

Inside-outside probabilities

Pawan Goyal

CSE, IIT Kharagpur

Week 5: Lecture 5

How to get the rule probabilities

Parsed Training Data

You can count!

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

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.

How to get the rule probabilities

Parsed Training Data

You can count!

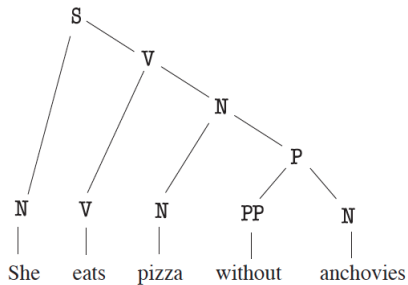
$$\hat{P}(N^j \rightarrow \delta) = \frac{C(N^j \rightarrow \delta)}{\sum_{\gamma} C(N^j \rightarrow \gamma)}$$

But what if the training data is not available?

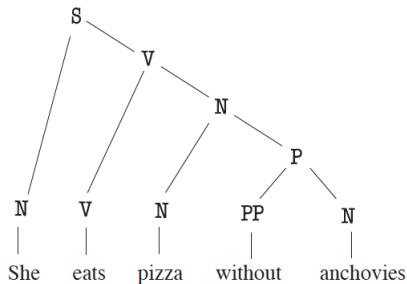
i.e. gold standard parse is not known.

- 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*

Example data



Example data



Rules of the form $A \rightarrow BC$

$S \rightarrow N V$

$V \rightarrow V N$

$N \rightarrow N P$

$P \rightarrow PP N.$

Rules of the form $A \rightarrow w$

$N \rightarrow \text{She}$

$V \rightarrow \text{eats}$

$N \rightarrow \text{pizza}$

$PP \rightarrow \text{without}$

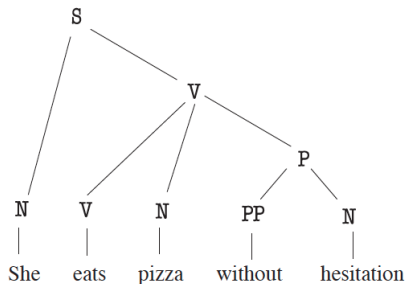
$N \rightarrow \text{anchovies.}$

Example data

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

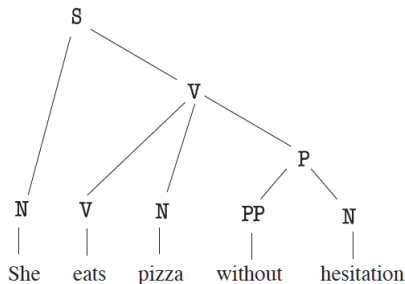
Example data

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



Example data

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



New Context-free rules:

$$V \rightarrow V N P$$
$$N \rightarrow \text{hesitation} .$$

Estimating the model parameters

We need to find probabilities such as

- $\phi(S \rightarrow N \ V)$
- $\phi(N \rightarrow pizza)$

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

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

Estimating the model parameters

We need to find probabilities such as

- $\phi(S \rightarrow N V)$
- $\phi(N \rightarrow pizza)$

Requirements

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

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

For the example grammar:

$$\begin{aligned}\phi(N \rightarrow N P) + \phi(N \rightarrow pizza) + \phi(N \rightarrow anchovies) &+ \\ &+ \phi(N \rightarrow hesitation) + \phi(N \rightarrow She) = 1 \\ \phi(V \rightarrow V N) + \phi(V \rightarrow V N P) + \phi(V \rightarrow eats) &= 1\end{aligned}$$

$$\begin{aligned}\phi(S \rightarrow N V) &= 1 \\ \phi(P \rightarrow PP N) &= 1 \\ \phi(PP \rightarrow without) &= 1\end{aligned}$$

Likelihood computation

W_1 = “She eats pizza without anchovies”

W_2 = “She eats pizza without hesitation”.

Likelihood computation

W_1 = “She eats pizza without anchovies”

W_2 = “She eats pizza without hesitation”.

$$\begin{aligned} P_\phi(W_1, T_1) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N) \phi(N \rightarrow N P) \times \\ &\times \phi(P \rightarrow PP N) \phi(N \rightarrow \text{She}) \phi(V \rightarrow \text{eats}) \times \\ &\times \phi(N \rightarrow \text{pizza}) \phi(PP \rightarrow \text{without}) \phi(N \rightarrow \text{anchovies}) \end{aligned}$$

$$\begin{aligned} P_\phi(W_2, T_1) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N P) \phi(P \rightarrow P PP) \times \\ &\times \phi(N \rightarrow \text{She}) \phi(V \rightarrow \text{eats}) \phi(N \rightarrow \text{pizza}) \times \\ &\times \phi(PP \rightarrow \text{without}) \phi(N \rightarrow \text{hesitation}) \end{aligned}$$

Likelihood computation

$$\begin{aligned}P_{\phi}(W_1, T_2) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N P) \phi(P \rightarrow P PP) \times \\&\times \phi(N \rightarrow \text{She}) \phi(V \rightarrow \text{eats}) \phi(N \rightarrow \text{pizza}) \times \\&\times \phi(PP \rightarrow \text{without}) \phi(N \rightarrow \text{anchovies})\end{aligned}$$

$$\begin{aligned}P_{\phi}(W_2, T_1) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N) \phi(N \rightarrow N P) \times \\&\times \phi(P \rightarrow PP N) \phi(N \rightarrow \text{She}) \phi(V \rightarrow \text{eats}) \times \\&\times \phi(N \rightarrow \text{pizza}) \phi(PP \rightarrow \text{without}) \phi(N \rightarrow \text{hesitation})\end{aligned}$$

Likelihood computation

$$\begin{aligned}P_{\phi}(W_1, T_2) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N P) \phi(P \rightarrow P PP) \times \\&\times \phi(N \rightarrow She) \phi(V \rightarrow eats) \phi(N \rightarrow pizza) \times \\&\times \phi(PP \rightarrow without) \phi(N \rightarrow anchovies)\end{aligned}$$

$$\begin{aligned}P_{\phi}(W_2, T_1) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N) \phi(N \rightarrow N P) \times \\&\times \phi(P \rightarrow PP N) \phi(N \rightarrow She) \phi(V \rightarrow eats) \times \\&\times \phi(N \rightarrow pizza) \phi(PP \rightarrow without) \phi(N \rightarrow hesitation)\end{aligned}$$

Likelihood of the corpus

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

Likelihood computation

$$\begin{aligned}P_{\phi}(W_1, T_2) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N P) \phi(P \rightarrow P PP) \times \\&\times \phi(N \rightarrow She) \phi(V \rightarrow eats) \phi(N \rightarrow pizza) \times \\&\times \phi(PP \rightarrow without) \phi(N \rightarrow anchovies)\end{aligned}$$

$$\begin{aligned}P_{\phi}(W_2, T_1) &= \phi(S \rightarrow N V) \phi(V \rightarrow V N) \phi(N \rightarrow N P) \times \\&\times \phi(P \rightarrow PP N) \phi(N \rightarrow She) \phi(V \rightarrow eats) \times \\&\times \phi(N \rightarrow pizza) \phi(PP \rightarrow without) \phi(N \rightarrow hesitation)\end{aligned}$$

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

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

Likelihood maximization

Approach

Starting at some initial parameters ϕ , re-estimate to obtain new parameters ϕ' for which $L(\phi') \geq L(\phi)$. Repeat until convergence

Parameter Estimation

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

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

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

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

Parameter Estimation

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

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

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

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

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

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

Parameter Estimation

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

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

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

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

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

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

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

Computing Expected counts

Inside probabilities

The nonterminal A derives the string of words $w_i, \dots w_j$ in the sentence :

$$\beta_{ij}(A) = P_{\phi}(A \Rightarrow^* w_i \dots w_j)$$

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

$$w_1 \dots w_{i-1} A w_{j+1} \dots w_n : \alpha_{ij}(A) = P_{\phi}(S \Rightarrow^* w_1 \dots w_{i-1} A w_{j+1} \dots w_n)$$

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

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

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

And how to compute inside-outside probabilities

Inductively, as discussed earlier

$$\beta_{ii}(A) = \phi(A \rightarrow w_i)$$

$$\alpha_{1n}(S) = 1$$