

The pumping lemma for contex-free languages

Definition (Unit rules and λ -rules)

A *unit rule* is a rule of the form $X \rightarrow Y$ where X and Y are variable symbols.

A λ -rule is a rule of the form $X \rightarrow \lambda$ where X is a variable symbol.

Lemma (Removing unit rules and λ -rules)

For a given context-free grammar G one can effectively construct a context-free grammar G' such that

$$L(G') = L(G) \setminus \{\lambda\}$$

and the set of rules of G' contains neither unit rules nor λ -rules.

Proof. See lecture.

The pumping lemma for contex-free languages

Definition (Chomsky normal form for context-free grammars)

A context-free grammar $G = (N, T, P, S)$ is in *Chomsky normal form* if all rules in P are of the form

$$X \rightarrow YZ \quad \text{or} \quad X \rightarrow a \quad \text{where } X, Y, Z \in N \text{ and } a \in T,$$

with the possible exception that P may contain the rule $S \rightarrow \lambda$, in which case S does not occur on the right-hand side of any rule in P .

Lemma (Transformation into Chomsky normal form)

For a given context-free grammar G one can effectively construct a context-free grammar G' in Chomsky normal form such that $L(G) = L(G')$. In addition, the grammar G' can be chosen such that all its variable symbols are useful.

The pumping lemma for contex-free languages

Definition (Useful variable symbols)

Let $G = (N, T, P, S)$ be a context-free grammar.

A variable symbol X is *useful* if there is a word $z \in T^*$ such that X occurs in some derivation of z , i.e., $z = uvw$ and

$$S \xRightarrow{G,*} uXw \xRightarrow{G,*} uvw;$$

otherwise, X is *useless*.

Lemma (Removing useless variable symbols)

If one cancels in a given context-free grammar G all useless variable symbols and all rules in which these variable symbols occur, one obtains a grammar G' that has only useful variable symbols and where $L(G) = L(G')$.

Sketch of proof. It suffices to observe that cancelling a useless variable symbol does not create new useless variable symbols. \square

The pumping lemma for contex-free languages

Proof. We construct context-free grammars G_1 , G_2 , and G_3 where G_3 is in Chomsky normal form and such that

$$L(G) = L(G_1) \quad \text{and} \quad L(G_1) \setminus \{\lambda\} = L(G_2) = L(G_3).$$

In case $\lambda \notin L(G)$, we can simply let $G' = G_3$.

Otherwise, we obtain G' as required by adding to G_3 a new start symbol S' and rules $S' \rightarrow \lambda$ and $S' \rightarrow S$, where S is the start symbol of G_3 .

By the lemma above, we can in addition ensure that all variable symbols of G' are useful by removing all useless variable symbols.

The pumping lemma for contex-free languages

Proof, cont.: (Replace terminal symbols a by Z_a)

The grammar G_1 is obtained from G as follows.

For every $a \in T$

replace a by Z_a in all rules of P ,

add a variable Z_a and a rule $Z_a \rightarrow a$,

where the Z_a are mutually distinct new variables.

By construction, we have $L(G) = L(G_1)$.

The rules in G_1 have the form $X \rightarrow a$ for a terminal symbol a or

$X \rightarrow Y_1 \cdots Y_t$ for some $t \geq 0$ and variable symbols Y_1, \dots, Y_t .

The pumping lemma for contex-free languages

Proof, cont.: (Remove multiple variable symbols)

The grammar G_3 is obtained from G_2 as follows.

Successively, each rule of the form $X \rightarrow Y_1 \cdots Y_t$ where $t \geq 3$ is replaced by the rules

$$\begin{aligned} X &\rightarrow Y_1 Z_1, \\ Z_1 &\rightarrow Y_2 Z_2, \\ &\vdots \\ Z_{t-1} &\rightarrow Y_{t-1} Y_t, \end{aligned}$$

where for each replacement the Z_i are chosen as mutually distinct new variable symbols.

By construction, we have $L(G_2) = L(G_3)$.

The rules in G_3 have the form $X \rightarrow a$ for a terminal symbol a or

$X \rightarrow Y_1 Y_2$ for variable symbols Y_1 and Y_2 . \square

The pumping lemma for contex-free languages

Proof, cont.: (Remove unit rules and λ -rules)

The grammar G_2 is obtained from G_1 as follows.

Let G_2 be a grammar equivalent to G_1 and without unit rules and λ -rules as in the proof of the corresponding lemma above.

By construction, we have $L(G_1) \setminus \{\lambda\} = L(G_2)$.

The rules in G_2 have the form $X \rightarrow a$ for a terminal symbol a or

$X \rightarrow Y_1 \cdots Y_t$ for some $t \geq 2$ and variable symbols Y_1, \dots, Y_t .

The pumping lemma for contex-free languages

In what follows, we derive a pumping lemma for context-free languages, as well as a variant for the subclass of linear languages.

Similar to the case of regular languages, these pumping lemmas are the standard tools for showing that a certain language is not context-free or is not linear.

Theorem (Pumping lemma for context-free languages)

Let L be a context-free language. Then there is a constant k such that every word $z \in L$ of length at least k can be written in the form

$$z = uvwxy$$

where the words u , v , w , x , and y have the following properties

- (i) $|vwx| \leq k$,
- (ii) $vx \neq \lambda$,
- (iii) $uv^iwx^iy \in L$ for all $i \geq 0$.

The pumping lemma for contex-free languages

Proof. Fix some context-free grammar $G = (N, T, P, S)$ in Chomsky normal form that generates L and let $k = 2^{|N|+1}$.

For any given word $z \in L$ such that $|z| \geq k$, consider a left derivation α of z in G and the corresponding parse tree $T(\alpha)$.

The parse tree $T(\alpha)$ has depth of at least $|N| + 2$ because a parse tree of depth at most $|N| + 1$ has at most $2^{|N|} < |z|$ leave nodes.

Fix a path between the root and a leave node of $T(\alpha)$ of maximum length among all such paths.

The length of the path is at least $|N| + 1$, hence there is a node K on the path at distance $|N| + 1$ from the leave node of the path.

The subtree of $T(\alpha)$ with root K has depth $|N| + 1$ because the leave node of the path is at this distance from K while there cannot be any node at larger distance, otherwise the path would not have maximum length.

The pumping lemma for contex-free languages

Proof, cont.: The part of the path between K and the leave node of the path contains $|N| + 2$ nodes, i.e., contains $|N| + 1$ nodes that are not leave nodes.

Each of these $|N| + 1$ nodes is marked with one of the $|N|$ variable symbols, i.e., there are two distinct nodes K' and K'' that are marked with the same variable symbol, say, with X .

So there are words u, v, w, x , and y over T such that

$$S \xRightarrow{G,*} uXy \xRightarrow{G,t} uvXxy \xRightarrow{G,*} uvwxy = z \quad \text{for some } t > 0.$$

From this derivation we obtain

- (i) $|vwx| \leq 2^{|N|+1} = k$ by choice of the node K ,
- (ii) $vx \neq \lambda$ by $t > 0$ and because P contains neither unit rules nor λ -rules,
- (iii) $uv^iwx^iy \in L$ for all $i \geq 0$ because $X \xRightarrow{G,*} v^iwy^i$ for all $i \geq 0$.

□

The pumping lemma for contex-free languages

Example (A language that is not context-free)

The language $L = \{0^n 1^n 0^n : n > 0\}$ is not context-free.

For a proof by contradiction, assume that L is context-free.

Choose both, a constant k and a partition $uvwxy$ of $0^k 1^k 0^k$, as in the pumping lemma for context-free languages, i.e.,

$$0^k 1^k 0^k = uvwxy, \quad |vwx| \leq k, \quad vx \neq \lambda, \quad \text{and } uwy \in L.$$

Then at least one of the words u and y must have length at least k .

So the word uwy has a prefix 0^k or a suffix 1^k but has length strictly less than $3k$.

Thus $uwy \notin L$, contradicting the choice of u, w , and y .

The pumping lemma for contex-free languages

Example (Words of prime length)

The language $L = \{w \in \{1\}^* : |w| \text{ is prime}\}$ of all words of prime length over the unary alphabet is not context-free.

For a proof by contradiction, assume that L_1 is regular and accordingly choose a constant k as in the pumping lemma for context-free languages and a prime number $p \geq k + 2$.

Let $z = uvwxy$ be a partition of $z = 1^p$ as in the pumping lemma.

For $m = |vx|$, it holds that $1 \leq m \leq |vwx| \leq k \leq p - 2$ (*).

Furthermore, the word $uv^{p-m}wx^{p-m}y$ is in L , however its length cannot be prime because of

$$|uv^{p-m}wx^{p-m}y| = \underbrace{|uwy|}_{=p-m} + m(p-m) = \underbrace{(m+1)}_{\geq 2 \text{ by } (*)} \underbrace{(p-m)}_{\geq 2 \text{ by } (*)},$$

which contradicts the definition of L .

□

The pumping lemma for context-free languages

Definition (Linear grammars and languages)

A grammar $G = (N, T, P, S)$ is *linear* if all its rules are of the form

$$X \rightarrow uYv \quad \text{where } X, Y \in N \text{ and } u, v \in T^*.$$

A language is *linear* if it is generated by a linear grammar.

Example (A linear and a nonlinear languages)

The language $L = \{0^n 1^n : n \geq 0\}$ is generated by the linear grammar $(\{S\}, \{0, 1\}, \{S \rightarrow 0S1 \mid \lambda\}, S)$, hence L is linear.

The language $L = \{0^m 1^m 0^n 1^n : m, n \geq 0\}$ is not linear, see below.

The pumping lemma for context-free languages

The following pumping lemma for linear languages is the standard tool for showing that a language is not linear.

Theorem (Pumping lemma for linear languages)

For every linear language L there is a constant k such that every word $z \in L$ of length at least k can be written in the form

$$z = uvwxy$$

where the words u, v, w, x , and y have the following properties

- (i) $|uv| \leq k$ and $|xy| \leq k$,
- (ii) $vx \neq \lambda$,
- (iii) $uv^i wx^i y \in L$ for all $i \geq 0$.

The pumping lemma for context-free languages

Definition (Chomsky normal form for linear grammars)

A linear grammar $G = (N, T, P, S)$ is in *Chomsky normal form* if all rules in P are of one of the forms

$$X \rightarrow aY, \text{ or } X \rightarrow Ya, \text{ or } X \rightarrow a \quad \text{where } X, Y \in N \text{ and } a \in T,$$

with the possible exception that P may contain the rule $S \rightarrow \lambda$, in which case S does not occur on the right-hand side of any rule in P .

Lemma (Transformation into Chomsky normal form)

For a given linear grammar G one can effectively construct a linear grammar G' in Chomsky normal form such that $L(G) = L(G')$.

In addition, the grammar G' can be chosen such that all its variable symbols are useful.

Proof. Similar to the context-free case, see lecture.

The pumping lemma for context-free languages

Proof. Fix some linear grammar $G = (N, T, P, S)$ in Chomsky normal form that generates L and let $k = |N| + 1$.

For any given word $z \in L$ such that $|z| \geq k$, consider a derivation α of z in G , which then must have length of at least k .

Each of the first k sentential forms that occur in α contain a single variable symbol, hence some variable symbol X occurs twice, i.e., there are words u, v, w, x , and y in T^* such that

$$S \xRightarrow{G, t_1} uXy \xRightarrow{G, t_2} uvXxy \xRightarrow{G, *} uvwxy = z,$$

$$0 \leq t_1, 0 < t_2 \text{ and } t_1 + t_2 \leq k - 1.$$

Since G is a linear grammar in Chomsky normal form we obtain

- (i) $|uv| + |xy| \leq t_1 + t_2 \leq k$,
- (ii) $xy \neq \lambda$ by $t_2 > 0$,
- (iii) $uv^i wx^i y \in L$ for all $i \geq 0$ because $X \xRightarrow{G, *} v^i w y^i$ for all $i \geq 0$.

□

The pumping lemma for context-free languages

Example (Linear and context-free languages)

The language $L = \{0^m 1^m 0^n 1^n : m, n \geq 0\}$ is context-free but not linear.

Obviously L can be generated by a context-free grammar.

In order to prove that L is not linear, assume otherwise and let k be a constant as in the pumping lemma for linear languages.

Fix a partition $uvwxy$ of $z = 0^k 1^k 0^k 1^k$ according to the pumping lemma, i.e.,

$$0^k 1^k 0^k 1^k = uvwxy.$$

Then we have by (i) $|uv| \leq k$ and $|xy| \leq k$ and by (ii) $vx \neq \lambda$, hence $v = 0^s$ and $x = 1^t$ where $s + t > 0$.

Consequently, $uw^s x \notin L$, contradicting the choice of $uvwxy$ and the pumping property (iii).

The pumping lemma for context-free languages

Corollary (Regular, linear, and context-free languages)

The classes of regular, linear, and context-free languages form a strict hierarchy in the sense that

$$\{L : L \text{ regular}\} \subsetneq \{L : L \text{ linear}\} \subsetneq \{L : L \text{ context-free}\}.$$

Proof. The inclusion relations hold because every right-linear grammar is linear, and every linear grammar is context-free.

That the inclusions are proper is witnessed by already discussed counterexamples, i.e., by the languages

$$\{0^n 1^n : n \geq 0\}, \text{ which is linear but not regular and } \{0^m 1^m 0^n 1^n : m, n \geq 0\}, \text{ which is context-free but not linear.}$$

□

The pumping lemma for context-free languages

Theorem (Closure properties of the context-free languages)

The class of context-free languages is closed under union, concatenation, and Kleene closure, i.e.,

- (i) if L_1 and L_2 are context-free, then $L_1 \cup L_2$ is context-free,
- (ii) if L_1 and L_2 are context-free, then $L_1 L_2$ is context-free,
- (iii) if L is context-free, then L^* is context-free,

Proof. (i), (ii) For given context-free languages L_1 and L_2 , let $G_1 = (N_1, T_1, P_1, S_1)$ and $G_2 = (N_2, T_2, P_2, S_2)$ be context-free grammars where $L_1 = L(G_1)$ and $L_2 = L(G_2)$.

The pumping lemma for context-free languages

Proof, cont.: Then the languages $L_1 \cup L_2$ and $L_1 L_2$ are generated by the context-free grammars

$$(N_1 \cup N_2 \cup \{S\}, T_1 \cup T_2, P_1 \cup P_2 \cup \{S \rightarrow S_1 | S_2\}, S) \text{ and } (N_1 \cup N_2 \cup \{S\}, T_1 \cup T_2, P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\}, S),$$

respectively, where S is a new variable symbol and for the latter grammar one has to assume that N_1 and N_2 are disjoint.

(iii) For a given context-free language L , let $G = (N, T, P, S)$ be a context-free grammar where $L = L(G)$.

Then the language L^* is generated by the context-free grammar

$$(N \cup \{S'\}, T, P \cup \{S' \rightarrow SS' | \lambda\}, S')$$

where S' is a new variable symbol.

□

The pumping lemma for context-free languages

Remark (Closure properties of the context-free languages)

The class of context-free languages is closed under transition to the mirror language, i.e., if a context-free language L is context-free, then the mirror language $L^R = \{w^R : w \in L\}$ is context-free, too.

For a proof, let the language L be context-free.

Choose a context-free grammar $G = (N, T, P, S)$ where $L = L(G)$.

Then the language L^R is generated by the context-free grammar $G^R = (N, T, P^R, S)$ where

$$P^R = \{X \rightarrow w^R : X \rightarrow w \in P\}.$$

It can be shown by induction over n for all words $w \in (N \cup T)^*$ and all n that we have $S \xRightarrow{G,n} w$ if and only if $S \xRightarrow{G^R,n} w^R$.

The pumping lemma for context-free languages

Theorem (Closure properties of the context-free languages)

The class of context-free languages is neither closed under intersection nor under complement, i.e.,

- (i) L_1 and L_2 context-free does not imply $L_1 \cap L_2$ context-free,
- (ii) L context-free does not imply that \bar{L} is context-free.

Proof. (i) The languages

$L_1 = \{0^m 1^m 0^n : m, n \geq 0\}$ and $L_2 = \{0^m 1^n 0^n : m, n \geq 0\}$,
are both context-free but have the non-context-free intersection

$$L_1 \cap L_2 = \{0^m 1^m 0^m : m \geq 0\}$$

(ii) The context-free languages are closed under union, hence closure under complement would imply closure under intersection by de Morgan's law $L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$. \square

It can be shown by specifying an appropriate grammar that the complement of the language $\{0^m 1^m 0^m : m \geq 0\}$ is context-free.