

NLP Midterm Exam

Álvaro Faúndez

November 2021

Question 1

We have a vocabulary V of size 300. How many parameters will you have to estimate for:

1. an MLE unigram language model
2. a bigram language model with add-one smoothing
3. a maximum likelihood trigram language model

Answer

1. an MLE unigram language model:
It needs to estimate $|V| = 300$ parameters.
2. a bigram language model with add-one smoothing
It needs to estimate $|V|^2 = 300^2$ parameters.
3. a maximum likelihood trigram language model
It needs to estimate $|V|^3 = 300^3$ parameters.

Question 2

Consider the following training corpus T of sentences:

1. $\langle s \rangle$ I am Sam $\langle /s \rangle$
2. $\langle s \rangle$ Sam I am $\langle /s \rangle$
3. $\langle s \rangle$ I do not like green eggs and ham $\langle /s \rangle$

Compute the following maximum likelihood parameters:

- $p(I \mid \langle s \rangle) =$
- $p(\langle /s \rangle \mid \text{Sam}) =$
- $p(I \mid \text{do}) =$
- $p(\text{ham} \mid \text{eggs, and}) =$

Answer

$$p(I \mid \langle s \rangle) = \frac{\text{count}(\langle s \rangle \text{ I})}{\text{count}(\langle s \rangle)} = \frac{2}{3} \quad (1)$$

$$p(\langle /s \rangle \mid \text{Sam}) = \frac{\text{count}(\text{Sam } \langle /s \rangle)}{\text{count}(\text{Sam})} = \frac{1}{2} \quad (2)$$

$$p(I \mid \text{do}) = \frac{\text{count}(\text{do I})}{\text{count}(\text{do})} = \frac{0}{1} \quad (3)$$

$$p(\text{ham} \mid \text{eggs, and}) = \frac{\text{count}(\text{eggs and ham})}{\text{count}(\text{eggs and})} = \frac{1}{1} \quad (4)$$

Question 3

We have the following training corpus:

1. the green book STOP
2. my blue book STOP
3. his green house STOP
4. book STOP

Assume we have a trigram language model with linear interpolation based on this corpus, with $\lambda_i = \frac{1}{3}$ for all i . Compute the value of the parameter $p(\text{book} \mid \text{the, green})$ under this model. Assume STOP as part of your unigram model.

Answer

The corpus C has 14 tokens. The probabilities needed are:

$$p_{ML}(\text{book} \mid \text{the green}) = \frac{\text{count}(\text{the green book})}{\text{count}(\text{the green})} = \frac{1}{1} \quad (5)$$

$$p_{ML}(\text{book} \mid \text{green}) = \frac{\text{count}(\text{green book})}{\text{count}(\text{green})} = \frac{1}{2} \quad (6)$$

$$p_{ML}(\text{book}) = \frac{\text{count}(\text{book})}{|C|} = \frac{3}{14} \quad (7)$$

Finally, the interpolation

$$\begin{aligned} p(\text{book} \mid \text{the green}) &= \lambda_3 \times p_{ML}(\text{book} \mid \text{the green}) + \lambda_2 \times p_{ML}(\text{book} \mid \text{green}) + \lambda_1 \times p_{ML}(\text{book}) \\ &= \frac{1}{3} \times \frac{\text{count}(\text{the green book})}{\text{count}(\text{the green})} + \frac{1}{3} \times \frac{\text{count}(\text{green book})}{\text{count}(\text{green})} + \frac{1}{3} \times \frac{\text{count}(\text{book})}{|C|} \\ &= \frac{1}{3} \left(\frac{1}{1} + \frac{1}{2} + \frac{3}{14} \right) \\ &= 0.57142857142 \end{aligned} \quad (8)$$

Question 4

Given a vocabulary V of size 1000 and a tag set T of size 30, suppose you wish to train a bigram HMM tagger. How many transition and emission parameters will your model have? Explain.

Answer

The emissions probabilities in a bigram HMM $e(w | t)$ are needed for all $w \in V$ and each $t \in T$. That means that there are $|V| \times |T| = 1000 \times 30 = 30000$ possible combinations of words and tags.

The transitions probabilities in a bigram HMM $p(t_i | t_k)$ are needed for each possible combination of two transitions. That means there are $|T| \times |T| = 30 \times 30 = 900$ possible combinations of transitions pairs.

In the WSJ corpus example we reviewed in class, the transitions matrix also included the start token $\langle s \rangle$. That would make the transitions matrix $31 \times 30 = 930$. If we also include the stop token $\langle /s \rangle$, the size would be $31 \times 31 = 961$. The emissions would stay the same.

Question 5

Let $V = \{Karlsson, lives, happily\}$ Let $T = \{N, V\}$

We have a trigram HMM tagger that has the following non-zero parameters (all other parameters are zero):

Transition parameters:

- $p(N \mid START) = 0.5$
- $p(V \mid START) = 0.5$
- $p(N \mid START, N) = 1.0$
- $p(N \mid START, V) = 0.9$
- $p(V \mid START, V) = 0.1$
- $p(STOP \mid N, N) = 1.0$
- $p(STOP \mid V, V) = 1.0$
- $p(STOP \mid V, N) = 1.0$
- $p(s \mid u, v) = 0$ for all other p parameters

Emission parameters:

- $e(Karlsson \mid N) = 0.8$
- $e(happily \mid N) = 0.2$
- $e(lives \mid V) = 0.7$
- $e(happily \mid V) = 0.3$

Under this model, how many pairs of sequences $x_1 \dots x_n, y_1 \dots y_n$ have non-zero probability $p(x_1, \dots x_n, y_1 \dots y_{n+1}) \geq 0$? Show all of these possible sequences.

Answer

Given the emissions:

$e(w \mid t)$	Karlsson	lives	happily
N	0.8	0	0.8
V	0	0.7	0.7

We can deduce that the possible N are **Karlsson** and **happily** and the possible V are **lives** and **happily**.

And from the transitions, we can deduce:

- then sentences can start either with an N or a V
 - If the sentence starts with an N, the only possibility for the next word is another N
 - If the sentence starts with a V, the only possibility for the next word is another N or V
- Since there are no transitions involving a combination of only N or V, the STOP must come after two words

That means that the only sequences with non-zero probabilities are:

- $(t_1, t_2, t_3) = (N, N, STOP)$
 - $(x_1, x_2, x_3) = (\text{Karlsson}, \text{Karlsson}, STOP)$
 - $(x_1, x_2, x_3) = (\text{Karlsson}, \text{happily}, STOP)$
 - $(x_1, x_2, x_3) = (\text{happily}, \text{happily}, STOP)$
 - $(x_1, x_2, x_3) = (\text{happily}, \text{Karlsson}, STOP)$
- $(t_1, t_2, t_3) = (V, N, STOP)$
 - $(x_1, x_2, x_3) = (\text{lives}, \text{Karlsson}, STOP)$
 - $(x_1, x_2, x_3) = (\text{lives}, \text{happily}, STOP)$
 - $(x_1, x_2, x_3) = (\text{happily}, \text{Karlsson}, STOP)$
 - $(x_1, x_2, x_3) = (\text{happily}, \text{happily}, STOP)$
- $(t_1, t_2, t_3) = (V, V, STOP)$
 - $(x_1, x_2, x_3) = (\text{lives}, \text{lives}, STOP)$
 - $(x_1, x_2, x_3) = (\text{lives}, \text{happily}, STOP)$
 - $(x_1, x_2, x_3) = (\text{happily}, \text{lives}, STOP)$
 - $(x_1, x_2, x_3) = (\text{happily}, \text{happily}, STOP)$

That makes a total of 12 possible sequences of words and tags with non-zero probabilities.

Question 6

Find one tagging error in each of the following sentences that are tagged with the Penn Treebank tagset:

1. I/PRP need/VBP a/DT flight/NN from/IN Atlanta/NN
2. Does/VBZ this/DT flight/NN serve/V B dinner/NNS
3. I/PRP have/VB a/DT friend/NN living/VBG in/IN Denver/NNP
4. Can/VBP you/PRP list/VB the/DT nonstop/JJ afternoon/NN flights/NNS

Answer

1. Atlanta should be NNP (proper noun, singular)
2. Dinner should be an NN (noun, singular or mass)
3. have should be a VBP (Verb, non-3rd ps. sing. present)
4. afternoon should be a JJ (adjective)

Question 7

We have a vocabulary $V = \{\text{Hello}\}$ and a constant $N \geq 1$. For any $x_1 \dots x_n$ such that $x_i \in V$ for $i = \dots (n-1)$ and $x_n = \text{STOP}$, we define

$$p(x_1, \dots, x_n) = \begin{cases} \frac{1}{N} & \text{if } n \leq N \\ 0 & \text{otherwise} \end{cases}$$

Is this a valid language model? Explain.

Answer

In order to be valid, given a corpus C a language model must:

1. $p(s) \geq 0$ for each sentence $s \in C$
2. $\sum_{s \in C} p(s) = 1$

Condition 1 : it is satisfied for each sentence because the two options are zero or a positive rational.

Condition 2 : it's not fulfilled. The probability function means that any sentence shorter or equal than N has probability $\frac{1}{N}$ and everything else zero. We can see that choosing $N = 3$. There are only two sentences possible shorter or equal than 3:

1. $x_1, x_2 = \text{Hello}, \text{STOP}$
2. $x_1, x_2, x_3 = \text{Hello}, \text{Hello}, \text{STOP}$

Those sentences have both a probability of $\frac{1}{3}$ and any other sentence has probability zero because they are longer than $N = 3$. This way, the sum of the probabilities is $\frac{2}{3}$ and not 1 as required to be a valid language model.

Question 8

Now we have a vocabulary $V = \{\text{Hello}, \text{Goodbye}\}$ and a constant $N \geq 1$. For any $x_1 \dots x_n$ such that $x_i \in V$ for $i = \dots (n-1)$ and $x_n = \text{STOP}$, we define

$$p(x_1, \dots, x_n) = \begin{cases} \frac{1}{2} & \text{if } n = 2 \\ 0 & \text{otherwise} \end{cases}$$

Is this a valid language model? Explain.

Answer

The difference with question 8 is in the probability function. Now, the probability function means that every sentence of size 2 has probability $\frac{1}{2}$ and everything else zero. We can see that choosing any N . there are always only two sentences of size 2:

1. $x_1, x_2 = \text{Hello}, \text{STOP}$
2. $x_1, x_2 = \text{Goodbye}, \text{STOP}$

Now, the sum of all the probabilities is 1, so it is a valid language model.

Question 3

Consider the task of classifying the word bass using the Naïve Bayes algorithm. The features used are bag-of-word features. Assume the following likelihoods for each word being part of a “fish” and “music” class, and equal prior probabilities for each class.

What class will Naïve Bayes assign to the sentence “I eat fresh bass after music lesson”? Show your work.

	fish	music
I	0.09	0.16
eat	0.29	0.06
fresh	0.10	0.05
after	0.07	0.06
music	0.04	0.15
lesson	0.08	0.11

Answer

Classes $C = \{C_{\text{fish}}, C_{\text{music}}\}$.

Bag of words $\{\text{I}, \text{eat}, \text{fresh}, \text{after}, \text{music}, \text{lesson}\}$

Features $\vec{x} = [1, 1, 1, 1, 1, 1]$

Priors:

$$\begin{aligned} P_{\text{prior}}(C_{\text{fish}}) &= \frac{1}{2} \\ P_{\text{prior}}(C_{\text{music}}) &= \frac{1}{2} \end{aligned} \tag{9}$$

$$\begin{aligned} P(C_{\text{fish}} | \vec{x}) &\propto P(C_{\text{fish}}) \times P(\vec{x} | C_{\text{fish}}) \\ &\propto P(\text{I} | C_{\text{fish}}) \times \\ &\quad P(\text{eat} | C_{\text{fish}}) \times P(\text{fresh} | C_{\text{fish}}) \times P(\text{after} | C_{\text{fish}}) \times \\ &\quad P(\text{music} | C_{\text{fish}}) \times P(\text{lesson} | C_{\text{fish}}) \\ &\propto 0.5 \times 0.09 \times 0.29 \times 0.10 \times 0.07 \times 0.04 \times 0.08 \\ &\propto 2.9232e - 7 \end{aligned} \tag{10}$$

$$\begin{aligned}
P(C_{\text{music}} | \vec{x}) &\propto P(C_{\text{music}}) \times P(\vec{x} | C_{\text{music}}) \\
&\propto P(\text{I} | C_{\text{music}}) \times \\
&\quad P(\text{eat} | C_{\text{music}}) \times P(\text{fresh} | C_{\text{music}}) \times P(\text{after} | C_{\text{music}}) \times \\
&\quad P(\text{music} | C_{\text{music}}) \times P(\text{lesson} | C_{\text{music}}) \\
&\propto 0.5 \times 0.16 \times 0.06 \times 0.05 \times 0.06 \times 0.15 \times 0.11 \\
&\propto 2.376e-7
\end{aligned} \tag{11}$$

Since $P(C_{\text{fish}} | \vec{x}) > P(C_{\text{music}} | \vec{x})$, the class assigned is C_{fish} .