows as necessary):

Vocabulary:			
Word	count	Word	count
get	2	Not	1
cheap	2	delivery	1
pills	8	this	1
online	6	take	1
pharmacy	2	your	2
Fast	4	order	2
shipment	2		
my	2		
prescription	1		

Total unique words in vocabulary : 15

Total number of words = 37

Total words of the class SPAM = 22

Total words of the class HAM = 15

B) [10 pt] Construct the model(s) (include add-1 smoothing) / derive all parameters. Include all formulas and detailed derivation.

Model / parameters:
<p>Prior Probabilities</p> <p>$P(y=SPAM) = 6/11$ $P(y=HAM) = 5/11$</p> <p><u>Likelihood Probabilities of words for class SPAM:</u> $P(x_i \text{class}) = \text{count}(x_i, SPAM) / \text{sum}(\text{Count}(x, y=SPAM))$</p> <p>Using the above likelihood equation :</p> <p> $P(x=\text{get} y=SPAM) = \text{count}(\text{get}, SPAM) = 2/22$ $P(x=\text{cheap} y=SPAM) = \text{count}(\text{cheap}, SPAM) = 2/22$ $P(x=\text{pills} y=SPAM) = \text{count}(\text{pills}, SPAM) = 6/22$ $P(x=\text{online} y=SPAM) = \text{count}(\text{online}, SPAM) = 4/22$ $P(x=\text{pharmacy} y=SPAM) = \text{count}(\text{pharmacy}, SPAM) = 1/22$ $P(x=\text{Fast} y=SPAM) = \text{count}(\text{Fast}, SPAM) = 3/22$ $P(x=\text{shipment} y=SPAM) = \text{count}(\text{shipment}, SPAM) = 1/22$ $P(x=\text{my} y=SPAM) = \text{count}(\text{my}, SPAM) = 0/22$ $P(x=\text{prescription} y=SPAM) = \text{count}(\text{prescription}, SPAM) = 0/22$ $P(x=\text{Not} y=SPAM) = \text{count}(\text{Not}, SPAM) = 0/22$ $P(x=\text{delivery} y=SPAM) = \text{count}(\text{delivery}, SPAM) = 1/22$ </p>

$P(x=\text{this} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 0/22$
 $P(x=\text{take} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 0/22$
 $P(x=\text{your} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/22$
 $P(x=\text{order} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/22$

Adding Laplace smoothing:
 $\text{Count}(w) + 1/22 + 15 = 1/37$

$P(x=\text{get} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 3/37$
 $P(x=\text{cheap} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 3/37$
 $P(x=\text{pills} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 7/37$
 $P(x=\text{online} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 5/37$
 $P(x=\text{pharmacy} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/37$
 $P(x=\text{Fast} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 4/37$
 $P(x=\text{shipment} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/37$
 $P(x=\text{my} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/37$
 $P(x=\text{prescription} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/37$
 $P(x=\text{Not} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/37$
 $P(x=\text{delivery} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/37$
 $P(x=\text{this} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/37$
 $P(x=\text{take} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/37$
 $P(x=\text{your} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/37$
 $P(x=\text{order} | y=\text{SPAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/37$

Likelihood Probabilities of words for class HAM:

$P(x=\text{get} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 0/15$
 $P(x=\text{cheap} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 0/15$
 $P(x=\text{pills} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/15$
 $P(x=\text{online} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/15$
 $P(x=\text{pharmacy} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 1/15$
 $P(x=\text{Fast} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 1/15$
 $P(x=\text{shipment} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 1/15$
 $P(x=\text{my} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/15$
 $P(x=\text{prescription} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/15$
 $P(x=\text{Not} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 1/15$
 $P(x=\text{delivery} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 0/15$
 $P(x=\text{this} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/15$
 $P(x=\text{take} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/15$
 $P(x=\text{your} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/15$
 $P(x=\text{order} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/15$

Adding Laplace smoothing: $\text{count}(\text{word})+1/15+15 = (c+1)/30$

$P(x=\text{get} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 1/30$

$P(x=\text{cheap} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 1/30$

$P(x=\text{pills} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 3/30$

$P(x=\text{online} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 3/30$

$P(x=\text{pharmacy} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/30$

$P(x=\text{Fast} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/30$

$P(x=\text{shipment} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/30$

$P(x=\text{my} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 3/30$

$P(x=\text{prescription} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/30$

$P(x=\text{Not} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/30$

$P(x=\text{delivery} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 1/30$

$P(x=\text{this} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/30$

$P(x=\text{take} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/30$

$P(x=\text{your} | y=\text{HAM}) = \text{count}(\text{get}, \text{SPAM}) = 2/30$

$P(x=\text{order} | y=\text{HAM}) = \text{count}(\text{get}, \text{HAM}) = 2/30$

C) [5 pt] Test your model (include all derivations | ignore “unknown” words if necessary):

Testing:

Test1(t1) : order pills fast

$P(y=\text{SPAM}) * P(x=\text{order} | \text{SPAM}) * P(\text{pills} | \text{SPAM}) * P(\text{fast} | \text{SPAM})$
 $= 6/11 * 2/37 * 7/37 * 4/37 = 336/557183 = 0.00060303 = 6.0303\text{E-}4$

$P(y=\text{HAM}) * P(x=\text{order} | \text{HAM}) * P(\text{pills} | \text{HAM}) * P(\text{fast} | \text{HAM})$
 $= 5/11 * 2/30 * 3/30 * 2/30 = 60/297000 = 0.00020202 = 2.0202\text{E-}4$

From above $P(Y=\text{SPAM} | t1) > P(Y=\text{HAM} | t1)$.

So **true positive(TP)**

Test2(t2) : no prescription online pharmacy

$P(y=\text{SPAM}) * P(\text{no} | \text{SPAM}) * P(\text{prescription} | \text{SPAM}) * P(\text{online} | \text{SPAM}) * P(\text{pharmacy} | \text{SPAM})$

$= 6/11 * 1/37 * 5/37 * 2/37 * 3/30 = 180/16715490 = 0.00001077$
 $= 1.077\text{E-}5$

$$\begin{aligned}
& P(y=\text{HAM}) * P(\text{no}|\text{HAM}) * P(\text{prescription}|\text{HAM}) * p(\text{online} | \text{HAM}) * p(\text{pharmacy}|\text{HAM}) \\
& = 5/11 * 2/30 * 3/30 * 2/30 \\
& = 60/297000 = 0.00020202 = 2.0202\text{E-}4
\end{aligned}$$

$P(y=\text{HAM}) > P(\text{SPAM})$
 So **False Negative (FP)**

Test3(t3): this order was online

$$P(y=\text{SPAM}) * P(\text{this}|\text{SPAM}) * P(\text{order}|\text{SPAM}) * p(\text{was} | \text{SPAM}) * p(\text{online} | \text{SPAM})$$

$$6/11 * 1/37 * 2/37 * 5/37 = 60/557183 = 0.00010768 = 1.0768\text{E-}5$$

$$P(y=\text{HAM}) * P(\text{this}|\text{HAM}) * P(\text{order}|\text{HAM}) * p(\text{was} | \text{HAM}) * p(\text{online} | \text{HAM})$$

$$= 5/11 * 2/30 * 2/30 * 3/30 = 60 / 297000 = 0.00020202 = 2.0202\text{E-}4$$

$$P(y=\text{HAM}) > P(y=\text{SPAM})$$

True Negative(TN)

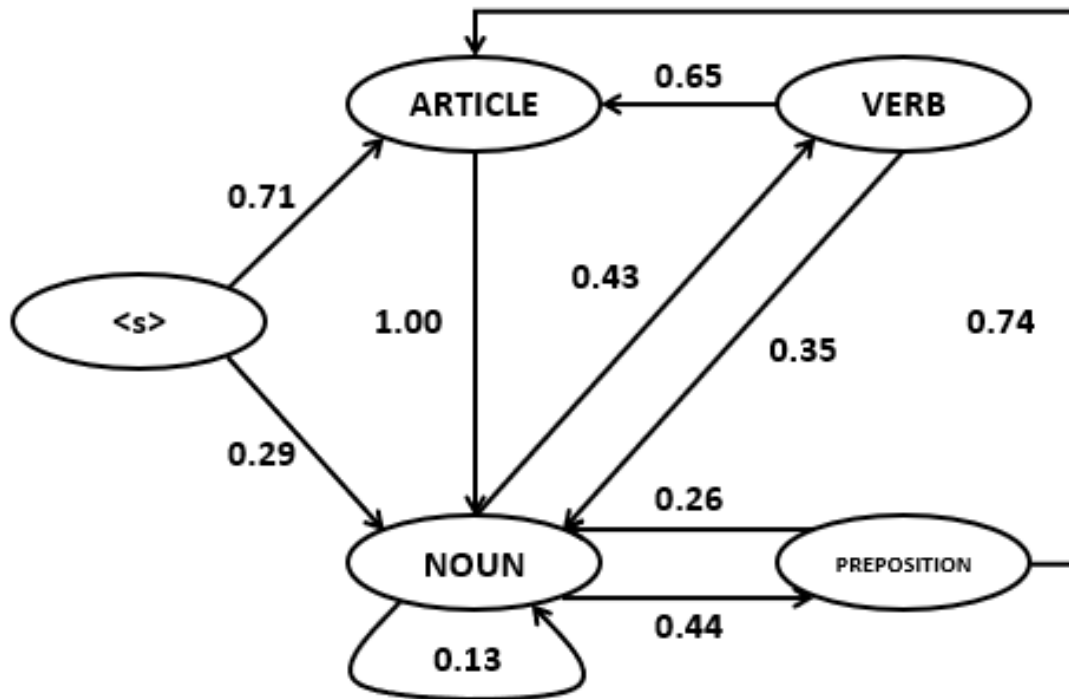
D) [2 pts] Evaluate your model (create confusion matrix, calculate accuracy, sensitivity (recall), precision, specificity, negative predictive value, accuracy and F-Score):

evaluation:			
	Positive SPAM	Negative HAM	
Positive: SPAM	TP = 1(t1)	FN = 1(t2)	Sensitivity TP/TP+FN =1/2 = 0.5
Negative: HAM	0	TN = 1(t3)	Specificity =TN/TN+FP =1/1 + 0 = 1.0
	Precision TP / (TP + FP) = 1/1+0 = 1.0	Negative Predictive = TN/TN+FN = 1/1+1 = 0.5	Accuracy: 1+1/1+1+1 =2/3

F-Score: 0.5

Problem 2 [10 pts]

Given the following Hidden Markov model (transition probabilities shown; emission probabilities to be determined by you using corpus C data) based on corpus C:



And the following table of selected word counts from some corpus C:

Word/Tag	N	V	ART	P	TOTAL
<i>flies</i>	21	23	0	0	44
<i>fruit</i>	49	5	1	0	55
<i>like</i>	10	30	0	21	61
<i>a</i>	1	0	201	0	202
<i>the</i>	1	0	300	2	303
<i>flower</i>	53	15	0	0	68
<i>flowers</i>	42	16	0	0	58
<i>birds</i>	64	1	0	0	65
others	592	210	56	284	1142

TOTAL	833	300	558	307	1998
--------------	-----	-----	-----	-----	------

Use Viterbi algorithm to tag (with Part of Speech tags) the following sentence:

birds like the fruit

Calculating Transition Probability Matrix:

	<s>	ARTICLE	NOUN	VERB	PREPOSITION
<s>	0	0.71	0.29	0	0
ARTICLE	0	0	1.00	0	0
NOUN	0	0	0.13	0.43	0.44
VERB	0	0.65	0.35	0	0
PREPOSITION	0	0.74	0.26	0	0

Computing Observation Matrix:

	birds	like	the	fruit
<s>	0.000	0.000	0.000	0.000
ARTICLE	0.000	0.000	0.537	0.001
NOUN	0.076	0.012	0.001	0.058
VERB	0.003	0.100	0.000	0.016
PREPOSITION	0.000	0.068	0.006	0.000

Computing the Viterbi values for First Column:

$V(1,) = \text{Viterbi}[s, \text{observation}]$:

$$V(1,1) = \text{viterbi}[\text{ARTICLE}, \text{birds}] = P(\text{ARTICLE} \mid \langle s \rangle) * P(\text{birds} \mid \text{Article}) = 0.71 * 0 = 0.000$$

$$V(1,2) = \text{viterbi}[\text{NOUN}, \text{birds}] = P(\text{NOUN} \mid \langle s \rangle) * P(\text{birds} \mid \text{NOUN}) = 0.29 * 0.076 = 0.0220$$

$$V(1,3) = \text{viterbi}[\text{VERB}, \text{birds}] = P(\text{VERB} \mid \langle s \rangle) * P(\text{birds} \mid \text{VERB}) = 0.0 * 0.003 = 0.000$$

$$V(1,4) = \text{viterbi}[\text{Preposition}, \text{birds}] = P(\text{Preposition} \mid \langle s \rangle) * P(\text{birds} \mid \text{Preposition}) = 0 * 0 = 0.000$$

PREPOSITION	V1,4 = 0.000	V2,4	V3,4	V4,4
VERB	V1,3 = 0.000	V2,3	V3,3	V4,3
NOUN	V1,2 = 0.0220	V2,2	V3,2	V4,2
ARTICLE	V1,1 = 0.000	V2,1	V3,1	V4,1
<S>	birds	likes	the	fruit

Computing the states in the second column :

$$V2(1) = \text{viterbi}[\text{ARTICLE}, \text{like}] = \max_s'(\text{viterbi}[\text{state } s', \text{like}] * a_{s',1} * b_1(\text{like})) = \max($$

$$V1(1) * P(\text{ARTICLE} \mid \text{ARTICLE}) * P(\text{like} \mid \text{ARTICLE}) = 0.000 * 0.00 * 0.000 = 0.000$$

$$V1(2) * P(\text{ARTICLE} \mid \text{NOUN}) * P(\text{like} \mid \text{ARTICLE}) = 0.0220 * 0.29 * 0.000 = 0.000$$

$$V1(3) * P(\text{ARTICLE} \mid \text{VERB}) * P(\text{like} \mid \text{ARTICLE}) = 0.000 * 0.00 * 0.000 = 0.000$$

$$V1(4) * P(\text{ARTICLE} \mid \text{PREPOSITION}) * P(\text{like} \mid \text{ARTICLE}) = 0.000 * 0.74 * 0.000 = 0.000$$

$$= 0.000$$

$$V2(2) = \text{viterbi}[\text{NOUN}, \text{like}] = \max_s'(\text{viterbi}[\text{state } s', \text{like}] * a_{s',2} * b_2(\text{like})) = \max($$

$$V1 (1) * P(NOUN | ARTICLE) * P(like | NOUN) = 0.000 * 1.00 * 0.012 = 0.000$$

$$V1 (2) * P(NOUN | NOUN) * P(like | NOUN) = 0.0220 * 0.13 * 0.012 = 0.0001131$$

$$V1 (3) * P(NOUN | VERB) * P(like | NOUN) = 0.000 * 0.35 * 0.012 = 0.000$$

$$V1 (4) * P(NOUN | PREPOSITION) * P(like | NOUN) = 0.000 * 0.26 * 0.012 = 0.000$$

$$V2 (2) = 0.00003432.$$

$$V2 (3) = \text{viterbi}[\text{VERB, like}] = \max_s'(\text{viterbi}[\text{state } s', \text{like}] * a_{s',3} * b_3(\text{like})) = \max($$

$$V1 (1) * P(\text{VERB} | \text{ARTICLE}) * P(\text{like} | \text{VERB}) = 0.000 * 0.00 * 0.100 = 0.000$$

$$V1 (2) * P(\text{VERB} | \text{NOUN}) * P(\text{like} | \text{VERB}) = 0.02204 * 0.43 * 0.100 = 0.0094772$$

$$V1 (3) * P(\text{VERB} | \text{VERB}) * P(\text{like} | \text{VERB}) = 0.000 * 0.00 * 0.100 = 0.000$$

$$V1 (4) * P(\text{VERB} | \text{PREPOSITION}) * P(\text{like} | \text{VERB}) = 0.000 * 0.00 * 0.100 = 0.000$$

$$V2 (3) = 0.0094772$$

3rd column:

$$V3 (1) = \text{viterbi}[\text{ARTICLE, the}] = \max_s'(\text{viterbi}[\text{state } s', \text{the}] * a_{s',1} * b_1(\text{the})) = \max($$

$$V2 (1) * P(\text{ARTICLE} | \text{ARTICLE}) * P(\text{the} | \text{ARTICLE}) = 0.000 * 0.00 * 0.360 = 0.000$$

$$V2 (2) * P(\text{ARTICLE} | \text{NOUN}) * P(\text{the} | \text{ARTICLE}) = 0.0001131 * 0.00 * 0.360 = 0.000$$

$$V2 (3) * P(\text{ARTICLE} | \text{VERB}) * P(\text{the} | \text{ARTICLE}) = 0.0094772 * 0.65 * 0.537 = 0.00330802$$

$$V2 (4) * P(\text{ARTICLE} | \text{PREPOSITION}) * P(\text{the} | \text{ARTICLE}) = 0.000 * 0.74 * 0.360 = 0.000$$

$$V3 (1) = 0.00330802$$

$$V3 (2) = \text{viterbi}[\text{NOUN, the}] = \max_s'(\text{viterbi}[\text{state } s', \text{the}] * a_{s',2} * b_2(\text{the})) = \max($$

$$V2 (1) * P(\text{NOUN} | \text{ARTICLE}) * P(a | \text{NOUN}) = 0.000 * 1.00 * 0.001 = 0.000$$

$$V2 (2) * P(\text{NOUN} | \text{NOUN}) * P(\text{the} | \text{NOUN}) = 0.00003432 * 0.10 * 0.001 = 3.432E-9$$

$$V2 (3) * P(\text{NOUN} | \text{VERB}) * P(\text{the} | \text{NOUN}) = 0.0094772 * 0.35 * 0.001 = 3.32E-7$$

$$V2 (4) * P(\text{NOUN} | \text{PREPOSITION}) * P(\text{the} | \text{NOUN}) = 0.000 * 0.26 * 0.001 = 0.000$$

$$V3 (2) = 3.32E-7$$

4th column:

$V4(1) = \text{viterbi}[\text{ARTICLE}, \text{fruit}] = \max_s'(\text{viterbi}[\text{state } s', \text{fruit}] * a_{s',1} * b_1(\text{fruit})) = \max($

$V3(1) * P(\text{ARTICLE} | \text{ARTICLE}) * P(\text{fruit} | \text{ARTICLE}) = 0.00330802 * 0.00 * 0.000 = 0.000$

$V3(2) * P(\text{ARTICLE} | \text{NOUN}) * P(\text{fruit} | \text{ARTICLE}) = 3.32E-7 * 0.00 * 0.000 = 0.000$

$V3(3) * P(\text{ARTICLE} | \text{VERB}) * P(\text{fruit} | \text{ARTICLE}) = 0.000 * 0.65 * 0.000 = 0.000$

$V3(4) * P(\text{ARTICLE} | \text{PREPOSITION}) * P(\text{fruit} | \text{ARTICLE}) = 0.000 * 0.74 * 0.000 = 0.000$

$V4(1) = 0.000$

$V4(2) = \text{viterbi}[\text{NOUN}, \text{fruit}] = \max_s'(\text{viterbi}[\text{state } s', \text{fruit}] * a_{s',2} * b_2(\text{fruit})) = \max($

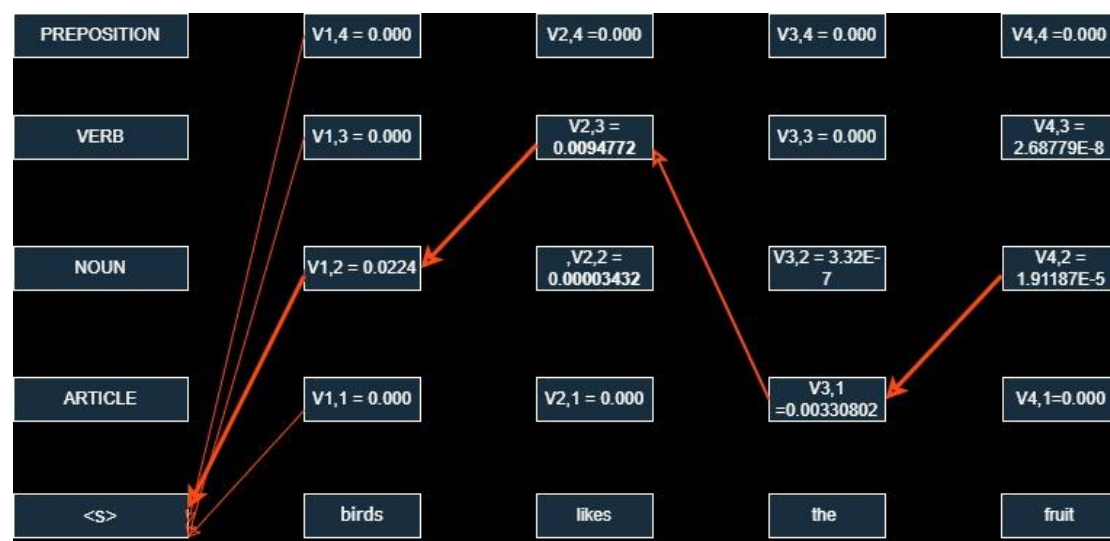
$V3(1) * P(\text{NOUN} | \text{ARTICLE}) * P(\text{fruit} | \text{NOUN}) = 0.00330802 * 1.00 * 0.058 = \mathbf{1.91187E-5}$

$V3(2) * P(\text{NOUN} | \text{NOUN}) * P(\text{fruit} | \text{NOUN}) = 3.32E-7 * 0.10 * 0.058 = 1.98256E-9$

$V3(3) * P(\text{NOUN} | \text{VERB}) * P(\text{fruit} | \text{NOUN}) = 0.000 * 0.35 * 0.058 = 0.000$

$V3(4) * P(\text{NOUN} | \text{PREPOSITION}) * P(\text{fruit} | \text{NOUN}) = 0.000 * 0.26 * 0.058 = 0.000$

$V4(2) = \mathbf{1.91187E-5}$



The most probable tags are : Noun, VERB, ARTICLE, NOUN for the words in the sentence: birds, like, the, fruits