

Theory of Computer Science

C3. Regular Languages: Regular Expressions, Pumping Lemma

Gabriele Röger

University of Basel

March 25, 2019

Theory of Computer Science

March 25, 2019 — C3. Regular Languages: Regular Expressions, Pumping Lemma

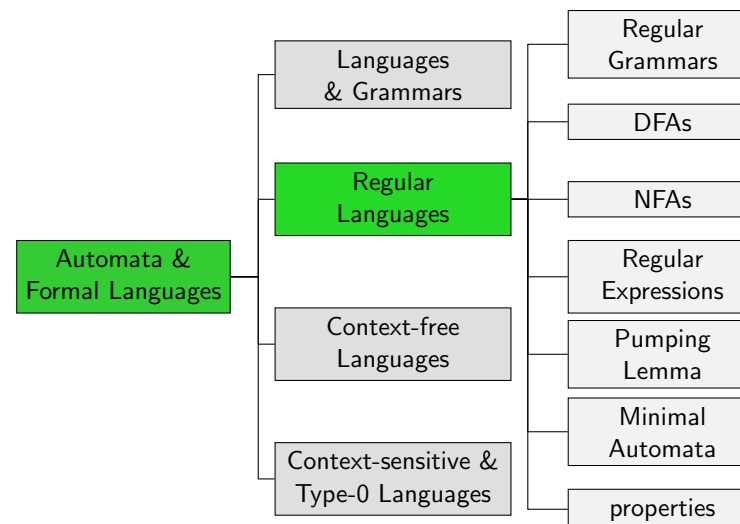
C3.1 Regular Expressions

C3.2 Pumping Lemma

C3.3 Summary

C3.1 Regular Expressions

Overview

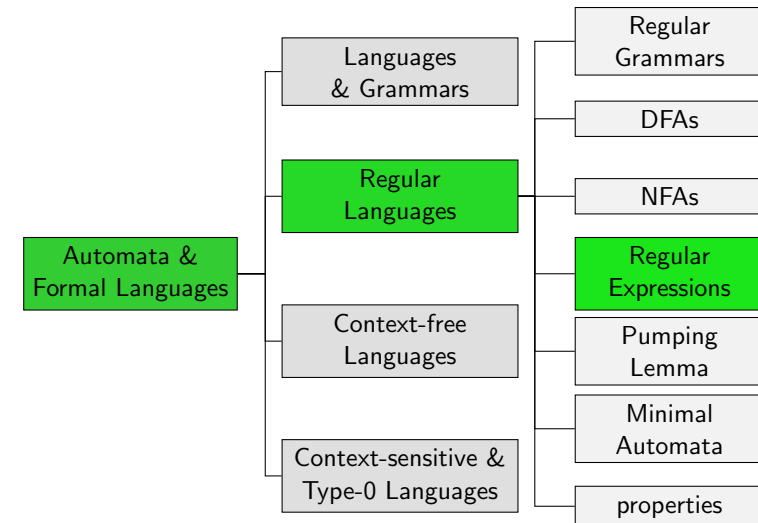


Formalisms for Regular Languages

- ▶ DFAs, NFAs and regular grammars can all describe exactly the regular languages.
- ▶ Are there other concepts with the same expressiveness?
- ▶ **Yes!** \rightsquigarrow regular expressions

Live demo

Overview



Regular Expressions: Definition

Definition (Regular Expressions)

Regular expressions over an alphabet Σ are defined inductively:

- ▶ \emptyset is a regular expression
- ▶ ε is a regular expression
- ▶ If $a \in \Sigma$, then a is a regular expression

If α and β are regular expressions, then so are:

- ▶ $(\alpha\beta)$ (**concatenation**)
- ▶ $(\alpha|\beta)$ (**alternative**)
- ▶ (α^*) (**Kleene closure**)

German: reguläre Ausdrücke, Verkettung, Alternative, kleenesche Hülle

Regular Expressions: Omitting Parentheses

omitted parentheses by convention:

- ▶ Kleene closure α^* binds more strongly than concatenation $\alpha\beta$.
- ▶ Concatenation binds more strongly than alternative $\alpha|\beta$.
- ▶ Parentheses for nested concatenations/alternatives are omitted (we can treat them as left-associative; it does not matter).

Example: $ab^*c|\varepsilon|abab^*$ abbreviates $((((a(b^*))c)|\varepsilon)|(((ab)a)(b^*)))$.

Regular Expressions: Examples

some regular expressions for $\Sigma = \{0, 1\}$:

- ▶ 0^*10^*
- ▶ $(0|1)^*1(0|1)^*$
- ▶ $((0|1)(0|1))^*$
- ▶ $01|10$
- ▶ $0(0|1)^*0|1(0|1)^*1|0|1$

Regular Expressions: Language

Definition (Language Described by a Regular Expression)

The **language described by a regular expression** γ , written $\mathcal{L}(\gamma)$, is inductively defined as follows:

- ▶ If $\gamma = \emptyset$, then $\mathcal{L}(\gamma) = \emptyset$.
- ▶ If $\gamma = \varepsilon$, then $\mathcal{L}(\gamma) = \{\varepsilon\}$.
- ▶ If $\gamma = a$ with $a \in \Sigma$, then $\mathcal{L}(\gamma) = \{a\}$.
- ▶ If $\gamma = (\alpha\beta)$, where α and β are regular expressions, then $\mathcal{L}(\gamma) = \mathcal{L}(\alpha)\mathcal{L}(\beta)$.
- ▶ If $\gamma = (\alpha|\beta)$, where α and β are regular expressions, then $\mathcal{L}(\gamma) = \mathcal{L}(\alpha) \cup \mathcal{L}(\beta)$.
- ▶ If $\gamma = (\alpha^*)$ where α is a regular expression, then $\mathcal{L}(\gamma) = \mathcal{L}(\alpha)^*$.

Examples: blackboard

Finite Languages Can Be Described By Regular Expressions

Theorem

Every finite language can be described by a regular expression.

Proof.

For every word $w \in \Sigma^*$, a regular expression describing the language $\{w\}$ can be built from regular expressions $a \in \Sigma$ by using concatenations.

(Use ε if $w = \varepsilon$.)

For every finite language $L = \{w_1, w_2, \dots, w_n\}$, a regular expression describing L can be built from the regular expressions for $\{w_i\}$ by using alternatives.

(Use \emptyset if $L = \emptyset$.) □

Regular Expressions Not More Powerful Than NFAs

Theorem

For every language that can be described by a regular expression, there is an NFA that accepts it.

Proof.

Let γ be a regular expression.

We show the statement by induction over the structure of regular expressions.

For $\gamma = \emptyset$, $\gamma = \varepsilon$ and $\gamma = a$,
NFAs that accept $\mathcal{L}(\gamma)$ are obvious. ...

Regular Expressions Not More Powerful Than NFAs

Theorem

For every language that can be described by a regular expression, there is an NFA that accepts it.

Proof (continued).

For $\gamma = (\alpha\beta)$, let M_α and M_β be NFAs that (by ind. hypothesis) accept $\mathcal{L}(\alpha)$ and $\mathcal{L}(\beta)$. W.l.o.g., their states are disjoint.

Construct NFA M for $\mathcal{L}(\gamma)$ by “daisy-chaining” M_α and M_β :

- ▶ states: union of states of M_α and M_β
- ▶ start states: those of M_α ; if $\varepsilon \in \mathcal{L}(\alpha)$, also those of M_β
- ▶ end states: end states of M_β
- ▶ state transitions: all transitions of M_α and of M_β ;
 additionally: for every transition to an end state of M_α ,
 an equally labeled transition to all start states of M_β

...

Regular Expressions Not More Powerful Than NFAs

Theorem

For every language that can be described by a regular expression, there is an NFA that accepts it.

Proof (continued).

For $\gamma = (\alpha|\beta)$, by the induction hypothesis let $M_\alpha = \langle Q_\alpha, \Sigma, \delta_\alpha, S_\alpha, E_\alpha \rangle$ and $M_\beta = \langle Q_\beta, \Sigma, \delta_\beta, S_\beta, E_\beta \rangle$ be NFAs that accept $\mathcal{L}(\alpha)$ and $\mathcal{L}(\beta)$.

W.l.o.g., $Q_\alpha \cap Q_\beta = \emptyset$.

Then the “union automaton”

$$M = \langle Q_\alpha \cup Q_\beta, \Sigma, \delta_\alpha \cup \delta_\beta, S_\alpha \cup S_\beta, E_\alpha \cup E_\beta \rangle$$

accepts the language $\mathcal{L}(\gamma)$.

...

German: Vereinigungsautomat

Regular Expressions Not More Powerful Than NFAs

Theorem

For every language that can be described by a regular expression, there is an NFA that accepts it.

Proof (continued).

For $\gamma = (\alpha^*)$, by the induction hypothesis let $M_\alpha = \langle Q_\alpha, \Sigma, \delta_\alpha, S_\alpha, E_\alpha \rangle$ be an NFA that accepts $\mathcal{L}(\alpha)$.

If $\varepsilon \notin \mathcal{L}(\alpha)$, add an additional state to M_α that is a start and end state and not connected to other states. M_α now recognizes $\mathcal{L}(\alpha) \cup \{\varepsilon\}$.

M is constructed from M_α by adding the following new transitions:
 whenever M_α has a transition from s to end state s' with symbol a ,
 add transitions from s to every start state with symbol a .

Then $\mathcal{L}(M) = \mathcal{L}(\gamma)$. □

DFAs Not More Powerful Than Regular Expressions

Theorem

Every language accepted by a DFA can be described by a regular expression.

Without proof.

Regular Languages vs. Regular Expressions

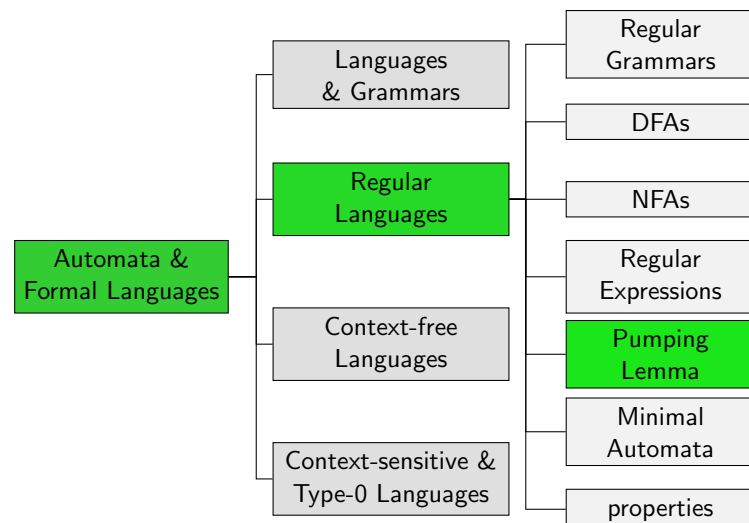
Theorem (Kleene)

The set of languages that can be described by regular expressions is exactly the set of regular languages.

This follows directly from the previous two theorems.

C3.2 Pumping Lemma

Overview



Pumping Lemma: Motivation



You can show that a language is regular by specifying an appropriate grammar, finite automaton, or regular expression. How can you show that a language is **not** regular?

- ▶ Direct proof that no regular grammar exists that generates the language
 \rightsquigarrow difficult in general
- ▶ **Pumping lemma**: use a necessary property that holds for all regular languages.

Pumping Lemma

Theorem (Pumping Lemma)

Let L be a regular language. Then there is an $n \in \mathbb{N}$ (a *pumping number* for L) such that all words $x \in L$ with $|x| \geq n$ can be split into $x = uvw$ with the following properties:

- 1 $|v| \geq 1$,
- 2 $|uv| \leq n$, and
- 3 $uv^i w \in L$ for all $i = 0, 1, 2, \dots$

Question: what if L is finite?

Pumping Lemma: Proof

Theorem (Pumping Lemma)

Let L be a regular language. Then there is an $n \in \mathbb{N}$ (a *pumping number* for L) such that all words $x \in L$ with $|x| \geq n$ can be split into $x = uvw$ with the following properties:

- 1 $|v| \geq 1$,
- 2 $|uv| \leq n$, and
- 3 $uv^i w \in L$ for all $i = 0, 1, 2, \dots$

Proof.

For regular L there exists a DFA $M = \langle Q, \Sigma, \delta, q_0, E \rangle$ with $\mathcal{L}(M) = L$. We show that $n = |Q|$ has the desired properties.

Consider an arbitrary $x \in \mathcal{L}(M)$ with length $|x| \geq |Q|$. Including the start state, M visits $|x| + 1$ states while reading x . Because of $|x| \geq |Q|$ at least one state has to be visited twice. ...

Pumping Lemma: Proof

Theorem (Pumping Lemma)

Let L be a regular language. Then there is an $n \in \mathbb{N}$ (a *pumping number* for L) such that all words $x \in L$ with $|x| \geq n$ can be split into $x = uvw$ with the following properties:

- 1 $|v| \geq 1$,
- 2 $|uv| \leq n$, and
- 3 $uv^i w \in L$ for all $i = 0, 1, 2, \dots$

Proof (continued).

Choose a split $x = uvw$ so M is in the same state after reading u and after reading uv . Obviously, we can choose the split in a way that $|v| \geq 1$ and $|uv| \leq |Q|$ are satisfied. ...

Pumping Lemma: Proof

Theorem (Pumping Lemma)

Let L be a regular language. Then there is an $n \in \mathbb{N}$ (a *pumping number* for L) such that all words $x \in L$ with $|x| \geq n$ can be split into $x = uvw$ with the following properties:

- 1 $|v| \geq 1$,
- 2 $|uv| \leq n$, and
- 3 $uv^i w \in L$ for all $i = 0, 1, 2, \dots$

Proof (continued).

The word v corresponds to a loop in the DFA after reading u and can thus be repeated arbitrarily often. Every subsequent continuation with w ends in the same end state as reading x . Therefore $uv^i w \in \mathcal{L}(M) = L$ is satisfied for all $i = 0, 1, 2, \dots$ \square

Pumping Lemma: Application

Using the pumping lemma (PL):

Proof of Nonregularity

- ▶ If L is regular, then the pumping lemma holds for L .
- ▶ By contraposition: if the PL does not hold for L , then L cannot be regular.
- ▶ That is: if there is no $n \in \mathbb{N}$ with the properties of the PL, then L cannot be regular.

Pumping Lemma: Caveat

Caveat:

The pumping lemma is a **necessary condition** for a language to be regular, but not a **sufficient one**.

- ↪ there are languages that satisfy the pumping lemma conditions but are not regular
- ↪ for such languages, other methods are needed to show that they are not regular (e.g., the [Myhill-Nerode theorem](#))

Pumping Lemma: Example

Example

The language $L = \{a^n b^n \mid n \in \mathbb{N}\}$ is not regular.

Proof.

Assume L is regular. Then let p be a pumping number for L .

The word $x = a^p b^p$ is in L and has length $\geq p$.

Let $x = uvw$ be a split with the properties of the PL.

Then the word $x' = uv^2w$ is also in L . Since $|uv| \leq p$, uv consists only of symbols a and $x' = a^{|u|} a^{2|v|} a^{p-|uv|} b^p = a^{p+|v|} b^p$.

Since $|v| \geq 1$ it follows that $p + |v| \neq p$ and thus $x' \notin L$.

This is a contradiction to the PL. $\rightsquigarrow L$ is not regular. □

Pumping Lemma: Another Example I

Example

The language $L = \{ab^n ac^{n+2} \mid n \in \mathbb{N}\}$ is not regular.

Proof.

Assume L is regular. Then let p be a pumping number for L .

The word $x = ab^p ac^{p+2}$ is in L and has length $\geq p$.

Let $x = uvw$ be a split with the properties of the PL.

From $|uv| \leq p$ and $|v| \geq 1$ we know that uv consists of one a followed by at most $p - 1$ b s.

We distinguish two cases, $|u| = 0$ and $|u| > 0$

Pumping Lemma: Another Example II

Example

The language $L = \{ab^nac^{n+2} \mid n \in \mathbb{N}\}$ is not regular.

Proof (continued).

If $|u| = 0$, then word v starts with an a.

Hence, $uv^0w = b^{p-|v|+1}ac^{p+2}$ does not start with symbol a and is therefore not in L . This is a contradiction to the PL.

If $|u| > 0$, then word v consists only of bs.

Consider $uv^0w = ab^{p-|v|}ac^{p+2}$. As $|v| \geq 1$, this word does not contain two more cs than bs and is therefore not in language L .

This is a contradiction to the PL.

We have in all cases a contradiction to the PL.

$\rightsquigarrow L$ is not regular. \square

C3.3 Summary

Summary

- ▶ **Regular expressions** are another way to describe languages.
- ▶ All regular languages can be described by regular expressions, and all regular expressions describe regular languages.
- ▶ Hence, they are equivalent to finite automata.
- ▶ The **pumping lemma** can be used to show that a language is **not regular**.