

3. Star-free languages

Goal 1: Show that FO[\leq] defines strict subset of regular languages.

Goal 2: Find alternative characterization

↳ Definable in FO[\leq]

↳ Represented by star-free regular expression

Recall: First-order formulas are WMSO-formulas without second-order variables/quantifiers.

Example:

$e = \forall x: P(x) \rightarrow \exists y: x < y \wedge P(y)$ over $\Sigma = \{a, b, c\}$ defines

"Every a is followed by a ."

$L(e) = \{a, b, c\}^* \cdot b \cdot \{b, c\}^* \cup \{b, c\}^*$

Note:

$\text{first}(x)$, $\text{last}(x)$, $x = y$ are still in FO[\leq]

Goal 1:

Known: FO[\leq]-languages are regular

(as FO[\leq] is syntactic restriction of WMSO).

Show: $(aa)^*$ is not FO[\leq]-definable (but regular).

↳ There is no φ in FO[\leq] with $L(\varphi) = (aa)^*$

⇒ FO[\leq]-definable languages form strict subset of regular languages.

3.1 Ehrenfeucht - Fraïssé Games

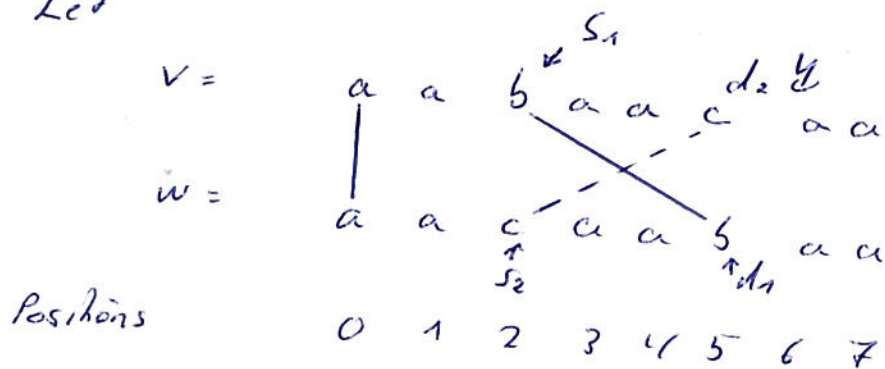
Tool from finite model theory (logic)

for proving inexpressibility results.

The game (informally):

- Two players: spoiler and duplicator
- Two words: v and w over Σ .
- Number of rounds: $k \in \mathbb{N}$.
- Potentially some existing edges

Let



Per round:

- Spoiler selects position in v or w
- Duplicator selects position in other word and connects them by a line

↳ positions have same letters (preserve P_a)

↳ new line does not cross existing lines (preserve $<$)

Next round.

Spoiler may reuse positions:

- then duplicator has to reply by corresponding position
- otherwise duplicator is forced to select new position (if spoiler does not reuse).

- Duplicator loses if cannot reply

- Duplicator wins if number of rounds passes without loss.

In the example:

• Duplicator wins 1 round

• Duplicator loses 2 rounds.

Definition (Partial isomorphism on word structures):

Let $S_v = (D_v, <_v, (P_a)_{a \in \Sigma})$ and

$S_w = (D_w, <_w, (P_a)_{a \in \Sigma})$.

A partial isomorphism (from S_v to S_w) is

a partial function $p: D_v \rightarrow D_w$ so that

(1) p is injective

(2) $\forall x \in \text{dom}(p) : \forall a \in \Sigma, P_{a,v}(x) \text{ iff } P_{a,w}(p(x))$.

// identically labelled.

(3) $\forall x, y \in \text{dom}(p): x \prec y \text{ iff } p(x) \prec_w p(y)$
 \parallel no crossing edges.

Let $\vec{s} = (s_1, \dots, s_n)$ and $\vec{t} = (t_1, \dots, t_n)$.

Typically write $\vec{s} \mapsto \vec{t}$ for $p = \{(s_1, t_1), \dots, (s_n, t_n)\}$

With this definition, interpret EF-game as follows:

- Let S_v, S_w two structures with designated elements \vec{s}, \vec{t} .
- Duplicator tries to establish partial isomorphism
- Spoiler tries to avoid this.

Overall goal of EF games:

- Compare the structures
- If duplicator wins k -rounds, S_v, S_w ^{cannot} be distinguished by formulas of quantifier-depth $\leq k$.

Definition (EF-game):

Let S_v, S_w two word structures
 and \vec{s}, \vec{t} two vectors of positions in S_v and S_w .
 Let $k \in \mathbb{N}$.

The EF-game $G_k((S_v, \vec{s}), (S_w, \vec{t}))$ consists of the following elements and rules:

- k rounds
- Initial configuration $\vec{s} \mapsto \vec{t}$
- Given a configuration v , a round consists of the following moves:
 - ↳ Spoiler chooses $s \in D_v$ or $t \in D_w$
 - ↳ Duplicator chooses $t \in D_w$ or $s \in D_v$.
 - ↳ The game continues with $v \cup \{(s, t)\}$ as new configuration.

Duplicator wins k -rounds, if last configuration is partial isomorphism.

Duplicator wins $G_k((S_v, \vec{s}), (S_w, \vec{t}))$ if has a winning strategy:
 whatever moves spoiler does, duplicator can win k -rounds.

Note:

Last configuration is partial isomorphism only if all intermediary configurations are partial isomorphisms.

How to check that duplicator wins $G_h((S_v, \vec{s}), (S_w, \vec{t}))$?

Lemma:

(1) Duplicator wins $G_0((S_v, \vec{s}), (S_w, \vec{t}))$

iff $\vec{s} \mapsto \vec{t}$ is partial isomorphism.

(2) Duplicator wins $G_{h+1}((S_v, \vec{s}), (S_w, \vec{t}))$ iff

(2a) $\forall s \in D_v : \exists t \in D_w : \text{Duplicator wins } G_h((S_v, s, \vec{s}), (S_w, t, \vec{t}))$

and

(2b) $\forall t \in D_w : \exists s \in D_v : \text{Duplicator wins } G_h((S_v, s, \vec{s}), (S_w, t, \vec{t}))$.

Intuition:

$G_h((S_v, s, \vec{s}), (S_w, t, \vec{t}))$ gives arbitrary first step in $G_{h+1}((S_v, \vec{s}), (S_w, \vec{t}))$.

Ehrenfeucht-Fraïssé Theorem (roughly) says that duplicator wins $G_h((S_v, \vec{s}), (S_w, \vec{t}))$ iff

v and w cannot be distinguished by FO[\exists]-formulas of quantifier depth $\leq h$.

Definition (Quantifier depth):

Let ℓ be a FO[\exists]-formula. The quantifier depth $q_d(\ell)$ is the nesting depth of quantifiers in ℓ .

It is defined inductively by

$$q_d(x < y) := 0 \quad q_d(P_n(x)) := 0$$

$$q_d(\ell_1 \vee \ell_2) := \max\{q_d(\ell_1), q_d(\ell_2)\} \quad q_d(\neg \ell) := q_d(\ell)$$

$$q_d(\exists x: \ell) := 1 + q_d(\ell).$$

Example:

Let $\ell = \forall x: P_n(x) \rightarrow \exists y: x < y \wedge P_n(x)$. Here, $q_d(\ell) = 2$.

Definition (k -Equivalence):

Let $v, w \in \Sigma^*$ and $\vec{s} = (s_1, \dots, s_n)$, $\vec{t} = (t_1, \dots, t_n)$ two vectors of positions.

Then $(Sv, \vec{s}), (Sw, \vec{t})$ are k -equivalent, $(Sv, \vec{s}) \equiv_{k,n} (Sw, \vec{t})$, if for all $\varphi(x_1, \dots, x_n)$ with $q_d(\varphi) \leq k$ we have

$$Sv, \llbracket \vec{s} / \vec{x} \rrbracket \models \varphi \text{ iff } Sw, \llbracket \vec{t} / \vec{x} \rrbracket \models \varphi.$$

Here, $\vec{x} = (x_1, \dots, x_n)$.

For the case of empty sequences, this means the two structures satisfy the same sentences of quantifier depth up to k .

Theorem (Ehrenfeucht - Fraïssé):

Duplicator wins $G_k((Sv, \vec{s}), (Sw, \vec{t}))$

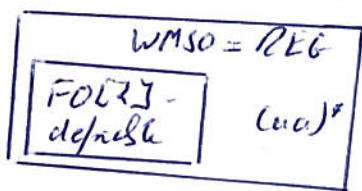
iff $(Sv, \vec{s}) \equiv_{k,n} (Sw, \vec{t})$.

Why is this cool?

Gives a "pumping argument".

Proposition:

Language $L = (aa)^*$ is not FO[\exists]-definable.



Proof:

Towards a contradiction, assume there was φ with $L = L(\varphi)$.

Then φ has some quantifier-depth $k \in \mathbb{N}$.

One can check that duplicator wins

$$G_k(a^{2k}, a^{2k+1})$$

Hence, the two words cannot be distinguished

by formula of q_d at most k (by Ehrenfeucht - Fraïssé Theorem)

Hence, a^{2^h} and $a^{2^{h+1}}$ cannot be distinguished by \mathcal{L} .

It is a contradiction to $a^{2^h} \in L(\mathcal{L})$ but $a^{2^{h+1}} \notin L(\mathcal{L})$.

Formula \mathcal{L} cannot exist.

□