

Lecture L&M 5

Chomsky Normal Form, CYK Parsing

5.1 Noncontracting grammars

For each of the following grammars, construct an equivalent grammar G_L with a nonrecursive start symbol that is “essentially noncontracting”, i.e. only the (new) start symbol may have a λ -rule.

(Extra:) Give a regular expression that describes the language accepted by each grammar.

A. Grammar G :

$$\begin{aligned} S &\rightarrow aS \mid bS \mid B \\ B &\rightarrow bb \mid C \mid \lambda \\ C &\rightarrow cC \mid \lambda \end{aligned}$$

First, we transform the grammar into an equivalent grammar in which the start symbol is no longer recursive.

$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow aS \mid bS \mid B \\ B &\rightarrow bb \mid C \mid \lambda \\ C &\rightarrow cC \mid \lambda \end{aligned}$$

The variables with λ -derivation according to the discussed algorithm:

iteration	<i>null</i>	<i>prev</i>
0	{ <i>B, C</i> }	
1	{ <i>S, B, C</i> }	{ <i>B, C</i> }
2	{ <i>S', S, B, C</i> }	{ <i>S, B, C</i> }
3	{ <i>S', S, B, C</i> }	{ <i>S', S, B, C</i> }

The requested grammar now becomes:

$$\begin{aligned} S' &\rightarrow S \mid \lambda \\ S &\rightarrow aS \mid bS \mid B \mid a \mid b \\ B &\rightarrow bb \mid C \\ C &\rightarrow cC \mid c \end{aligned}$$

The corresponding regular expression is $(a \cup b)^*(bb \cup c^*)$.

B. Grammar G :

$$\begin{aligned} S &\rightarrow ABC \mid aBC \\ A &\rightarrow aA \mid BC \\ B &\rightarrow bB \mid \lambda \\ C &\rightarrow cC \mid \lambda \end{aligned}$$

The grammar already had a non-recursive start symbol.

$$\begin{aligned} S &\rightarrow ABC \mid aBC \\ A &\rightarrow aA \mid BC \\ B &\rightarrow bB \mid \lambda \\ C &\rightarrow cC \mid \lambda \end{aligned}$$

The variables with λ -derivation:

Iteration	<i>null</i>	<i>prev</i>
0	{ <i>B, C</i> }	
1	{ <i>B, C, A</i> }	{ <i>B, C</i> }
2	{ <i>B, C, A, S</i> }	{ <i>B, C, A</i> }
3	{ <i>B, C, A, S</i> }	{ <i>B, C, A, S</i> }

The resulting noncontracting grammar:

$$\begin{aligned} S &\rightarrow ABC \mid AB \mid BC \mid AC \mid A \mid B \mid C \mid \lambda \mid aBC \mid aB \mid aC \mid a \\ A &\rightarrow aA \mid BC \mid B \mid C \mid a \\ B &\rightarrow bB \mid b \\ C &\rightarrow cC \mid c \end{aligned}$$

The corresponding regular expression is $a^*b^*c^*b^*c^*$.

5.2 Chain rules

For each of the following grammars, construct an equivalent grammar G_C without chain rules. (Note that the grammars have no λ -rules, but you may need to make the start symbol nonrecursive.)

C. Grammar G :

$$\begin{aligned} S &\rightarrow AS \mid A \\ A &\rightarrow aA \mid bB \mid C \\ B &\rightarrow bB \mid b \\ C &\rightarrow cC \mid B \end{aligned}$$

Adding non-recursive start symbol:

$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow AS \mid A \\ A &\rightarrow aA \mid bB \mid C \\ B &\rightarrow bB \mid b \\ C &\rightarrow cC \mid B \end{aligned}$$

The chain variables are computed iteratively again:

Iter	$chain(S')$	$chain(S)$	$chain(A)$	$chain(B)$	$chain(C)$
0	$\{S'\}$	$\{S\}$	$\{A\}$	$\{B\}$	$\{C\}$
1	$\{S', S\}$	$\{S, A\}$	$\{A, C\}$	$\{B\}$	$\{C, B\}$
2	$\{S', S, A\}$	$\{S, A, C\}$	$\{A, C, B\}$	$\{B\}$	$\{C, B\}$
3	$\{S', S, A, C\}$	$\{S, A, C, B\}$	$\{A, C, B\}$	$\{B\}$	$\{C, B\}$
4	$\{S', S, A, C, B\}$	$\{S, A, C, B\}$	$\{A, C, B\}$	$\{B\}$	$\{C, B\}$
4	$\{S', S, A, C, B\}$	$\{S, A, C, B\}$	$\{A, C, B\}$	$\{B\}$	$\{C, B\}$

So the grammar G_C is (after removing duplicate right-hand sides):

$$\begin{aligned} S' &\rightarrow AS \mid aA \mid bB \mid b \mid cC \\ S &\rightarrow AS \mid aA \mid bB \mid b \mid cC \\ A &\rightarrow aA \mid bB \mid cC \mid b \\ B &\rightarrow bB \mid b \\ C &\rightarrow cC \mid bB \mid b \end{aligned}$$

D. (extra) Grammar G :

$$\begin{aligned} S &\rightarrow AB \mid C \\ A &\rightarrow aA \mid B \\ B &\rightarrow bB \mid C \\ C &\rightarrow cC \mid a \mid A \end{aligned}$$

The grammar already has a non-recursive start symbol.

$$\begin{aligned} S &\rightarrow AB \mid C \\ A &\rightarrow aA \mid B \\ B &\rightarrow bB \mid C \\ C &\rightarrow cC \mid a \mid A \end{aligned}$$

The chain variables are:

$$\begin{aligned} \text{chain}(S) &= \{S, C, A, B\} \\ \text{chain}(A) &= \{A, B, C\} \\ \text{chain}(B) &= \{B, C, A\} \\ \text{chain}(C) &= \{C, A, B\} \end{aligned}$$

So the grammar G_C is:

$$\begin{aligned} S &\rightarrow AB \mid aA \mid bB \mid cC \mid a \\ A &\rightarrow aA \mid bB \mid cC \mid a \\ B &\rightarrow bB \mid aA \mid cC \mid a \\ C &\rightarrow cC \mid a \mid aA \mid bB \end{aligned}$$

It is easy to see that this is equivalent to

$$\begin{aligned} S &\rightarrow AA \mid aA \mid bA \mid cA \mid a \\ A &\rightarrow aA \mid bA \mid cA \mid a \end{aligned}$$

5.3 Useless nonterminals

For each of the following grammars, construct an equivalent grammar G_T without nonterminating symbols, and then an equivalent grammar G_U without unreachable (and without nonterminating) symbols. Use the algorithms given in the lecture (and in the book).

E. Grammar G :

$$\begin{aligned} S &\rightarrow AA \mid CD \mid bB \\ A &\rightarrow aA \mid a \\ B &\rightarrow bB \mid bC \\ C &\rightarrow cB \\ D &\rightarrow dD \mid d \end{aligned}$$

Given the grammar

$$\begin{aligned} S &\rightarrow AA \mid CD \mid bB \\ A &\rightarrow aA \mid a \\ B &\rightarrow bB \mid bC \\ C &\rightarrow cB \\ D &\rightarrow dD \mid d \end{aligned}$$

The terminating variables can be determined using the strategy from the book.

Iteration	<i>term</i>	<i>prev</i>
0	$\{A, D\}$	
1	$\{A, D, S\}$	$\{A, D\}$
2	$\{A, D, S\}$	$\{A, D, S\}$

So B and C are non-terminating. This results in

$$\begin{aligned} S &\rightarrow AA \\ A &\rightarrow aA \mid a \\ D &\rightarrow dD \mid d \end{aligned}$$

Now the following applies

Iteration	<i>reach</i>	<i>prev</i>
0	$\{S\}$	
1	$\{S, A\}$	$\{S\}$
2	$\{S, A\}$	$\{S, A\}$

So D is not reachable; we get

$$\begin{aligned} S &\rightarrow AA \\ A &\rightarrow aA \mid a \end{aligned}$$

F. (extra) Grammar G :

$$\begin{aligned} S &\rightarrow ACH \mid BB \\ A &\rightarrow aA \mid aF \\ B &\rightarrow CFH \mid b \\ C &\rightarrow aC \mid DH \\ D &\rightarrow aD \mid BD \mid Ca \\ F &\rightarrow bB \mid b \\ H &\rightarrow dH \mid d \end{aligned}$$

Given the grammar

$$\begin{aligned} S &\rightarrow ACH \mid BB \\ A &\rightarrow aA \mid aF \\ B &\rightarrow CFH \mid b \\ C &\rightarrow aC \mid DH \\ D &\rightarrow aD \mid BD \mid Ca \\ F &\rightarrow bB \mid b \\ H &\rightarrow dH \mid d \end{aligned}$$

The terminating variables:

Iteration	<i>term</i>	<i>prev</i>
0	$\{B, F, H\}$	
1	$\{B, F, H, S, A\}$	$\{B, F, H\}$
2	$\{B, F, H, S, A\}$	$\{B, F, H, S, A\}$

So C and D are non-terminating. This results in

$$\begin{aligned} S &\rightarrow B B \\ A &\rightarrow a A \mid a F \\ B &\rightarrow b \\ F &\rightarrow b B \mid b \\ H &\rightarrow d H \mid d \end{aligned}$$

Now the following applies

Iteration	<i>reach</i>	<i>prev</i>
0	{ S }	
1	{ S, B }	{ S }
2	{ S, B }	{ S, B }

So A, F and H are not reachable. This results in:

$$\begin{aligned} S &\rightarrow B B \\ B &\rightarrow b \end{aligned}$$

5.4 Finalisation of CNF rules

For each of the following grammars, construct an equivalent one in Chomsky normal form (CNF). (Note that the grammars already have a nonrecursive start symbol, are essentially noncontracting, and contain neither chain rules nor useless symbols.)

G. Grammar G :

$$\begin{aligned} S &\rightarrow aA \mid ABa \\ A &\rightarrow AA \mid a \\ B &\rightarrow AbB \mid bb \end{aligned}$$

Given the grammar

$$\begin{aligned} S &\rightarrow aA \mid ABa \\ A &\rightarrow AA \mid a \\ B &\rightarrow AbB \mid bb \end{aligned}$$

The Chomsky normal form is as follows

$$\begin{aligned} S &\rightarrow A'A \mid AT_1 \\ A' &\rightarrow a \\ T_1 &\rightarrow BA' \\ A &\rightarrow AA \mid a \\ B &\rightarrow AT_2 \mid B'B' \\ T_2 &\rightarrow B'B \\ B' &\rightarrow b \end{aligned}$$

H. Grammar G :

$$\begin{aligned} S &\rightarrow aAbB \mid ABC \mid a \\ A &\rightarrow aA \mid a \\ B &\rightarrow bBcC \mid b \\ C &\rightarrow abc \end{aligned}$$

Given the grammar

$$\begin{aligned} S &\rightarrow aAbB \mid ABC \mid a \\ A &\rightarrow aA \mid a \\ B &\rightarrow bBcC \mid b \\ C &\rightarrow abc \end{aligned}$$

The Chomsky normal form is as follows

$$\begin{aligned} S &\rightarrow A'T_1 \mid AT_3 \mid a \\ A' &\rightarrow a \\ T_1 &\rightarrow AT_2 \\ T_2 &\rightarrow B'B \\ B' &\rightarrow b \\ T_3 &\rightarrow BC \\ A &\rightarrow A'A \mid a \\ B &\rightarrow B'T_4 \mid b \\ T_4 &\rightarrow BT_5 \\ T_5 &\rightarrow C'C \\ C' &\rightarrow c \\ C &\rightarrow A'T_6 \\ T_6 &\rightarrow B'C' \end{aligned}$$

Note: The right-hand sides of the grammar in Exercise G have at most three symbols. The main difference in Exercise H is that you have to deal with right-hand sides of more than three symbols. During the tutorial, after doing G, we recommend that you start with H, but as soon as you feel that you have understood the principle, move on to the next exercise.

5.5 The CYK parsing algorithm

- I. Given the context-free (CNF) grammar $G = (V, \Sigma, P, S)$ with $V = \{S, A, B, C, D, E\}$, $\Sigma = \{a, b, c\}$ and P consisting of the rules

$$\begin{aligned} S &\rightarrow CA \mid BD \\ A &\rightarrow a \\ B &\rightarrow b \mid c \\ C &\rightarrow AE \\ D &\rightarrow EB \mid b \\ E &\rightarrow BD \mid CA \end{aligned}$$

- (a) Construct the X_{ij} -matrix for the words $acba$ and $baca$ according to the CYK-algorithm.

- (b) Use the X_{ij} -matrix to determine whether $acba$ and $baca$ belong to $\mathcal{L}(G)$, and if yes, what the corresponding derivation tree is.
- (c) (*extra*) Would you call the CYK-algorithm a “top-down” or a “bottom-up” parsing algorithm, or neither? Motivate your answer.

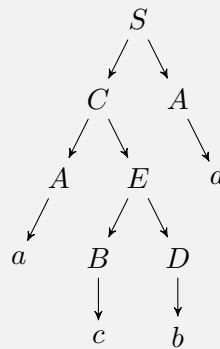
a. The X_{ij} -matrices for $acba$ and $baca$ according to CYK:

	1	2	3	4
1	{A}	–	{C}	{S, E}
2		{B}	{S, E}	–
3			{B, D}	–
4				{A}

	1	2	3	4
1	{B, D}	–	–	–
2		{A}	–	–
3			{B}	–
4				{A}

So $baca$ is not recognized!

- b. Because S is on the top right of the matrix for $acba$, this word belongs to the language $\mathcal{L}(G)$. This is not the case for $baca$. The derivation tree for $acba$ can be read directly from the matrix, using the corresponding rules:



- c. The decomposition algorithm follows a bottom-up approach. Indeed, we start with the decomposition of all substrings of length 1, next all substrings of length 2, etc., until finally, we arrive at the string itself.

J. (*extra*) Given the context-free grammar $G = (V, \Sigma, P, S)$ with $V = \{S, A, B, C, D\}$, $\Sigma = \{a, b, c\}$ and P consisting of the rules

$$\begin{aligned}
 S &\rightarrow AC \mid DB \\
 A &\rightarrow a \\
 B &\rightarrow b \mid c \\
 C &\rightarrow SA \\
 D &\rightarrow BS \mid b
 \end{aligned}$$

(You can ignore the recursion in the start symbol S for this exercise.)

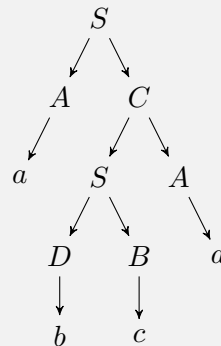
- (a) Construct for the words $abca$ and $acab$ the X_{ij} -matrix.
- (b) Use the X_{ij} -matrix to determine whether $abca$ and $acab$ belong to $\mathcal{L}(G)$, and if yes, what the corresponding derivation tree is.

a. The X_{ij} -matrices for $abca$ and $acab$ according to CYK:

	1	2	3	4
1	{A}	–	–	{S}
2		{B, D}	{S}	{C}
3			{B}	–
4				{A}

	1	2	3	4
1	{A}	–	–	–
2		{B}	–	–
3			{A}	–
4				{B, D}

b. Because S is on the top right of the matrix for $abca$, this word belongs to the language $\mathcal{L}(G)$. This is not the case for $acab$. The derivation tree for $abca$ can be read directly from the matrix, using the corresponding rules:



K. (extra) Give the upper diagonal matrix produced by the CYK algorithm when run with the Chomsky normal form grammar of G below and the input strings $abbb$ and $aabbb$.

$$S \rightarrow aSb \mid ab$$

For convenience, the (CNF) grammar of G is given here:

$$\begin{aligned} S &\rightarrow AT \mid AB \\ T &\rightarrow XB \\ X &\rightarrow AT \mid AB \\ A &\rightarrow a \\ B &\rightarrow b \end{aligned}$$

The X_{ij} -matrices for $abbb$ and $aabbb$:

	1	2	3	4
1	{A}	{S, X}	{T}	–
2		{B}	–	–
3			{B}	–
4				{B}

	1	2	3	4	5
1	{A}	–	–	{S, X}	{T}
2		{A}	{S, X}	{T}	–
3			{B}	–	–
4				{B}	–
5					{B}

So neither of these words are accepted. Looking at the original grammar and language of G , we can see this directly: the language was $\{a^i b^i \mid i \geq 1\}$, and neither $abbb$ nor $aabbb$ are elements of this.