

Unit 4 Unsupervised Learning (2

Course > weeks)

> Homework 5 >

2. Maximum Likelihood Estimation

Audit Access Expires May 11, 2020

You lose all access to this course, including your progress, on May 11, 2020.

2. Maximum Likelihood Estimation

Consider a general multinomial distribution with parameters θ . Recall that the likelihood of a dataset $\mathcal D$ is given by:

$$P\left(\mathcal{D}; heta
ight)=\prod_{i=1}^{| heta|} heta_{i}^{c_{i}}$$

where c_i is the occurrence count of the i-th event.

The MLE of θ is the setting of θ that maximizes $P\left(\mathcal{D};\theta\right)$. In lecture we derived this to be

$$heta_i^* = rac{c_i}{\sum_{j=1}^{| heta^*|} c_j}$$

Unigram Model

4/4 points (graded)

Consider the sequence:

ABABBCABAABCAC

A unigram model considers just one character at a time and calculates $p\left(w\right)$ for $w\in\{A,B,C\}.$

What is the MLE estimate of heta? Give your result to three decimal places.

 θ_A^* 6/14 **Answer:** 0.4285714286

 θ_B^* 5/14 \checkmark Answer: 0.3571428571

Using the MLE estimate of θ on \mathcal{D} , which of the following sequences is most likely?

ABC

BBB

ABB

AAC



Solution:

We calculate the MLE as $rac{\mathrm{count}(w)}{N}$ where N=14 and the counts are 6, 5, and 3.

For comparing probabilities in part two, we simply multiply. We only need to compare the numerators: $6\times5\times3$, 5^3 , 6×5^2 , and $6^2\times3$.

Submit

You have used 1 of 3 attempts

1 Answers are displayed within the problem

Bigram Model 1

1.0/1 point (graded)

A bigram model computes the probability $p\left(\mathcal{D};\theta\right)$ as:

$$p\left(\mathcal{D}; heta
ight) = \prod_{w_1,w_2 \in \mathcal{D}} p\left(w_2|w_1
ight)$$

where w_2 is a word that follows w_1 in the corpus.

This is also a multinomial model. Assume the vocab size is N. How many parameters are there?

Grading note: The formula above contains an error: the probability $p\left(\mathcal{D};\theta\right)$ in a bigram model is generally:

$$p\left(\mathcal{D}; heta
ight)=p\left(w_{0}
ight)\prod_{w_{1},w_{2}\in\mathcal{D}}p\left(w_{2}|w_{1}
ight)$$

where w_0 is the first word, and (w_1,w_2) is a pair of consecutive words in the document. In this case, the number of parameters is $(N-1)+(N^2-N)=N^2-1$. However, with the model as written above, there are only parameters N^2-N .

The grader is now fixed to accept both as correct and regrading is happening.

N*(N-1)

✓ Answer: N^2 - 1

STANDARD NOTATION

Solution:

Recall the likelihood of D in bigram model is (though this is not what written):

$$p\left(\mathcal{D}; heta
ight)=p\left(w_{0}
ight)\prod_{w_{1},w_{2}\in\mathcal{D}}p\left(w_{2}|w_{1}
ight)$$

where w_0 is the first word, and (w_1,w_2) is a pair of consecutive words in the document.

Denote the set of all N words by V. The set of parameters is

$$\{p\left(w_{0}
ight):w_{0}\in V\}\;\cup\;\{p\left(w_{1}|w_{2}
ight):w_{1}\in V,w_{2}\in V\}$$

and the only constraints on these parameters are

$$egin{array}{lll} \sum_{w_{0}\in V}p\left(w_{0}
ight)&=&1\ & \sum_{w_{1}\in V}p\left(w_{1}|w_{2}
ight)&=&1 & ext{for all }w_{2}\in V. \end{array}$$

Hence, the number of parameters is $(N-1)+(N^2-N)=N^2-1$. (Note that this is also the number of parameters $p\left(w_1,w_2\right)$ where $w_1\in V,w_2\in V$, which determine the joint distribution.

Solution to the problem as written:

The likelihood of ${\cal D}$ in bigram model was given as

$$p\left(\mathcal{D}; heta
ight) = \prod_{w_1,w_2 \in \mathcal{D}} p\left(w_2|w_1
ight)$$

without taking into account the likelihood $p\left(w_{0}
ight)$ of the first word. In this case, the parameters are

$$\left\{ p\left(w_{1}|w_{2}
ight):w_{1}\in V,w_{2}\in V
ight\}$$

where $\sum_{w_1 \in V} p\left(w_1|w_2
ight) = 1$ for all $w_2 \in V$. Hence, the number of parameters is $N^2 - N$.

Submit

You have used 1 of 3 attempts

1 Answers are displayed within the problem

Bigram Model 2

0/1 point (graded)

Which of the following represents the MLE for the **conditional probability** $p\left(w_2 \mid w_1\right)$?

$$\bigcap \frac{\operatorname{count}(w_1, w_2)}{\sum_{w_1', w_2' \in \mathcal{D}} \operatorname{count}(w_1', w_2')}$$

$$\bigcirc \frac{\operatorname{count}(w_1, w_2)}{\sum_{w_1, w_2' \in \mathcal{D}} \operatorname{count}(w_1, w_2')} \checkmark$$

$$\frac{\sum_{w_1',w_2 \in \mathcal{D}} \operatorname{count}(w_1',w_2)}{\sum_{w_1,w_2' \in \mathcal{D}} \operatorname{count}(w_1,w_2')}$$



Solution:

This is a simple application of Bayes Rule:

$$p\left(w_{2}|w_{1}
ight)=rac{p\left(w_{1},w_{2}
ight)}{p\left(w_{1}
ight)}$$

To compute $p(w_1)$, we marginalize out w_2 .

Submit

You have used 3 of 3 attempts

1 Answers are displayed within the problem

Bigram Model 3

1 point possible (graded)

Consider the same sequence from the unigram model:

ABABBCABAABCAC

If you estimate θ on this, what probability will be assigned to the following test sequence? Assume the starting probabilities of all characters $p\left(w|\text{null}\right)$ is uniform. Give your answer to three decimal places.

AABCBAB

Answer: 0

Solution:

There is no need to compute the actual probability. Since the transition $C \to B$ does not appear in $\mathcal D$, the probability assigned to this new sequence will be 0. This is why techniques like smoothing are important in practice for small datasets.

Submit

You have used 0 of 3 attempts

1 Answers are displayed within the problem

Discussion

Hide Discussion

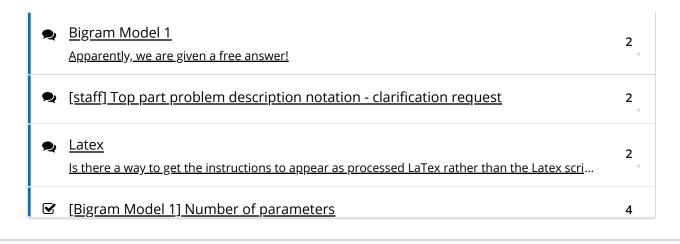
Topic: Unit 4 Unsupervised Learning (2 weeks) :Homework 5 / 2. Maximum Likelihood Estimation

Add a Post

Show all posts by	recent activity
bigram model 3 yes look and there is the answer	2
 ☑ Bigram model 2 simplification and model 3? ≜ Community TA 	4
? <u>Bigram Model 3 - unknown bigram</u>	7
Definition: count(w1,w2) Where is count(w1,w2) defined? Its meaning is not clear.	4
■ Bigram 3 Nice.	5 new_ 11
Create a dict from two numpy arrays to store the thetas Can be useful for storing the thetas of the bigram model: a = np.array(['AA','AB']) b = np.array	1 zeros
? How to approach the last question - Bigram model 3 Hi Should we calculate the number of AA sequences, then AB, then BC and so on in the company to the sequences.	6 origin
Pigram Model 2: the first 3 options mean the same (to me) To me, the first 3 options mean the same since you're just changing the name of the index	3 new_ exes
Pigram Model 1 Apparently from the answer we are not considering words w1 and w2 can be the same.	1 Why

2. Maximum Likelihood Estimation | Homework 5 ...

https://courses.edx.org/courses/course-v1:MITx+...



© All Rights Reserved