

**Natural Language Processing**  
**Prof. Pawan Goyal**  
**Department of Computer Science and Engineering**  
**Indian Institute of Technology, Kharagpur**

**Lecture - 25**  
**Inside-Outside Probabilities**

Welcome back for the final lecture of this week. So, in the last lecture, we had started with the concept of inside outside probabilities and how do you use them for answering certain questions like what is the probability of a sentence as per my grammar and what is the most likely parse. So, I have use the inside algorithm specifically to find out what is the probabilities of this sentence as per my grammar and then in the end I was saying that we will also use this concept to find out the rule probabilities of my grammar and how do we do that again exactly what we will discuss in this lecture.

(Refer Slide Time: 01:00)

*How to get the rule probabilities*

*Parsed Training Data*  
You can count!

$$\hat{P}(N^j \rightarrow \delta) = \frac{C(N^j \rightarrow \delta)}{\sum_{\gamma} C(N^j \rightarrow \gamma)}$$

*But what if the training data is not available?*  
i.e. gold standard parse is not known.

- Underlying CFG is known and we are given a set of sentences
- For each sentence, we can find out all the possible parses
- Maximize the likelihood of the sentences in the data under the PCFG constraints

Pawan Goyal (IIT Kharagpur)      Syntax      Week 5: Lec 25

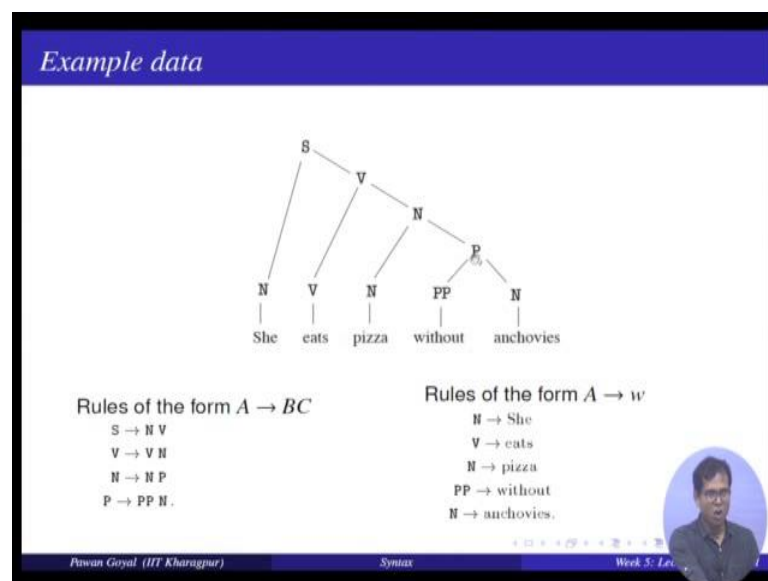
Now in general, how do you obtain the rule probabilities and remember this is very similar to what we were also talking about in the case of HMMs, when I have to learn the parameters of my HMM, I can do that in 2 manner, one where I am given already the labeled data set; that means, I am given what are the sentences and what is their past reach. If I am given all that so, then from there I can compute how many times a particular non terminal derives a particular sequence divide by number of times a particular non terminal has been used that will give me the probability of this particular

rule. If there is only one rule possible from  $N_j$  this is always you want, but if there multiple rules are possible I can find out what fraction of times this particular instance has been used in my labeled data sets and by that I can compute all the rule probabilities.

Now, so this is easy, but what about the case where the training data is not given to us that is I have no way to find out in which sentence what is the rule that has been used. So, what is in fact, given to us? So, we are; assume that we are given the underlying context free grammar. So, I am given all the possible rules, but what I am not given what is the probabilities for each individual rules so; that means, the particular PCGFs part the rule probabilities is not given to us and I am given a lot of sentences and I am the given the grammar that will generate these sentences, but not the rule probabilities and my task is how do I find out these rules probabilities.

So, in other words how do I find out the parameters of my PCGFs and we will use the same sort of idea that we did earlier that is you are given the observation of sentences find out the parameters of a model that maximize the likelihood of observing the sentences. So, this is what we are going to do maximize the likelihood of the sentences in the data under the PCGF constraints and for that we will use some sort of expectation expression algorithm.

(Refer Slide Time: 03:25)



Let us take a simple example and what is the intuition for using this. So, we start with this simple sentence like she eats pizza without anchovies and a particular parse; tree is

given to us. So, what kind of rules do you see here? So, you have rules of the form a single non terminal derives to 2 different non terminals like S derives N and V, V derives V and N and so on and there are rules of the form a non terminal derives a terminal like here and derives anchovies PP derives without N derives pizza V derives eats and so on. So, these are different sort of non terminals that are given the rules that are given to me I do not know the probabilities for each individual rule.

Now, can you think of any other parse for this sentence, she eats pizza without anchovies, so one parse that we saw was she eats pizza and pizza without anchovies modifying pizza what is the other possibilities. So, other could be the without anchovies does not modify pizza it is a separate phrase altogether starting from V, so something like she eats pizza with fog.

(Refer Slide Time: 05:00)

*Example data*

Is any other parse possible for *She eats pizza without anchovies* syntactically?  
Consider *She eats pizza without hesitation*

*New Context-free rules:*

$V \rightarrow VNP$   
 $H \rightarrow hesitation$

Pawan Goyal (IIT Kharagpur)      Syntax      Week 5: Lec.

In that sort of meaning so yeah so, she eats pizza without hesitation that is another sort of possibility, now a single V here is giving me V and N P. So, this is the different parse tree for the same sentence

Now, I can what are the new rules that I have added V gives me V N and P and N gives me hesitation these are 2 new rules I have added. So, now, what do you see I have 2 different sentences; she eats pizza without anchovies and she eats pizza without hesitation and both have 2 possible parse edge and I know what are the all possible rules

now my task is given these sentences how do I find out what should be the ideal rule probabilities.

(Refer Slide Time: 05:49)

*Estimating the model parameters*

We need to find probabilities such as

- $\phi(S \rightarrow N V)$
- $\phi(N \rightarrow \text{pizza})$

*Requirements*

For each non-terminal  $A$ , the derivation probabilities sum up to 1

$$\sum_{\alpha} \phi(A \rightarrow \alpha) = 1$$

For the example grammar:

$$\begin{aligned} \phi(N \rightarrow NP) + \phi(N \rightarrow \text{pizza}) + \phi(N \rightarrow \text{anchovies}) + \\ + \phi(N \rightarrow \text{hesitation}) + \phi(N \rightarrow \text{She}) &= 1 \\ \phi(V \rightarrow VN) + \phi(V \rightarrow VP) + \phi(V \rightarrow \text{cats}) &= 1 \\ \phi(S \rightarrow NV) &= 1 \\ \phi(P \rightarrow PP N) &= 1 \\ \phi(PP \rightarrow \text{without}) &= 1 \end{aligned}$$

Pawan Goyal (IIT Kharagpur) Syntax Week 5: Lecture 5 5/11

Now for that let me first define so we need to find out all these probabilities. So, probability of this rule  $S$  gives me  $N$  and  $V$  probability of rule  $N$  gives me pizza and so on, now what is the PCGF requirement? The PCGF requirement is that it starting from a single non terminal all the possible sets of right hand side that I can generate (Refer Time: 06:11). So, that is so, all the possible rules starting from  $A$  to any  $\alpha$  the probability should also added to 1 so now, if I look at my grammar so, I have 5 rules pair and occurs in the left hand side  $N$  gives me  $NP$ ,  $N$  gives me pizza anchovies and so on. So, all these should also add to 1.

Similarly, 3 rules where  $V$  occurs in the left hand side so, all these should also added to 1 and then there are some other rules like  $S$  gives me  $NV$ ,  $NV$  nothing else, this would be 1 and so on. So, this constraints I can obtain from my from my grammar, I know what are the rules and what is the constraint starting from the left hand side all possible rule should have a probability adding up to 1.

(Refer Slide Time: 06:58)

*Likelihood computation*

$W_1 = \text{"She eats pizza without anchovies"}$

$W_2 = \text{"She eats pizza without hesitation"}$

$$P_\phi(W_1, T_1) = \phi(S \rightarrow N V) \phi(V \rightarrow V N) \phi(N \rightarrow N P) \times$$

$$\times \phi(P \rightarrow PP N) \phi(N \rightarrow \text{She}) \phi(V \rightarrow \text{eats}) \times$$

$$\times \phi(N \rightarrow \text{pizza}) \phi(PP \rightarrow \text{without}) \phi(N \rightarrow \text{anchovies})$$
  

$$P_\phi(W_2, T_1) = \phi(S \rightarrow N V) \phi(V \rightarrow V N P) \phi(P \rightarrow P PP) \times$$

$$\times \phi(N \rightarrow \text{She}) \phi(V \rightarrow \text{eats}) \phi(N \rightarrow \text{pizza}) \times$$

$$\times \phi(PP \rightarrow \text{without}) \phi(N \rightarrow \text{hesitation})$$

Pawan Goyal (IIT Kharagpur) Syntax Week 5: L

Now, I have 2 sentences, can I compute the likelihood of the sentences W 1 and W 2. So, what do I mean by likelihood what is the probability of generating the sentence as of my grammar? Right now I am not giving the rule probabilities, but I can write down in terms of the variable rule probabilities. So, what will the likelihood of W 1? So, for that I have to take the 2 possible parse trees. So, here T 1 is the first parse tree.

So, here I am giving the probabilities of both the sentences as per the first parse tree. So, probability of sentence W 1 as per the first parse tree T 1 and probability of sentence W 2 as per the first parse tree T 1, how do I compute that it is very similar to what we saw in the case of PCGF? How do we compute the probability of a parse tree if the rule probabilities are given that was very easy here the possibilities are not given, but you can parameterize. So, you will say what is the probability of S giving me N V and so on up to you go to the leaves and I do not know these whole probabilities similarly I can write down this likelihood of the sentence W 2 as per my first parse tree.

(Refer Slide Time: 08:11)

*Likelihood computation*

$$P_{\phi}(W_1, T_2) = \phi(S \rightarrow N V) \phi(V \rightarrow V N P) \phi(P \rightarrow P PP) \times$$

$$\times \phi(N \rightarrow She) \phi(V \rightarrow eats) \phi(N \rightarrow pizza) \times$$

$$\times \phi(PP \rightarrow without) \phi(N \rightarrow anchovies)$$
  

$$P_{\phi}(W_2, T_1) = \phi(S \rightarrow N V) \phi(V \rightarrow V N) \phi(N \rightarrow N P) \times$$

$$\times \phi(P \rightarrow PP N) \phi(N \rightarrow She) \phi(V \rightarrow eats) \times$$

$$\times \phi(N \rightarrow pizza) \phi(PP \rightarrow without) \phi(N \rightarrow hesitation)$$
  

*Likelihood of the corpus*

Probability of a sentence  $W$ :  $P_{\phi}(W) = \sum_T P_{\phi}(W, T)$

If the training data comprises of sentences  $W_1, W_2, \dots, W_N$ , then the likelihood is

$$L(\phi) = P_{\phi}(W_1) P_{\phi}(W_2) \dots P_{\phi}(W_N)$$

Pawan Goyal (IIT Kharagpur) Syntax Week 5: Lecture 5 7/11

Similarly, I can do further second parse tree also for both the sentences. So, this tells me. So, if I know all the possible parse trees because my CFGS given to me, I will know all the possible parse tree, I can put my all the rule probabilities as variables and define what is the likelihood of various parse tree now what is the likelihood of the sentence, it will be summation of the likelihood as per different possible parse trees. So, probability sentence such summation over all the possible trees that can generate this  $P \phi W T$ , in this case for both sentences  $W_1 W_2$ , I had 2 different parse trees.

So, I will just add the 2; 2 probabilities to get the likelihood of the sentence and how do I get the likelihood of the whole corpus that has multiple sentences for that I will come to the likelihood of each sentence and multiply those. So, if I have a sentences  $W_1$  to  $W_N$ , I compute the likelihood of each and keep on multiplying now we know that how do I express the likelihood of my corpus in terms of the rule probabilities right the only variable here are all the rule probabilities now I can further define my problem.

(Refer Slide Time: 09:33)

The slide has a blue header with the text "Likelihood maximization". Below the header is a large white area. In the center of this area is a pink rounded rectangle containing the text "Approach" in red, followed by "Starting at some initial parameters  $\phi$ , re-estimate to obtain new parameters  $\phi'$  for which  $L(\phi') \geq L(\phi)$ . Repeat until convergence". At the bottom of the slide is a blue footer bar containing the text "Pawan Goyal (IIT Kharagpur)" on the left, "Syntax" in the center, and "Week 5: Lecture 5 8 / 11" on the right.

My problem would be, so, this is some sort of E M approach, I will start with some initial parameters  $\phi$ ,  $\phi$  means the rule probabilities I want to re estimate so that I obtain some new parameters  $\phi'$  such that the likelihood of my corpus increases now. So,  $L(\phi')$  will be greater than equal to  $L(\phi)$  and I keep on doing that until I converge. So, now, here we have to apply (Refer Time: 10:04) algorithm so that we can keep on updating our rule probabilities, this is the parameter of my system and how do we do that if you remember that like what we did in the case of learning parameters for G S analogous to that what we will do here we will start with some arbitrary rule probabilities  $\phi$  and use that to compute something intermediate.

In this case what we will compute? What is the expected number of times a particular rule has been used if the rule probabilities are as per the current parameters? I will compute the expected value; again use the expected value to compute the probabilities. So, I will compute my  $\phi'$  again use the  $\phi'$  to compute the expected number of times each rule has been used and again compute  $\phi'$  and keep on doing that until you converge and that is why we will be using the inside outside probabilities. So, let us see.

(Refer Slide Time: 11:02)

**Parameter Estimation**

Given some rule probabilities  $\phi$  and training corpus  $W_1, W_2 \dots W_n$ , the new parameters are obtained as:

$$\phi'(A \rightarrow B \ C) = \frac{\text{count}(A \rightarrow B \ C)}{\sum_{\alpha} \text{count}(A \rightarrow \alpha)}$$

$$\phi'(A \rightarrow w) = \frac{\text{count}(A \rightarrow w)}{\sum_{\alpha} \text{count}(A \rightarrow \alpha)}$$

What is  $\text{count}(\cdot)$ ?

$$\text{count}(A \rightarrow B \ C) = \sum_{i=1}^N c_{\phi}(A \rightarrow B \ C, W_i)$$

$$\text{count}(A \rightarrow w) = \sum_{i=1}^N c_{\phi}(A \rightarrow w, W_i)$$

$c_{\phi}(A \rightarrow \alpha, W_i)$  is the expected number of times  $(A \rightarrow \alpha)$  is used in generating the sentence  $W_i$ , when the rule probabilities are given by  $\phi$ .

Pawan Goyal (IIT Kharagpur) Syntax Week 5: Lecture 5 9/11

Idea is that I start with some rules probabilities  $\phi$  and I am given a corpus that that what are sentences that I observing  $W_1, W_2, \dots, W_n$  and I will obtain the new parameters  $\phi'$  using the simple idea. So, this is something that we were saying if we are given the labeled data that is why I will compute the rule probabilities. So, I am saying I can always define probability of the rule  $A \rightarrow B \ C$  as the number of times the rule  $A \rightarrow B \ C$  is used in my corpus divide by all the possible all the different times where  $A$  derives  $\alpha$  for all possible  $\alpha$  and this gives me the probability for a deriving  $B \ C$ .

Similarly, I can compute probability  $A \rightarrow w$  by saying how many times, this rule has been used divide by number of times  $A$  gives me  $\alpha$  has been used in my corpus in my actual corpus, but we do not have any labeled corpus, we only know what parse are possible and for each parse, we can compute the probabilities using the previous parameter  $\phi$ . So, how do I write down this count  $A$  derives  $B \ C$  number of times, this rule has been used for that I use the idea that I have multiple sentences any sentence I can find the expected number of times this rule  $A$  deriving  $B \ C$  has been used. So, that is count  $A$  deriving  $B \ C$  is nothing but summation over all, the sentences number of times  $A$  deriving  $B \ C$  has been used for the particular sentence same 1 for the count of  $A$  deriving  $w$  each sentence find out the expected number of times a particular rule has been used, now how do I actually come up with this formulation expected number of times, a

particular rule has been used in a sentence and for that we use the inside probabilities and outside probabilities.

(Refer Slide Time: 13:17)

**Computing Expected counts**

*Inside probabilities*

The nonterminal  $A$  derives the string of words  $w_i \dots w_j$  in the sentence :  
 $\beta_{ij}(A) = P_\phi(A \Rightarrow^* w_i \dots w_j)$

*Outside probabilities*

Beginning with the start symbol  $S$  we can derive the string  
 $w_1 \dots w_{i-1} A w_{j+1} \dots w_n : \alpha_{ij}(A) = P_\phi(S \Rightarrow^* w_1 \dots w_{i-1} A w_{j+1} \dots w_n)$

*Expected count*

$$c_\phi(A \rightarrow BC, W) = \frac{\phi(A \rightarrow BC)}{P_\phi(W)} \sum_{1 \leq i \leq j \leq k \leq n} \alpha_{ik}(A) \beta_{ij}(B) \beta_{j+1,k}(C)$$

$$c_\phi(A \rightarrow w, W) = \frac{\phi(A \rightarrow w)}{P_\phi(W)} \sum_{1 \leq i \leq n} \alpha_{ii}(A)$$

Pawan Goyal (IIT Kharagpur)
Syntax
Week 5: Lecture 5 10/11

Now coming back to this inside and outside probabilities and how do we use that to compute the expected number of times a rule has been used in a sentence now this must be clear by the previous slide that if I can compute the expected number of times, a rule has been used in a sentence, I can keep on updating my parameters, this is the only bottleneck, in the previous computation and we will see how to do that using the inside and outside probabilities. So, let me give the definitions again. So, inside probability is starting from non terminal  $A$ , I derive the words  $W_i$  to  $W_j$  in the sentence so that is probability that  $A$  derives  $W_i$  to  $W_j$  as per my grammar.

And the outside probabilities starting from the symbol  $S$ , I can derive the string  $W_1$  to  $W_{i-1}$   $A$  and  $W_{j+1}$  to  $W_N$ . So, it starting from  $S$ , I can derive  $W_1$  to  $W_{i-1}$   $A$   $W_{j+1}$  to  $W_N$ , now once we are given the inside and outside probabilities, we can actually compute the expected number of times, the rule has been used and the expression comes out to be this one expected number of times a rule  $A$  has been  $A$  derives  $BC$  has been used in my sentence  $W$  each the rule probability a given  $BC$  divide by the probability of the sentence and this very peculiar term that you see that you are seeing here. So, you are seeing here  $\alpha_{ik}(A) \beta_{ij}(B)$  and  $\beta_{j+1,k}(C)$ , now how do I actually come up with a term like that and how do I come up with this

expression that is expected count is given by this. So, for that let us go back to what we were discussing in the last lecture that I can multiply inside and outside probabilities to know something about the probability for the sentence.

(Refer Slide Time: 15:16)

Handwritten mathematical derivation on a blue background:

Top left:  $\frac{A \rightarrow BC}{A \rightarrow W}$

Top right:  $\alpha_j(p, q) \beta_j(p, q) = P(N^1 \Rightarrow W_{1m}, N^j \Rightarrow W_{pq} | G)$   
 $= \frac{P(N^1 \Rightarrow W_{1m} | G)}{P_\phi(W)} \cdot \frac{P(N^j \Rightarrow W_{pq} | N^1 \Rightarrow W_{1m}, G)}{P_\phi(W)}$

Middle:  $\text{Exp. \# times } (N^j \Rightarrow W_{pq}) \text{ is used} = \sum_{p=1}^m \sum_{q=p}^m \frac{\alpha_j(p, q) \beta_j(p, q)}{P_\phi(W)}$

Bottom left:  $\text{Exp. \# times } (N^j \Rightarrow N^p N^s) \text{ is used} =$

Bottom right:  $\sum_{p=1}^{q-1} P(N^j \Rightarrow N^p N^s) = \sum_{p=1}^{q-1} \beta_r(p, d) \cdot \beta_s(d+1, e)$

Diagram below the bottom left equation shows a tree structure for  $N^j$  deriving  $W_p$  and  $W_q$ , and another tree structure for  $N^j \Rightarrow N^p N^s$  deriving  $W_p$  and  $W_q$ .

Let us go back to that. So, what we were saying if I multiply alpha j p q and beta j p q what does it give me? It gives me the probability that it starting from N 1, I can derive W 1 m and it starting from N j, I can derive W p q as per my grammar. So, now, I can use the chain rule here to write it like that. So, its probability N 1 derives W 1 m. So, this means it derives in any number of steps given by grammar times probability N j derives W p q given N 1 derives W 1 m and my grammar. So, now, what is this probability that N one derives W 1 m given by grammar second write as the P phi W probability of the sentence and what is this say what is the probability that this rule N j has been used to derive W Pp q given that the sentences there and my grammar is there.

And how many times this has been used in this particular context? Only 1 time, so, can I write down expected number of times N j derives W p q is used that would be alpha j p q beta j p q given divided by P phi W and suppose because I do not want to fix this p q, I just want to say expected number of times the rule N j has been used. So, each time it has been used only once for deriving W to W q. So, here I will have to sum over all the possible p and q. So, I will say p can go from one to m suppose there are W 1 to W m.

So, these are number of words and  $q$  can go from  $p$  up to  $m$ . So, this is the expected number of times my rule and  $j$  has been used.

Now, what is something that I have to express I want to find out for example, expected number of times a rule like  $A$  goes to  $B\ C$  has been used or  $a$  goes to  $W$  has been used what is the expected number of times these has been used now for that suppose let us take the easy case expected number of times the rule  $a$  goes to  $W$  or  $a$  derives  $W$  has been used. So, in this case what I am saying  $N_j$  derives a particular terminal here that is some  $p$ . So, I can write the  $\beta_j p$  for that case and  $\beta_j p$  is simply the rule probability that is what is the probability that in non terminal derives this word  $W_p$ , so we will see the expression for that so, this one is easy.

But what about this case when the rule  $A$  derives  $B\ C$  so, in the particular notation that I have written suppose you want to say expected number of times  $N_j$  derives  $N_r\ N_s$  is used. So, now, what would happen? So, this  $\beta_j p\ q$  is when the terminal  $N_j$  derives the whole sequence  $p\ q$ ,  $W_p$  to  $W_q$  and can use any possible rules yes  $N_r\ N_s$  or  $N_j\ N_z$  whatever it can use any non terminals now what I am limiting I am saying this rule should only use this. So, this non terminal should only derive  $N_r\ N_s$ . So, then I am saying so; that means, my  $N_j$  will derive  $N_r\ N_s$  and this  $N_r\ N_s$  will again derive say  $W_p$  to some  $W_d$  and this will derive  $W_d + 1$  to  $W_q$ .

Now how do I modify this equation? So,  $\alpha_j p\ q$  is the outside probability that will remain the same nothing has changed for outside probability, but inside probability because I am saying this should be the situation. So, I further express it like with a particular path. So, I will write in place of  $\beta_j p\ q$ , I will write probability of the rule  $N_j$  derives  $N_r\ N_s$  times this beta probability that is  $\beta_r p\ d$  times  $\beta_s d + 1$  to  $q$ , but now the  $d$  can vary, I have already been given that  $N_j$  gives me  $N_r\ N_s$ , but they can take different possible  $d$ s. So, this will be summation over  $d$  and  $d$  can vary from  $p$  to  $q$  minus 1. So, these between  $p$  and  $q$  now if you put that can you see that you can actually obtain the same expression that was given in the slide.

If you go back to the slide that is what we have been doing here. So, you see the expression we have 3 parameters  $\beta_i j\ b$ . So, 3 parameters  $i\ j$  and  $k$  that corresponds to  $p\ d$  and  $q$  and this was the outside probability and this is the inside probability for the two children;  $i\ j$  and  $j + 1\ k$  and then you have the rule probability here. So, this

expected number of count has been derived in this particular form and same thing you can try with the next formula the expected number of times the rule a deriving the W has been used in my graph and you will obtain this particular expression.

What we are seeing here suppose I start with some initial rule probabilities. So, I can use the inside outside probability. So, all the recursive formulation to compute all the inside outside probabilities for my various rules and stages once I do that I can compute what is the expected number of times each and each individual rule has been used in my corpus as per the current parameters. Once I have the expected number of times rule has been used I can further estimate my parameters by number of times the rule used divide by number of times any particular rule starting with that non terminal has been used and that will give me the new parameters, again I will compute inside outside probabilities expected count, re-estimate the parameters and this I will continue until this converges.

(Refer Slide Time: 23:47)

And how to compute inside-outside probabilities

Inductively, as discussed earlier

$$\beta_{ii}(A) = \phi(A \rightarrow w_i)$$
$$\alpha_{1n}(S) = 1$$

Prasen Goyal (IIT Khartagpur) Syntax Week 5: Lect...

And yeah, so computing inside outside probabilities is as we discussed earlier by this inductive manner. So, what we discussed in this module was that what is parsing in terms of a constituency structure and how do you use the formulation of context free grammar to do parsing, how do we incorporate the rule probabilities there? How do we learn the rule probabilities using this interesting concept of inside outside probabilities? So, I hope by the example that we did in the class, you will be able to understand how it is actual exactly works in practice.

In the next week, we will be starting with this different notion of parsing. So right now, we have done a constituency parsing. So, we will see there is a different notion of parsing called dependency parsing. So, what is the formulation that that dependency parsing follows? How it is different from this constituency parsing and what are different methods we can use for that that will be a topic for the next week.

Thank you.