

PONTIFICIA UNIVERSIDAD CATÓLICA DE CHILE

DEPARTAMENTO DE CIENCIA DE LA COMPUTACIÓN

APUNTE IIC2223

Teoría de Autómatas y Lenguajes Formales

Autor

Cristóbal ROJAS

En base a apuntes de

Prof. Cristian RIVEROS

3 de diciembre de 2022



Índice

1. Lenguajes regulares	3
1.1. Palabras y autómatas	3
1.1.1. Palabras	3
1.1.2. Autómatas	4
1.2. Construcciones de autómatas	6
1.2.1. Autómatas con función parcial de transición	6
1.2.2. Operaciones de conjuntos	7
1.3. No-determinismo	9
1.3.1. Definición de un NFA	9
1.3.2. Comparación con DFA	11
1.4. Expresiones regulares	13
1.5. Autómatas con transiciones sin lectura	15
1.5.1. e-NFA	15
1.5.2. NFA versus e-NFA	17
1.6. Teorema de Kleene	19
1.6.1. Desde Expresiones a Autómatas	19
1.6.2. Desde Autómatas a Expresiones	20
2. Propiedades de lenguajes regulares	23
2.1. Lema de bombeo	23
2.2. Minimización de autómatas	26
2.2.1. Colapsar estados	27
2.2.2. Algoritmo de minimización	28
2.3. Teorema de Myhill-Nerode	30
2.3.1. Relaciones de Myhill-Nerode	30
2.3.2. Camino al teorema	31
2.4. Autómatas en dos direcciones	33
2.4.1. Definición de un 2DFA	33
2.4.2. 2DFA vs DFA	34
3. Algoritmos para lenguajes regulares	37
3.1. Evaluación de expresiones regulares	37
3.2. Transductores	40
3.3. Análisis léxico	43
3.4. Algoritmo de Knuth-Morris-Prat	43
3.4.1. Autómata de un patrón	43
3.4.2. Autómata finito con k-lookahead	45
3.4.3. Algoritmo KMP	46
4. Lenguajes libres de contexto	48
4.1. Gramáticas libres de contexto	48
4.1.1. Gramáticas	48
4.1.2. Árboles y derivaciones	50
4.2. Simplificación de gramáticas	50
4.3. Forma normal de Chomsky	50
4.4. Lema de bombeo para lenguajes libres de contexto	50
4.5. Algoritmo CKY	50

5. Algoritmos para lenguajes libres de contexto	51
5.1. Autómatas apiladores	51
5.1.1. Versión normal	51
5.1.2. Versión alternativa	53
5.2. Autómatas apiladores vs gramáticas libres de contexto	55
5.2.1. Desde CFG a PDA	56
5.2.2. Desde PDA a CFG	57
5.3. Parsing: cómputo de First y Follow	59
5.3.1. Prefijos	61
5.3.2. First y Follow	62
5.3.3. Calcular First	63
5.3.4. Calcular Follow	64
5.4. Gramáticas LL	65
5.4.1. Definición Gramáticas LL	66
5.4.2. Caracterización LL	68
5.5. Parsing con gramáticas LL(k)	70
5.5.1. Algunas consideraciones	70
5.5.2. Parsing de LL(k)	73
6. Extracción de información	77
6.1. Extracción	77
6.2. Enumeración de resultados: Autómatas con anotaciones	77

1. Lenguajes regulares

1.1. Palabras y autómatas

1.1.1. Palabras

Definiciones. Consideremos que:

- ♦ Un **alfabeto** Σ es con conjunto finito.
- ♦ Un elemento de Σ lo llamaremos una **letra** o **símbolo**.
- ♦ Una **palabra** o **string** sobre Σ es una secuencia finita de letras en Σ .

Ejemplo 1.1

- $\Sigma = \{a, b, c\}$
- Palabras sobre Σ :

$$\text{aaaaabb, bcaabab, } a, \text{ bbbbb, } \dots$$

- ♦ El largo $|w|$ de una palabra w es el número de letras.

$$|w| \stackrel{\text{def}}{=} \# \text{ de letras en } w$$

- ♦ Denotaremos ϵ como la **palabra sin símbolos** de largo 0.

$$|\epsilon| \stackrel{\text{def}}{=} 0$$

- ♦ Denotaremos por Σ^* como el **conjunto de todas las palabras** sobre Σ .

Ejemplo 1.2

- Para $\Sigma = \{0, 1\}$:
- $|00011001| = 8$
 - $\Sigma^* = \text{todas las palabras posibles formadas por 0s y 1s.}$

Definición. Dados dos palabras $u, v \in \Sigma^*$ tal que $u = a_1 \dots a_n$ y $v = b_1 \dots b_m$:

$$u \cdot v \stackrel{\text{def}}{=} a_1 \dots a_n b_1 \dots b_m$$

Decimos que $u \cdot v$ es la palabra “ **u concatenada con v** ”.

Ejemplo 1.3

Para $\Sigma = \{0, 1, 2, \dots, 9\}$:

- ♦ $0123 \cdot 9938 = 01239938$ y $3493 \cdot \epsilon = 3493$

Propiedades de concatenación. La concatenación cumple:

- ♦ **Asociatividad:** $(u \cdot v) \cdot w = u \cdot (v \cdot w)$
- ♦ **Largo:** $|u \cdot v| = |u| + |v|$
- ♦ ¿Cumple conmutatividad? No. Por ejemplo, si $u = ab$ y $v = bb$, entonces $u \cdot v = abbb \neq bbab = v \cdot u$.

Definición. Sea Σ un alfabeto y $L \subseteq \Sigma^*$. Decimos que L es un **lenguaje** sobre el alfabeto Σ .

Ejemplo 1.4

Sea $\Sigma = \{a, b\}$:

- ♦ $L_0 = \{\epsilon, a, aa, b, aa\}$
- ♦ $L_1 = \{\epsilon, b, bb, bbb, bbbb, \dots\}$
- ♦ $L_2 = \{w \mid \exists u \in L_1, w = a \cdot u\}$
- ♦ $L_3 = \{w \mid \exists u, v \in \Sigma^*, w = u \cdot abba \cdot v\}$
- ♦ $L_4 = \{w \mid \exists u \in \Sigma^*, w = u \cdot u\}$

Un **lenguaje** puede ser visto como una **propiedad** de palabras.

Convenciones. Durante todo este texto:

- ♦ Para **letras** se usarán los símbolos: a, b, c, d, e, \dots
- ♦ Para **palabras** se usarán los símbolos: w, u, v, x, y, z, \dots
- ♦ Para **alfabetos** se usarán los símbolos: Σ, Γ, \dots
- ♦ Para **lenguajes** se usarán los símbolos: L, M, N, \dots
- ♦ Para **números** se usarán los símbolos: i, j, k, l, m, n, \dots

1.1.2. Autómatas

Una autómata **finito** es:

- ♦ Un modelo de computación sencillo, basado en una cantidad **finita** de memoria.
- ♦ Procesa cada **palabra** de principio a fin en **una sola pasada**.
- ♦ Al terminar, el autómata decide si **acepta** o **rechaza** el input.

Usaremos los autómatas finitos para definir **lenguajes**.

Definición. Un autómata finito determinista (DFA) es una tupla:

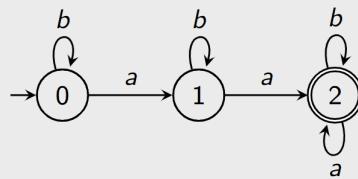
$$\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$$

- ♦ Q es un conjunto finito de **estados**.
- ♦ Σ es el alfabeto del **input**.
- ♦ $\delta : Q \times \Sigma \rightarrow Q$ es la función de **transición**.
- ♦ $q_0 \in Q$ es el **estado inicial**.
- ♦ $F \subseteq Q$ es el conjunto de **estados finales** (o **aceptación**).

Ejemplo 1.5

- ♦ $Q = \{0, 1, 2\}$
- ♦ $\Sigma = \{a, b\}$
- ♦ $\delta : Q \times \Sigma \rightarrow Q$ se define como:

$$\begin{aligned}\delta(0, a) &= 1 \\ \delta(1, a) &= 2 \\ \delta(2, a) &= 2 \\ \delta(q, b) &= q \quad \forall q \in \{0, 1, 2\}\end{aligned}$$
- ♦ $q_0 = 0$
- ♦ $F = \{2\}$

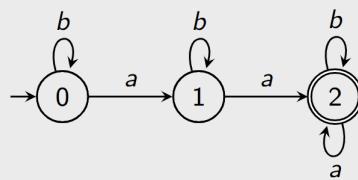


Ejecución. Sea $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ un DFA y $w = a_1a_2\dots a_n \in \Sigma^*$ un input. Una **ejecución** (o *run*) ρ de \mathcal{A} sobre w es una secuencia:

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} p_2 \xrightarrow{a_3} \dots \xrightarrow{a_n} p_n$$

donde $p_0 = q_0$ y para todo $i \in \{0, 1, \dots, n - 1\}$ se cumple que $\delta(p_i, a_{i+1}) = p_{i+1}$.

Una ejecución ρ de \mathcal{A} sobre w es de **aceptación** si $p_n \in F$.

Ejemplo 1.6

- ♦ ¿Cuál es la ejecución de \mathcal{A} sobre $bbab$?

$\rho : 0 \xrightarrow{b} 0 \xrightarrow{b} 0 \xrightarrow{a} 1 \xrightarrow{b} 1$. La ejecución **no** es de aceptación ya que no termina en un estado final.

- ♦ ¿Cuál es la ejecución de \mathcal{A} sobre $abab$?

$\rho : 0 \xrightarrow{a} 1 \xrightarrow{b} 1 \xrightarrow{a} 2 \xrightarrow{b} 2$. La ejecución **sí** es de aceptación ya que termina en un estado final.

Aceptación. Sea $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ un DFA y $w \in \Sigma^*$. Decimos que \mathcal{A} **acepta** w si la ejecución de \mathcal{A} sobre w es de aceptación. Al contrario, decimos que \mathcal{A} **rechaza** w si la ejecución de \mathcal{A} sobre w NO es de aceptación.

El **lenguaje aceptado** por \mathcal{A} se define como:

$$\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ acepta } w\}$$

Un lenguaje $L \subseteq \Sigma^*$ se dice **regular** si, y sólo si, **existe** un autómata finito determinista \mathcal{A} tal que

$$L = \mathcal{L}(\mathcal{A})$$

1.2. Construcciones de autómatas

Veremos una definición alternativa al autómata visto en la sección anterior.

1.2.1. Autómatas con función parcial de transición

Definición. Un autómata finito determinista con **función parcial de transición** (DFAp) es una tupla:

$$\mathcal{A} = (Q, \Sigma, \gamma, q_0, F)$$

- ♦ Q es un conjunto finito de estados.
- ♦ Σ es el alfabeto del input.
- ♦ $\gamma : Q \times \Sigma \rightharpoonup Q$ es una **función parcial de transición**.
- ♦ $q_0 \in Q$ es el estado inicial.
- ♦ $F \subseteq Q$ es el conjunto de estados finales (o aceptación).

Ejecución. Sea $w = a_1a_2 \dots a_n \in \Sigma^*$. De igual manera que un DFA, una **ejecución** (*o run*) ρ de \mathcal{A} sobre w es una secuencia:

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} p_2 \dots \xrightarrow{a_n} p_n$$

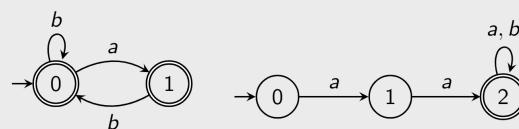
donde $p_0 = q_0$ y para todo $i \in \{0, \dots, n-1\}$ **está definido** $\gamma(p_i, a_{i+1}) = p_{i+1}$.

Una ejecución ρ de \mathcal{A} sobre w es de **aceptación** si $p_n \in F$. Notemos que ahora una palabra puede NO tener una ejecución.

Aceptación. Sea \mathcal{A} un DFAp y $w \in \Sigma^*$. Decimos que \mathcal{A} **acepta** w si **existe una ejecución** de \mathcal{A} sobre w que es de aceptación. También, el **lenguaje aceptado** por \mathcal{A} se define como:

$$\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ acepta } w\}$$

Ejemplo 1.7



¿DFA \neq DFAp? Establezcamos una proposición. Para todo autómata \mathcal{A} con función parcial de transición, existe un autómata \mathcal{A}' (con función total de transición) tal que:

$$\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}')$$

En otras palabras, $\text{DFA} \equiv \text{DFAp}$.

Demostración. Sea $\mathcal{A} = (Q, \Sigma, \gamma, q_0, F)$ un autómata con función parcial de transición. Sea q_s un **nuevo estado** tal que $q_s \notin Q$. Construimos el DFA $\mathcal{A}' = (Q \cup \{q_s\}, \Sigma, \delta', q_0, F)$ tal que:

$$\delta'(p, a) = \begin{cases} \gamma(p, a) & \text{si } p \neq q_s \text{ y } (p, a) \in \text{dom}(\gamma) \\ q_s & \text{si no} \end{cases}$$

para todo $p \in Q \cup \{q_s\}$ y $a \in \Sigma$. Queremos demostrar que $\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}')$.

Dem. $\mathcal{L}(\mathcal{A}) \subseteq \mathcal{L}(\mathcal{A}')$. Sea $w = a_1 \dots a_n \in \mathcal{L}(\mathcal{A})$. Entonces, existe una ejecución de aceptación ρ de \mathcal{A} sobre w :

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} p_2 \dots \xrightarrow{a_n} p_n$$

donde $p_0 = q_0$, para todo $i \in \{0, \dots, n-1\}$ está definido $\gamma(p_i, a_{i+1}) = p_{i+1}$ y $p_n \in F$. Como $\delta'(p_i, a_{i+1}) = \gamma(p_i, a_{i+1})$ para todo $i \in \{0, \dots, n-1\}$ (por la definición de δ'), entonces ρ es también una ejecución de aceptación de \mathcal{A}' sobre w . Por lo tanto, $w \in \mathcal{L}(\mathcal{A}')$.

Dem. $\mathcal{L}(\mathcal{A}') \subseteq \mathcal{L}(\mathcal{A})$. Sea $w = a_1 \dots a_n \in \mathcal{L}(\mathcal{A}')$. Entonces, existe una ejecución de aceptación ρ de \mathcal{A}' sobre w :

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} p_2 \dots \xrightarrow{a_n} p_n$$

donde $p_0 = q_0$, para todo $i \in \{0, \dots, n-1\}$ tenemos $\gamma'(p_i, a_{i+1}) = p_{i+1}$ y $p_n \in F$. Demostraremos que $p_i \neq q_s$ para todo $i \in \{0, \dots, n\}$. Por **contradicción**, suponga que existe i tal que $p_i = q_s$, entonces, tenemos que $p_{i+1} = q_s$. Por **inducción**, podemos demostrar que $p_j = q_s$ para todo $j \geq i$, y así, podemos concluir que $p_n = q_s$, llevándonos a una contradicción. Como $p_i \neq q_s$ para todo $i \in \{0, \dots, n\}$, tenemos que:

$$\delta'(p_i, a_{i+1}) = \gamma(p_i, a_{i+1}) \quad \forall i \in \{0, 1, \dots, n-1\}$$

y entonces ρ es una **ejecución de aceptación** de \mathcal{A} sobre w . Por lo tanto, concluimos que $w \in \mathcal{L}(\mathcal{A})$. ■

Advertencia. Desde ahora, se utilizarán autómatas con funciones **totales** de transición, pero sin pérdida de generalidad, en algunos ejemplos habrán autómatas con funciones **parciales** de transición por simplicidad.

1.2.2. Operaciones de conjuntos

Definiciones. Dado dos lenguajes $L, L' \subseteq \Sigma^*$ se define:

$$\begin{aligned} L^C &= \{w \in \Sigma^* \mid w \notin L\} \\ L \cap L' &= \{w \in \Sigma^* \mid w \in L \wedge w \in L'\} \\ L \cup L' &= \{w \in \Sigma^* \mid w \in L \vee w \in L'\} \end{aligned}$$

Dado dos autómatas \mathcal{A} y \mathcal{A}' :

1. ¿Existe un autómata \mathcal{B} tal que $\mathcal{L}(\mathcal{B}) = \mathcal{L}(\mathcal{A})^C$?
2. ¿Existe un autómata \mathcal{B} tal que $\mathcal{L}(\mathcal{B}) = \mathcal{L}(\mathcal{A}) \cap \mathcal{L}(\mathcal{A}')$?
3. ¿Existe un autómata \mathcal{B} tal que $\mathcal{L}(\mathcal{B}) = \mathcal{L}(\mathcal{A}) \cup \mathcal{L}(\mathcal{A}')$?

Construcción de $\mathcal{L}(\mathcal{A})^C$. Dado un autómata $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$, definimos el autómata:

$$\mathcal{A}^C = (Q, \Sigma, \delta, q_0, Q \setminus F)$$

Teorema 1

Para todo autómata \mathcal{A} , se tiene que $\mathcal{L}(\mathcal{A})^C = \mathcal{L}(\mathcal{A}^C)$.

Producto de autómatas. Suponga que:

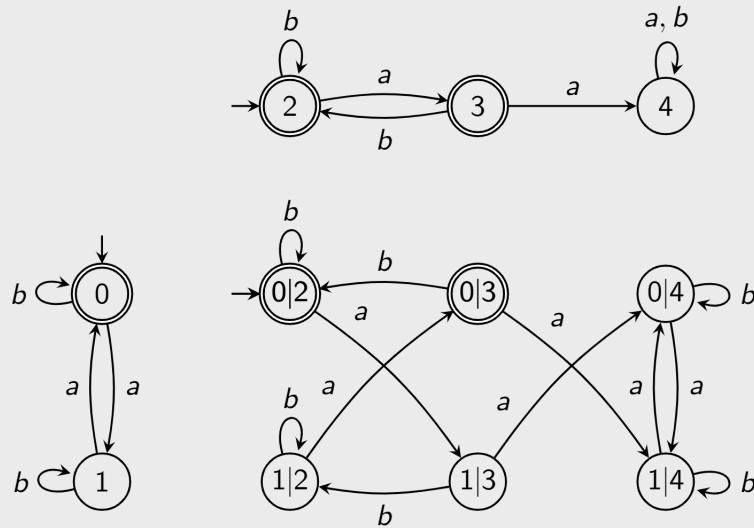
$$\begin{aligned}\mathcal{A} &= (Q, \Sigma, \delta, q_0, F) \\ \mathcal{A}' &= (Q', \Sigma, \delta', q'_0, F')\end{aligned}$$

y considere una palabra $w \in \Sigma^*$. ¿Cómo ejecutamos ambos autómatas sobre w al **mismo tiempo**? La idea es ejecutar \mathcal{A} y \mathcal{A}' en paralelo. Así, definimos el **producto** entre \mathcal{A} y \mathcal{A}' como el autómata $\mathcal{A} \times \mathcal{A}' = (Q^\times, \Sigma, \delta^\times, q_0^\times, F^\times)$ tal que:

- ♦ $Q^\times = Q \times Q' = \{(q, q') \mid q \in Q \wedge q' \in Q'\}$
- ♦ $\delta^\times((q, q'), a) = (\delta(q, a), \delta'(q', a))$
- ♦ $q_0^\times = (q_0, q'_0)$
- ♦ $F^\times = F \times F'$

Ejemplo 1.8

Todas las palabras sobre $\{a, b\}$ con una cantidad par de a -letras tal que no hay dos a -letras seguidas.



Teorema 2

Para todo par de autómatas \mathcal{A} y \mathcal{A}' se tiene que

$$\mathcal{L}(\mathcal{A} \times \mathcal{A}') = \mathcal{L}(\mathcal{A}) \cap \mathcal{L}(\mathcal{A}')$$

Demostración teorema 2. Solo se demostrará que $\mathcal{L}(\mathcal{A}) \cap \mathcal{L}(\mathcal{A}') \subseteq \mathcal{L}(\mathcal{A} \times \mathcal{A}')$, la otra dirección queda propuesta para el lector.

Sea $w = a_1 \dots a_n \in \mathcal{L}(\mathcal{A}) \cap \mathcal{L}(\mathcal{A}')$. Entonces $w \in \mathcal{L}(\mathcal{A})$ y $w \in \mathcal{L}(\mathcal{A}')$. Existen ejecuciones de aceptación ρ y ρ' de \mathcal{A} y \mathcal{A}' sobre w , respectivamente:

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} p_n \quad \rho' : p'_0 \xrightarrow{a_1} p'_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} p'_n$$

- ♦ $p_0 = q_0$ y $p'_0 = q'_0$.
- ♦ $\delta(p_{i-1}, a_i) = p_i$ y $\delta'(p'_{i-1}, a_i) = p'_i$ para todo $i \in \{1, \dots, n\}$.
- ♦ $p_n \in F$ y $p'_n \in F'$.

Por definición, tenemos que: $\rho^\times : (p_0, p'_0) \xrightarrow{a_1} (p_1, p'_1) \xrightarrow{a_2} \dots \xrightarrow{a_n} (p_n, p'_n)$

- ♦ $(p_0, p'_0) = (q_0, q'_0)$.
- ♦ $(p_i, p'_i) = (\delta(p_{i-1}, a_i), \delta'(p'_{i-1}, a_i)) = \delta^\times((p_{i-1}, p'_{i-1}), a_i) \forall i \in \{1, \dots, n\}$.
- ♦ $(p_n, p'_n) \in F \times F'$.

Por lo tanto, ρ^\times es una ejecución de $\mathcal{A} \times \mathcal{A}'$ sobre w y $w \in \mathcal{L}(\mathcal{A} \times \mathcal{A}')$. ■

Unión de autómatas. Sabemos que

$$\mathcal{L}(\mathcal{A}) \cup \mathcal{L}(\mathcal{A}') = \left(\mathcal{L}(\mathcal{A})^C \cap \mathcal{L}(\mathcal{A}')^C \right)^C$$

Para calcular el autómata que acepta el lenguaje $\mathcal{L}(\mathcal{A}) \cup \mathcal{L}(\mathcal{A}')$:

1. Complementamos \mathcal{A} y \mathcal{A}' .
2. Intersectamos \mathcal{A}^C y $(\mathcal{A}')^C$.
3. Complementamos $\mathcal{A}^C \times (\mathcal{A}')^C$.

1.3. No-determinismo

*"Indeterminism is the concept that events (certain events, or events of certain types) are not caused deterministically (cf. causality) by prior events. It is the opposite of **determinism** and related to chance. It is highly relevant to the philosophical problem of **free will**."* - Wikipedia.

1.3.1. Definición de un NFA

Definición. Un autómata finito **no-determinista** (NFA) es una estructura:

$$\boxed{\mathcal{A} = (Q, \Sigma, \Delta, I, F)}$$

- ♦ Q es un conjunto finito de estados.
- ♦ Σ es el alfabeto del input.
- ♦ $F \subseteq Q$ es el conjunto de estados finales (o aceptación).
- ♦ $\Delta \subseteq Q \times \Sigma \times Q$ es la **relación de transición**.
- ♦ $I \subseteq Q$ es un **conjunto de estados iniciales**.

Ejemplo 1.9

♦ $Q = \{0, 1, 2\}$, $\Sigma = \{a, b\}$, $I = \{0, 1\}$, $F = \{2\}$

♦ $\Delta \subseteq Q \times \Sigma \times Q$ se define como:

$$(0, a, 0) \in \Delta$$

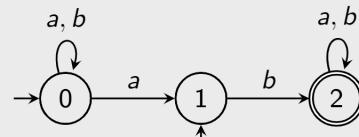
$$(0, a, 1) \in \Delta$$

$$(0, b, 0) \in \Delta$$

$$(1, b, 2) \in \Delta$$

$$(2, a, 2) \in \Delta$$

$$(2, b, 2) \in \Delta$$



Ejecución. Sea $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ un NFA y $w = a_1 a_2 \dots a_n \in \Sigma^*$ el input. Una **ejecución** (o *run*) ρ de \mathcal{A} sobre w es una secuencia:

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} p_n$$

donde $p_0 \in I$ y para todo $i \in \{0, \dots, n-1\}$, se tiene que $(p_i, a_{i+1}, p_{i+1}) \in \Delta$.

Una ejecución ρ de \mathcal{A} sobre w es de **aceptación** si $p_n \in F$.

Aceptación. Decimos que \mathcal{A} **acepta** w si existe una ejecución de \mathcal{A} sobre w que es de aceptación. Por otro lado, decimos que \mathcal{A} **rechaza** si todas las ejecuciones de \mathcal{A} sobre w NO son de aceptación. Además, el lenguaje **aceptado** por \mathcal{A} se define como:

$$\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ acepta } w\}$$

Interpretación. Las siguientes interpretaciones pueden ayudar a entender mejor un NFA:

1. $\Delta \subseteq Q \times \Sigma \times Q$ es la **relación de transición**.

“(q, a, p) ∈ Δ entonces existe una transición desde q a p al leer a”.

2. $I \subseteq Q$ es un **conjunto de estados iniciales**.

“p ∈ I entonces p es un posible estado inicial del autómata”

1'. $\Delta : Q \times \Sigma \rightarrow 2^Q$ es una **función de transición**.

“q ∈ Δ(p, a) entonces q es un posible estado que puedo llegar desde p al leer a”.

Esta interpretación es más común encontrarla en libros sobre teoría de autómatas.

Además, el **no-determinismo** puede ser visto como:

1. Paralelización infinita, es decir, cada ejecución es un *thread* distinto.

2. “Guessing and Verifying” (adivinar y verificar).

El no-determinismo NO debe ser visto como:

♦ Explicitamente como el **indeterminismo** o “libre albedrío”. Para un input, un NFA siempre produce el mismo resultado.

♦ Comportamiento **aleatorio** del autómata.

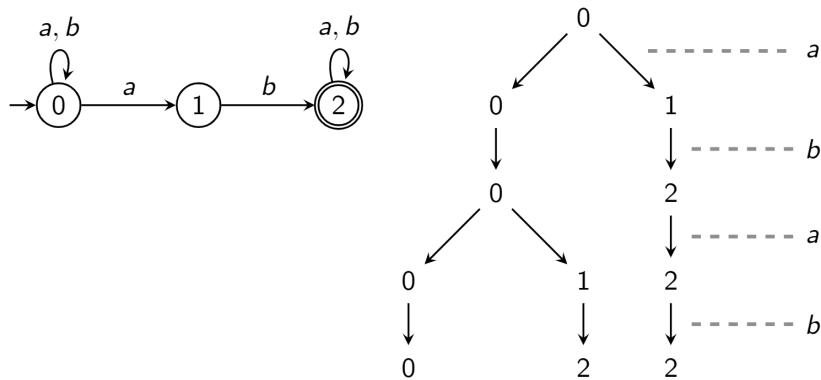


Figura 1: Interpretación del no-determinismo

1.3.2. Comparación con DFA

A continuación, veremos que autómata finito determinista (DFA) puede almacenar **todas** las ejecuciones de un NFA. A este proceso se le conoce como **determinización**.

Teorema 3

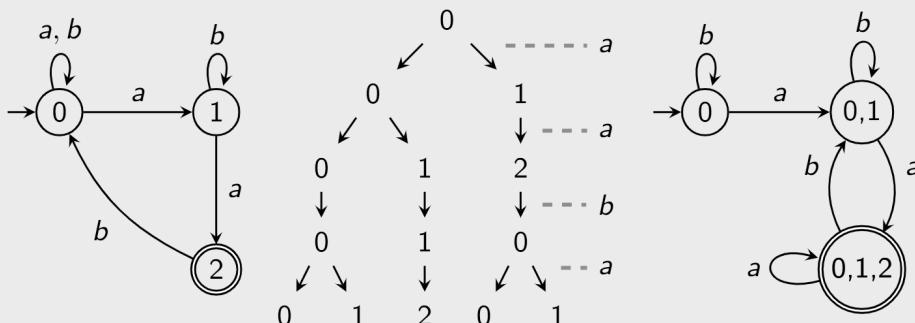
Para todo autómata finito no-determinista \mathcal{A} , existe un autómata determinista \mathcal{A}' tal que

$$\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}')$$

En otras palabras, $DFA \equiv NFA$.

Idea. Primero, pensemos en la idea de determinización: “almacenar en el autómata determinista todos los estados actuales de las ejecuciones en curso (sin repetidos)”.

Ejemplo 1.10



Formalización. Para un autómata no-determinista $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$, definimos el autómata determinista (**determinización** de \mathcal{A}):

$$\mathcal{A}^{\text{det}} = (2^Q, \Sigma, \delta^{\text{det}}, q_0^{\text{det}}, F^{\text{det}})$$

- ♦ $2^Q = \{S \mid S \subseteq Q\}$ es el conjunto potencia de Q .

- ♦ $q_0^{\text{det}} = I$.

- ♦ $\delta^{\text{det}} : 2^Q \times \Sigma \rightarrow 2^Q$ tal que:

$$\delta^{\text{det}}(S, a) = \{q \in Q \mid \exists p \in S. (p, a, q) \in \Delta\}$$

- ♦ $F^{\text{det}} = \{S \in 2^Q \mid S \cap F \neq \emptyset\}$, es decir, todos los conjuntos que tengan al menos un estado final.

Demostración teorema 3. La determinización puede verse como un **subset construction**. Partamos con $\mathcal{L}(\mathcal{A}) \subseteq \mathcal{L}(\mathcal{A}^{\text{det}})$.

Sea $w = a_1 \dots a_n \in \mathcal{L}(\mathcal{A})$. Existe una ejecución ρ de \mathcal{A} sobre w :

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} p_n$$

donde $p_0 = I$, $(p_i, a_{i+1}, p_{i+1}) \in \Delta$ para todo $i \in \{0, \dots, n-1\}$ y $p_n \in F$.

Como \mathcal{A}^{det} es determinista, entonces existe una ejecución ρ' de \mathcal{A}^{det} sobre w :

$$\rho' : S_0 \xrightarrow{a_1} S_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} S_n$$

donde $S_0 = I$ y $\delta^{\text{det}}(S_i, a_{i+1}) = S_{i+1}$ para todo $i \in \{0, \dots, n-1\}$. Luego, queremos demostrar que $p_i \in S_i$ para todo $i \in \{0, \dots, n-1\}$.

Por **inducción** sobre i , tenemos que:

- ♦ **Caso base:** $p_0 \in S_0$ por definición de \mathcal{A}^{det} .

- ♦ **Inducción:** Suponemos que $p_i \in S_i$ y demostramos para $i+1$. Como sabemos que:

- $\delta^{\text{det}}(S_i, a_{i+1}) = S_{i+1} = \{q \in Q \mid \exists p \in S_i. (p, a, q) \in \Delta\}$ y
- $(p_i, a_{i+1}, p_{i+1}) \in \Delta$

Entonces $p_{i+1} \in S_{i+1}$, ya que, si estamos en p_i leyendo a_{i+1} , la transición nos dice que pasaremos al estado p_{i+1} que pertenece a S_{i+1} por la hipótesis de inducción.

Luego, como $p_n \in S_n$, entonces $S_n \cap F \neq \emptyset$ y así $S_n \in F^{\text{det}}$. Por lo tanto, $w \in \mathcal{L}(\mathcal{A}^{\text{det}})$.

Ahora, demostramos la otra dirección: $\mathcal{L}(\mathcal{A}^{\text{det}}) \subseteq \mathcal{L}(\mathcal{A})$.

Sea $w = a_1 \dots a_n \in \mathcal{L}(\mathcal{A}^{\text{det}})$. Existe una ejecución ρ de \mathcal{A}^{det} sobre w :

$$\rho : S_0 \xrightarrow{a_1} S_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} S_n$$

donde $S_0 = I$, $\delta^{\text{det}}(S_i, a_{i+1}) = S_{i+1}$ para todo $i \in \{0, \dots, n-1\}$ y $S_n \in F^{\text{det}}$, con $S_n \cap F \neq \emptyset$. Buscamos demostrar entonces que para todo $i \leq n$ y para todo $p \in S_i$, existe una ejecución:

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} \dots \xrightarrow{a_i} p_i = p$$

tal que:

1. $p_0 \in I$.

2. $(p_j, a_{j+1}, p_{j+1}) \in \Delta$ para todo $j \in \{0, \dots, i-1\}$.

Por **inducción** sobre i , tenemos que:

- ♦ **Caso base:** Si $p \in S_0 = I$, entonces la ejecución $\rho : p$ cumple 1. y 2.

- ♦ **Inducción:** Supongamos que se cumple para todo $p \in S_i$. Sea $q \in S_{i+1}$. Como $\delta^{\text{det}}(S_i, a_{i+1}) = S_{i+1} = \{q \in Q \mid \exists p \in S_i. (p, a, q) \in \Delta\}$ y $q \in S_{i+1}$, entonces existe $p \in S_i$ tal que $(p, a_{i+1}, q) \in \Delta$.

Por hipótesis de inducción, existe $\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} \dots \xrightarrow{a_i} p_i = p$ que satisface 1. y 2.

Por lo tanto, $\rho' : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} \dots \xrightarrow{a_i} p_i \xrightarrow{a_{i+1}} q$ también satisface 1. y 2.

Como lo anterior queda demostrado y como $S_n \cap F \neq \emptyset$, para $p \in S_n \cap F$ existe una ejecución de aceptación de \mathcal{A} sobre w . Por lo tanto, $w \in \mathcal{L}(\mathcal{A})$. ■

1.4. Expresiones regulares

Definición. R es una **expresión regular** sobre Σ si R es igual a:

1. a , para alguna letra $a \in \Sigma$.
2. ϵ
3. \emptyset
4. $(R_1 + R_2)$, donde R_1 y R_2 son expresiones regulares.
5. $(R_1 \cdot R_2)$, donde R_1 y R_2 son expresiones regulares.
6. (R_1^*) , donde R_1 es una expresión regular. Esta expresión se conoce como **clausura de Kleene**.

Denotaremos como ExpReg el conjunto de **todas las expresiones regulares** sobre Σ .

Ejemplo 1.11

Las siguientes son expresiones regulares sobre $\Sigma = \{a, b\}$:

- ♦ $(a + b)$
- ♦ $((a \cdot b) \cdot c)$
- ♦ (a^*)
- ♦ $(b \cdot (a^*))$
- ♦ $((a + b)^*)$
- ♦ $((a \cdot ((b \cdot a)^*)) + \epsilon)$
- ♦ $((a \cdot ((b \cdot a)^*)) + \emptyset)$

Para reducir la cantidad de paréntesis, se define el **orden de precedencia**:

1. estrella $(\cdot)^*$
2. concanetación \cdot
3. unión $+$

Semántica. Para una expresión regular R cualquiera, se define el lenguaje $\mathcal{L}(R) \subseteq \Sigma^*$ **inductivamente** como:

1. $\mathcal{L}(a) = \{a\}$, para toda letra $a \in \Sigma$.
2. $\mathcal{L}(\epsilon) = \{\epsilon\}$.
3. $\mathcal{L}(\emptyset) = \emptyset$.
4. $\mathcal{L}(R_1 + R_2) = \mathcal{L}(R_1) \cup \mathcal{L}(R_2)$, donde R_1 y R_2 son expresiones regulares.
5. $\mathcal{L}(R_1 \cdot R_2) = \mathcal{L}(R_1) \cdot \mathcal{L}(R_2)$, donde R_1 y R_2 son expresiones regulares.
6. $\mathcal{L}(R_1^*) = \bigcup_{k=0}^{\infty} \mathcal{L}(R_1)^k$, donde R_1 es una expresión regular.

Para el punto 5. y 6., definimos para dos lenguajes $L_1, L_2 \subseteq \Sigma^*$ el **producto** de L_1 y L_2 :

$$L_1 \cdot L_2 = \{w_1 \cdot w_2 \mid w_1 \in L_1 \wedge w_2 \in L_2\}$$

Además, para un lenguaje $L \subseteq \Sigma^*$ se define la **potencia** a la $n \geq 0$:

$$L^n = \{w_1 \cdot w_2 \dots w_n \mid \forall i \leq n. w_i \in L\}$$

La **potencia** a la 0 se define como $L^0 = \{\epsilon\}$.

Ejemplo 1.12

Se muestran a continuación lenguajes definidos por algunas ExpReg:

- ♦ $\mathcal{L}((a + b)^*) = \{a, b\}^*$
- ♦ $\mathcal{L}((a \cdot b) \cdot (b \cdot a)) = \{abba\}$
- ♦ $\mathcal{L}(a \cdot (b \cdot a) + b \cdot a + (a \cdot b) \cdot a) = \{aba, ba\}$

Definición. Decimos que R_1 es **equivalente** a R_2 si, y sólo si, $\mathcal{L}(R_1) = \mathcal{L}(R_2)$. Si R_1 es equivalente a R_2 , escribiremos $R_1 \equiv R_2$.

Lema. Los operadores de unión + y producto · son **asociativos**.

$$\begin{aligned} (R_1 + R_2) + R_3 &\equiv R_1 + (R_2 + R_3) \\ (R_1 \cdot R_2) \cdot R_3 &\equiv R_1 \cdot (R_2 \cdot R_3) \end{aligned}$$

La demostración de este lema queda como ejercicio propuesto al lector.

Ejemplo 1.13

Más lenguajes definidos por algunas ExpReg:

- ♦ $\mathcal{L}(a^* \cdot b \cdot a^*) =$ todas las palabras con una sola b .
- ♦ $\mathcal{L}((a + b)^* \cdot b \cdot (a + b)^*) =$ todas las palabras con una o más $b's$.

Definición. Usamos las siguientes **abreviaciones** de expresiones regulares:

$$\begin{aligned} R^+ &\equiv R \cdot R^* \\ R^k &\equiv R \cdot \underset{k}{\dots} \cdot R \\ R^? &\equiv R + \epsilon \\ \Sigma &\equiv a_1 + \dots + a_n \end{aligned}$$

para $R \in \text{ExpRegs}$ y $\Sigma = \{a_1, \dots, a_n\}$.

Ejemplo 1.14

Más lenguajes definidos por algunas ExpReg:

- ♦ $\mathcal{L}(\Sigma^* \cdot b \cdot \Sigma^*) =$ todas las palabras con una sola b .
- ♦ $\mathcal{L}(b^* \cdot (a \cdot b)^5) =$ todas las palabras con 5 a 's.
- ♦ $\mathcal{L}(a^* \cdot (b + c)^?) =$ todas las palabras de a 's y terminadas en b o c .
- ♦ $\mathcal{L}((a \cdot b^+)^+) =$ todas las palabras que empiezan con a y donde cada a esta seguida de al menos una b .

1.5. Autómatas con transiciones sin lectura

Hasta ahora, hemos visto lo siguiente:

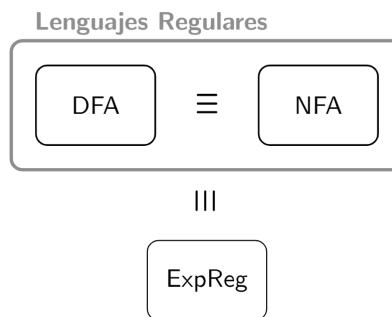


Figura 2: Mapa actual de nuestros modelos de computación

Podemos demostrar que $\text{ExpReg} \subseteq \text{NFA}$, pero para eso necesitamos un nuevo modelo.

1.5.1. ϵ -NFA

Lo nuevo de este autómata:

1. tiene transiciones no deterministas y
2. tiene transiciones leyendo la palabra vacía ϵ :

$$p \xrightarrow{\epsilon} q$$

La importancia de un ϵ -NFA es que es un modelo **muy útil** para construir nuevos autómatas y NO agrega más poder de computación a los NFA.

Definición. Un autómata finito no-determinista con ϵ -transiciones (ϵ -NFA) es una tupla:

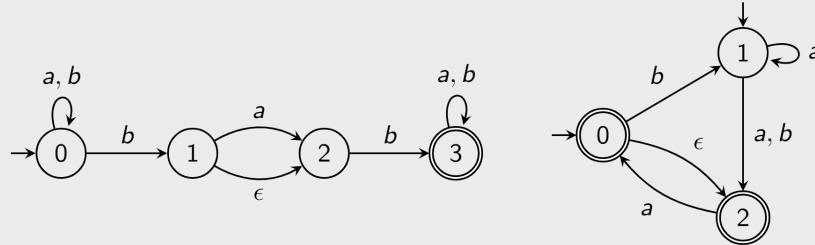
$$\mathcal{A} = (Q, \Sigma, \Delta, I, F)$$

- ♦ Q es un conjunto finito de estados.
- ♦ Σ es el alfabeto del input.
- ♦ $I \subseteq Q$ es un conjunto de estados iniciales.

- ♦ $F \subseteq Q$ es el conjunto de estados finales (o aceptación).
- ♦ $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$ es la relación de transición.

Ejemplo 1.15

Algunos ejemplos de ϵ -NFA's:



Para ϵ -NFA veremos una **forma alternativa** para definir las nociones de ejecución y aceptación.

Ejecución. Sea $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ un ϵ -NFA. Definimos:

- ♦ Un par $(q, w) \in Q \times \Sigma^*$ es una **configuración** de \mathcal{A} .
- ♦ Una configuración (q, w) es **inicial** si $q \in I$.
- ♦ Una configuración (q, w) es **final** si $q \in F$ y $w = \epsilon$.

"Intuitivamente, una configuración (q, aw) representa que \mathcal{A} se encuentra en el estado q procesando la palabra aw y leyendo a ".

- ♦ Se define la relación $\vdash_{\mathcal{A}} \subseteq (Q \times \Sigma^*) \times (Q \times \Sigma^*)$ de **siguiente-paso** entre configuraciones de \mathcal{A} :

$$(p, u) \vdash_{\mathcal{A}} (q, v)$$

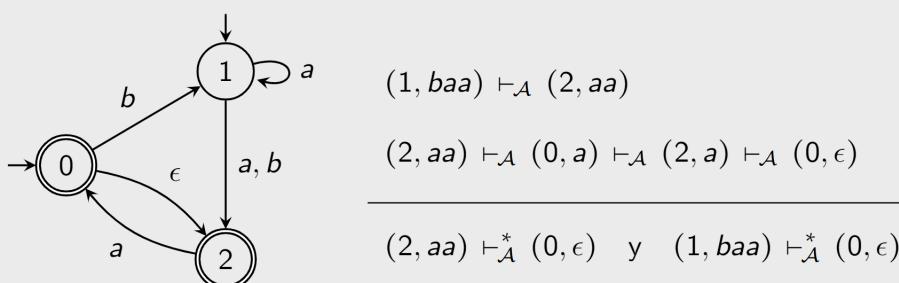
si, y sólo si, existe $(p, c, q) \in \Delta$ con $c \in \Sigma \cup \{\epsilon\}$ tal que $u = c \cdot v$.

- ♦ Se define $\vdash_{\mathcal{A}}^*$ como la clausura **refleja** y **transitiva** de $\vdash_{\mathcal{A}}$:

para toda configuración $(q, w) :$	$(q, w) \vdash_{\mathcal{A}}^* (q, w)$
si $(p, u) \vdash_{\mathcal{A}} (p', w)$ y $(p', w) \vdash_{\mathcal{A}}^* (q, v) :$	$(p, u) \vdash_{\mathcal{A}}^* (q, v)$

Decimos que $(p, u) \vdash_{\mathcal{A}}^* (q, v)$ si uno puede ir de (p, u) a (q, v) en **0 o más pasos**.

Ejemplo 1.16



Aceptación. Sea $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ un ϵ -NFA y $w \in \Sigma^*$. Decimos que \mathcal{A} **acepta** w si existe una configuración **inicial** (q_0, w) y una configuración **final** (q_f, ϵ) tal que:

$$(q_0, w) \vdash_{\mathcal{A}}^* (q_f, \epsilon)$$

Además, el **lenguaje aceptado** por \mathcal{A} se define como:

$$\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ acepta } w\}$$

Nota. Si \mathcal{A} no tiene ϵ -transiciones o no-determinismo, esta es una forma **alternativa** para definir ejecución y aceptación para NFA y DFA.

1.5.2. NFA versus ϵ -NFA

Partimos enunciando el siguiente teorema:

Teorema 4

Para todo autómata finito no-determinista con ϵ -transiciones \mathcal{A} , existe un autómata no-determinista \mathcal{A}' tal que:

$$\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}')$$

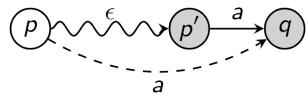
En otras palabras, NFA $\equiv \epsilon$ -NFA.

Para demostrar este teorema, mostraremos como construir un autómata no-determinista a partir de un ϵ -NFA removiendo las ϵ -transiciones.

Construcción desde ϵ -NFA a NFA. Dado un ϵ -NFA $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ se define el NFA:

$$\mathcal{A}' = (Q, \Sigma, \Delta', I, F')$$

- ♦ para todo $p, q \in Q$, $(p, a, q) \in \Delta'$ si, y sólo si, $\exists p' \in Q$ tal que:
 - $(p, \epsilon) \vdash_{\mathcal{A}}^* (p', \epsilon)$ y
 - $(p', a, q) \in \Delta$.
- ♦ $F' = \{p \in Q \mid \exists q \in F. (p, \epsilon) \vdash_{\mathcal{A}}^* (q, \epsilon)\}$



Por definición, si $(p, a, q) \in \Delta$, entonces $(p, a, q) \in \Delta'$ para todo $a \in \Sigma$.

Teorema 5

Dado un ϵ -NFA $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ se tiene que:

$$\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}')$$

Demostración teorema 5. Demostrar el teorema anterior es equivalente a demostrar que para todo $p \in Q$ y $w \in \Sigma^*$:

$$\exists q \in F. (p, w) \vdash_{\mathcal{A}}^* (q, \epsilon) \quad \text{si, y sólo si,} \quad \exists q' \in F'. (p, w) \vdash_{\mathcal{A}'}^* (q', \epsilon)$$

De aquí, podemos **concluir** que \mathcal{A} acepta w si, y sólo si, \mathcal{A}' acepta w .

Por **inducción** sobre el largo de w :

♦ **Caso base:** Para $w = \epsilon$:

(\Rightarrow) Sea $q \in F$ tal que $(p, \epsilon) \vdash_{\mathcal{A}}^* (q, \epsilon)$. Por definición de $F^{\not\epsilon}$, se tiene que $p \in F^{\not\epsilon}$. Por lo tanto, $(p, \epsilon) \vdash_{\mathcal{A}^{\not\epsilon}}^* (p, \epsilon)$.

(\Leftarrow) Sea $q' \in F^{\not\epsilon}$ tal que $(p, \epsilon) \vdash_{\mathcal{A}^{\not\epsilon}}^* (q', \epsilon)$. Como $\mathcal{A}^{\not\epsilon}$ no tiene ϵ -transiciones, entonces $p = q'$ y $p \in F^{\not\epsilon}$. Por definición de $F^{\not\epsilon}$, existe $q \in F$ tal que $(p, \epsilon) \vdash_{\mathcal{A}}^* (q, \epsilon)$.

♦ **Caso inductivo:** Sea $w = a \cdot u$ y $p \in Q$:

(\Leftarrow) Sea $q' \in F^{\not\epsilon}$ tal que $(p, au) \vdash_{\mathcal{A}^{\not\epsilon}}^* (q', \epsilon)$. Por definición de $\vdash_{\mathcal{A}^{\not\epsilon}}^*$ existen $p' \in Q$ tal que:

$$(p, au) \stackrel{(1)}{\vdash}_{\mathcal{A}^{\not\epsilon}} (p', u) \stackrel{(2)}{\vdash}_{\mathcal{A}^{\not\epsilon}}^* (q', \epsilon)$$

Por (1) sabemos que $(p, au) \vdash_{\mathcal{A}}^* (p', u)$. (3)

Como $|u| < |au|$ y por (2), por **HI** existe $q \in F$: $(p', u) \vdash_{\mathcal{A}}^* (q, \epsilon)$. (4)

Juntando (3) y (4), tenemos que $(p, au) \vdash_{\mathcal{A}}^* (q, \epsilon)$.

(\Rightarrow) Sea $q \in F$ tal que $(p, au) \vdash_{\mathcal{A}}^* (q, \epsilon)$. Por definición de $\vdash_{\mathcal{A}}^*$ existen $p', p'' \in Q$ tal que:

$$(p, au) \stackrel{(1)}{\vdash}_{\mathcal{A}}^* (p', au) \stackrel{(2)}{\vdash}_{\mathcal{A}}^* (p'', u) \stackrel{(3)}{\vdash}_{\mathcal{A}}^* (q, \epsilon)$$

Por (1) tenemos que $(p, \epsilon) \vdash_{\mathcal{A}}^* (p', \epsilon)$. (4)

Por (2) tenemos que $(p', a, p'') \in \Delta$. (5)

Por (4) y (5), sabemos que $(p, a, p'') \in \Delta^{\not\epsilon}$ y $(p, a \cdot u) \vdash_{\mathcal{A}^{\not\epsilon}}^* (p'', u)$. (6)

Como $|u| < |au|$ y (3), por **HI** existe $q' \in F^{\not\epsilon}$: $(p'', u) \vdash_{\mathcal{A}^{\not\epsilon}}^* (q', \epsilon)$. (7)

Juntando (6) y (7), tenemos que $(p, au) \vdash_{\mathcal{A}^{\not\epsilon}}^* (q', \epsilon)$. ■

Con el teorema 5 demostrado, nuestro mapa de modelos se ve así:

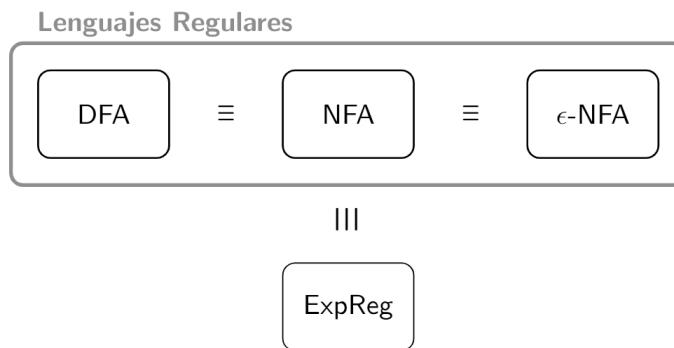


Figura 3: Mapa actual de nuestros modelos de computación

En la siguiente sección mostraremos que todos definen el mismo conjunto de lenguajes.

1.6. Teorema de Kleene

1.6.1. Desde Expresiones a Autómatas

Veremos a continuación que toda ExpReg se puede transformar en un autómata.

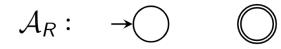
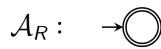
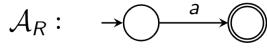
Construcción inductiva. Para cada $R \in \text{ExpReg}$, construimos un ϵ -NFA \mathcal{A}_R :

$$\boxed{\mathcal{A}_R = (Q, \Sigma, \Delta, \{q_0\}, \{q_f\})}$$

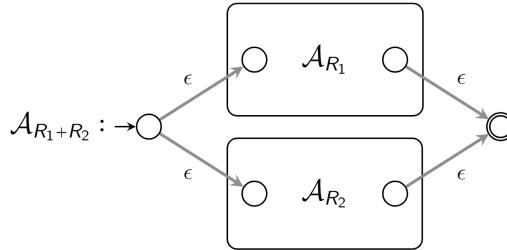
tal que $\mathcal{L}(R) = \mathcal{L}(\mathcal{A}_R)$.

Casos bases. Vemos primero los casos base de la construcción inductiva:

1. si $R = a$,
 2. si $R = \epsilon$:
 3. si $R = \emptyset$:
- para alguna letra $a \in \Sigma$:



4. si $R = (R_1 + R_2)$, donde R_1 y R_2 son expresiones regulares:



Hacemos la construcción inductiva de $R = (R_1 + R_2)$ por **inducción**. Sea \mathcal{A}_{R_1} y \mathcal{A}_{R_2} los ϵ -NFA para R_1 y R_2 , respectivamente, tal que:

- ♦ $\mathcal{A}_{R_1} = (Q_1, \Sigma, \Delta_1, \{q_0^1\}, \{q_f^1\})$
- ♦ $\mathcal{A}_{R_2} = (Q_2, \Sigma, \Delta_2, \{q_0^2\}, \{q_f^2\})$

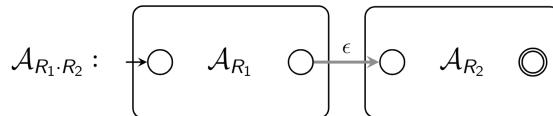
Definimos el ϵ -NFA $\mathcal{A}_{R_1+R_2} = (Q, \Sigma, \Delta, \{q_0\}, \{q_f\})$ tal que:

- ♦ $Q = Q_1 \uplus Q_2 \uplus \{q_0, q_f\}$
- ♦ $\Delta = \Delta_1 \uplus \Delta_2 \uplus \{(q_0, \epsilon, q_0^1), (q_0, \epsilon, q_0^2), (q_0^1, \epsilon, q_f), (q_0^2, \epsilon, q_f)\}$

Proposición. Si $R = (R_1 + R_2)$, entonces $\mathcal{L}(R_1 + R_2) = \mathcal{L}(\mathcal{A}_{R_1+R_2})$.

La demostración de esta proposición queda como ejercicio propuesto para el lector.

5. si $R = (R_1 \cdot R_2)$, donde R_1 y R_2 son expresiones regulares:



Hacemos la construcción inductiva de $R = (R_1 \cdot R_2)$ por **inducción**. Sea \mathcal{A}_{R_1} y \mathcal{A}_{R_2} los ϵ -NFA para R_1 y R_2 , respectivamente, tal que:

- ♦ $\mathcal{A}_{R_1} = (Q_1, \Sigma, \Delta_1, \{q_0^1\}, \{q_f^1\})$
- ♦ $\mathcal{A}_{R_2} = (Q_2, \Sigma, \Delta_2, \{q_0^2\}, \{q_f^2\})$

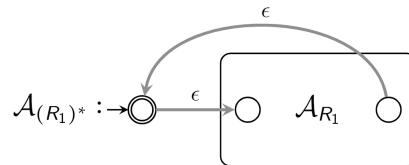
Definimos el ϵ -NFA $\mathcal{A}_{R_1 \cdot R_2} = (Q, \Sigma, \Delta, \{q_0^1\}, \{q_f^2\})$ tal que:

- ♦ $Q = Q_1 \uplus Q_2$
- ♦ $\Delta = \Delta_1 \uplus \Delta_2 \uplus \{(q_f^1, \epsilon, q_0^2)\}$

Proposición. Si $R = (R_1 \cdot R_2)$, entonces $\mathcal{L}(R_1 \cdot R_2) = \mathcal{L}(\mathcal{A}_{R_1 \cdot R_2})$.

La demostración de esta proposición queda como ejercicio propuesto para el lector.

6. si $R = (R_1^*)$, donde R_1 es una expresión regular:



Hacemos la construcción inductiva de $R = (R_1^*)$ por **inducción**. Sea \mathcal{A}_{R_1} el ϵ -NFA para R_1 , tal que:

- ♦ $\mathcal{A}_{R_1} = (Q_1, \Sigma, \Delta_1, \{q_0^1\}, \{q_f^1\})$

Definimos el ϵ -NFA $\mathcal{A}_{(R_1)^*} = (Q, \Sigma, \Delta, \{q_0\}, \{q_0\})$ tal que:

- ♦ $Q = Q_1 \uplus \{q_0\}$
- ♦ $\Delta = \Delta_1 \uplus \{(q_0, \epsilon, q_0), (q_f^1, \epsilon, q_0)\}$

Proposición. Si $R = (R_1^*)$, entonces $\mathcal{L}(R_1^*) = \mathcal{L}(\mathcal{A}_{(R_1^*)})$.

La demostración de esta proposición queda como ejercicio propuesto para el lector.

Teorema 6

Para todo $R \in \text{ExpReg}$, existe un ϵ -NFA \mathcal{A}_R tal que:

$$\mathcal{L}(R) = \mathcal{L}(\mathcal{A}_R)$$

En otras palabras, $\text{ExpReg} \subseteq \epsilon\text{-NFA}$.

1.6.2. Desde Autómatas a Expresiones

Dado un autómata finito no-determinista (que, sin pérdida de generalidad, tiene un estado inicial):

$$\mathcal{A} = (Q, \Sigma, \Delta, q_0, F)$$

Para $X \subseteq Q$ y $p, q \in Q$, considerar el conjunto:

$$\alpha_{p,q}^X \subseteq \Sigma^*$$

tal que $w = a_1 \dots a_n \in \alpha_{p,q}^X$ si, y sólo si, existe una **ejecución**:

$$\rho : p_0 \xrightarrow{a_1} p_1 \xrightarrow{a_2} \dots \xrightarrow{a_n} p_n$$

1. $(p_i, a_{i+1}, p_{i+1}) \in \Delta$ para todo $i \in [0, n - 1]$,
2. $p_0 = p$,
3. $p_n = q$, y
4. $p_i \in X$ para todo $i \in [1, n - 1]$.

*“Intuitivamente, $\alpha_{p,q}^X$ es el conjunto de todas las palabras w tal que existe un **camino** (i.e. ejecución) desde p a q etiquetado por w y **todos los estados** en este camino están en X , con la posible excepción de p y q ”.*

¿Cómo definimos $\mathcal{L}(\mathcal{A})$ en términos de $\alpha_{p,q}^X$? Establecemos el siguiente lema:

$$\mathcal{L}(\mathcal{A}) = \bigcup_{q \in F} \alpha_{q_0, p}^Q$$

Estrategia. Conocida como el algoritmo de McNaughton-Yamada:

1. Para cada $\alpha_{p,q}^X$, definir **inductivamente** una expresión regular $R_{p,q}^X$:

$$\mathcal{L}(R_{p,q}^X) = \alpha_{p,q}^X$$

2. Para $F = \{p_1, \dots, p_k\}$ definir la **expresión regular**:

$$R_{\mathcal{A}} = R_{q_0, p_1}^Q + R_{q_0, p_2}^Q + \dots + R_{q_0, p_k}^Q$$

3. Demostrar que:

$$\mathcal{L}(\mathcal{A}) = \mathcal{L}(R_{\mathcal{A}})$$

Definición inductiva de $R_{p,q}^X$. Tenemos que:

♦ **Caso base:** $X = \emptyset$

Sea $a_1, \dots, a_k \in \Sigma$ todas las letras tal que:

$$(p, a_i, q) \in \Delta$$

- Si $p \neq q$, entonces:

$$R_{p,q}^{\emptyset} \stackrel{\text{def}}{=} \begin{cases} a_1 + \dots + a_k & \text{si } k \geq 1 \\ \emptyset & \text{si } k = 0 \end{cases}$$

- Si $p = q$, entonces:

$$R_{p,q}^{\emptyset} \stackrel{\text{def}}{=} \begin{cases} a_1 + \dots + a_k + \epsilon & \text{si } k \geq 1 \\ \epsilon & \text{si } k = 0 \end{cases}$$

♦ **Caso general:** $X \neq \emptyset$

Por **inducción**, suponemos que para todo $r, s \in Q$ y para todo $Y \subset X$, $R_{r,s}^Y$ es una expresión regular tal que:

$$\mathcal{L}(R_{r,s}^Y) = \alpha_{r,s}^Y$$

Demostramos la construcción para $R_{p,q}^X$ con $p, q \in Q$. Sea $r \in X$ cualquiera:

$$R_{p,q}^X \stackrel{\text{def}}{=} R_{p,q}^{X-\{r\}} + R_{p,r}^{X-\{r\}} \cdot \left(R_{r,r}^{X-\{r\}} \right)^* \cdot R_{r,q}^{X-\{r\}}$$

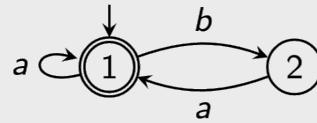
Proposición. Para todo $X \subseteq Q$ y $p, q \in Q$:

$$\mathcal{L}(R_{p,q}^X) = \alpha_{p,q}^X$$

Corolario. $\mathcal{L}(\mathcal{A}) = \mathcal{L}(R_{\mathcal{A}})$

Ejemplo 1.17: Algoritmo MNY

Considere el autómata:



R^\emptyset	1	2
1	$a^?$	b
2	a	ϵ

$R^{\{1\}}$	1	2
1	$a^? + a^? \cdot (a^?)^* \cdot a^?$	$b + a^? \cdot (a^?)^* \cdot b$
2	$a + a \cdot (a^?)^* \cdot a^?$	$\epsilon + a \cdot (a^?)^* \cdot b$

\equiv

$R^{\{1\}}$	1	2
1	a^*	a^*b
2	a^+	$\epsilon + a^+b$

$R^{\{1,2\}}$	1	2
1	$a^* + a^*b(\epsilon + a^+b)^* a^+$	$a^*b + a^*b(\epsilon + a^+b)^*(\epsilon + a^+b)$
2	$a^+ + (\epsilon + a^+b)(\epsilon + a^+b)^* a^+$	$(\epsilon + a^+b) + (\epsilon + a^+b)(\epsilon + a^+b)^*(\epsilon + a^+b)$

\equiv

$R^{\{1,2\}}$	1	2
1	$a^* (ba^+)^*$	$(a^*b)(a^+b)^*$
2	$(a^+b)^* a^+$	$(a^+b)^*$

2. Propiedades de lenguajes regulares

2.1. Lema de bombeo

Supongamos que deseamos aceptar el siguiente lenguaje:

$$L = \{a^i b^i \mid i \geq 0\} = \{\epsilon, ab, aabb, aaabbb, aaaaabbbb, \dots\}$$

con un **autómata finito determinista**. ¿Es posible? La respuesta es que no, ya que nuestros autómatas no tienen la capacidad de “contar”. Por ende, L sería un lenguaje NO regular ya que no podría ser definido por un autómata.

Enunciado. Sea $L \subseteq \Sigma^*$. Si L es **regular**, entonces:

(LB) existe un $N > 0$ tal que
 para toda palabra $x \cdot y \cdot z \in L$ con $|y| \geq N$
 existen palabras $u \cdot v \cdot w = y$ con $v \neq \epsilon$ tal que
 para todo $i \geq 0$, $x \cdot u \cdot v^i \cdot w \cdot z \in L$.

El contrapositivo del lema de bombeo nos servirá para demostrar que un lenguaje L NO es regular.

Sea $L \subseteq \Sigma^*$. Si:

(\neg LB) para todo $N > 0$
 existe una palabra $x \cdot y \cdot z \in L$ con $|y| \geq N$ tal que
 para todo $u \cdot v \cdot w = y$ con $v \neq \epsilon$
 existe un $i \geq 0$, $x \cdot u \cdot v^i \cdot w \cdot z \notin L$.
 entonces L NO es regular.



“ L NO es regular”



“ L es regular”

El escoge un $N > 0$

Uno escoge $x \cdot y \cdot z \in L$ con $|y| \geq N$

El escoge $u \cdot v \cdot w = y$ con $v \neq \epsilon$

Uno escoge $i \geq 0$

Uno gana si $xuv^i wz \notin L$

El gana si $xuv^i wz \in L$

LB versión juego. “Dado un lenguaje $L \subseteq \Sigma^*$, si **UNO** tiene una estrategia ganadora en el juego (\neg LB) para toda estrategia posible del demonio, entonces L NO es regular”. Con **estrategia**, nos referimos a todas las movidas posibles que podría ejecutar el **demonio** (considerar todos los casos posibles de sus elecciones).

Ejemplo 2.1

Considere el lenguaje $L = \{a^i b^i \mid i \geq 0\}$:



" $a^n b^n$ NO es regular"



" $a^n b^n$ es regular"

Escojo $N > 0$

Yo escojo $\underbrace{a^N}_x \cdot \underbrace{b^N}_y \cdot \underbrace{\epsilon}_z \in L$

Entonces escojo $\underbrace{b^n}_u \cdot \underbrace{b^m}_v \cdot \underbrace{b^l}_w = \underbrace{b^N}_y$ con $m > 0$

Yo escojo $i = 2$

Ganamos el juego ya que con $i = 2$ estaremos bombeando más b -letras y entonces la palabra no tendrá la misma cantidad de a -letras que de b -letras ($i \neq j$), por ende, L NO es regular.

Ejemplo 2.2

Considere el lenguaje $L = \{a^n b^m \mid n \geq m\}$:



" $a^n b^m$ NO es regular"



" $a^n b^m$ es regular"

Escojo $N > 0$

Yo escojo $\underbrace{a^N}_x \cdot \underbrace{b^N}_y \cdot \underbrace{\epsilon}_z \in L$

Entonces escojo $\underbrace{b^j}_u \cdot \underbrace{b^k}_v \cdot \underbrace{b^l}_w = \underbrace{b^N}_y$ con $k > 0$

Yo escojo $i = 2$

Ganamos el juego ya que, nuevamente, con $i = 2$, estaremos bombeando más b -letras y entonces la palabra puede tener más b -letras que a -letras ($n < m$), y así L NO es regular.

Ejemplo 2.3

Considere el lenguaje $L = \{w \cdot w \mid w \in \{a, b\}^*\}$



" L NO es regular"



" L es regular"

Escojo $N > 0$

Yo escojo $\underbrace{a^N}_x \underbrace{b}_y \cdot \underbrace{a^N}_z \cdot \underbrace{b}_y \in L$

Entonces escojo $\underbrace{a^j}_u \cdot \underbrace{a^k}_v \cdot \underbrace{a^l}_w = \underbrace{a^N}_y$ con $k > 0$

Yo escojo $i = 0$

Ganamos el juego ya que con la elección de $i = 0$ estamos haciendo que una de las mitades de la palabra sea distinta a su otra mitad, por ende, L NO es regular.

Ejemplo 2.4

Considere el lenguaje $L = \{a^{2^n} \mid n > 0\}$



" a^{2^n} NO es regular"



" a^{2^n} es regular"

Escojo $N > 0$

Yo escojo $\underbrace{a^{2^n-N}}_x \cdot \underbrace{a^N}_y \cdot \underbrace{\epsilon}_z \in L$ con $N < 2^n$

Entonces escojo $\underbrace{a^j}_u \cdot \underbrace{a^k}_v \cdot \underbrace{a^l}_w = \underbrace{a^N}_y$ con $k > 0$

Yo escojo $i = 2$

Ganamos el juego ya que con la elección de $i = 2$, tenemos que en la elección de y tendremos una mayor cantidad de a -letras bombeadas y se romperá el equilibrio $2^N - N + N$, por ende, L NO es regular.

2.2. Minimización de autómatas

¿Cómo minimizamos un autómata finito?

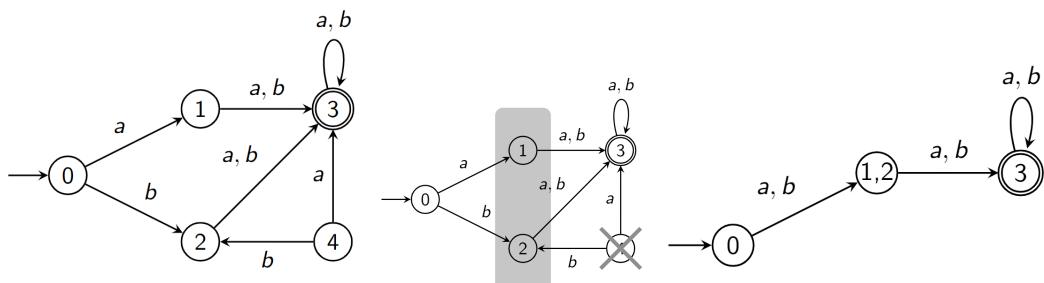
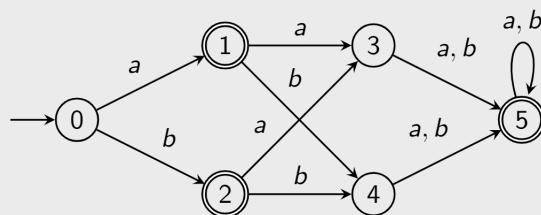


Figura 4: Idea de minimización

1. Eliminar estados inaccesibles.
 - ♦ Fácil de realizar y no cambia el lenguaje del autómata finito.
2. Colapsar estados “equivalentes”.
 - ♦ ¿Cómo sabemos cuáles estados colapsar y cuáles no?

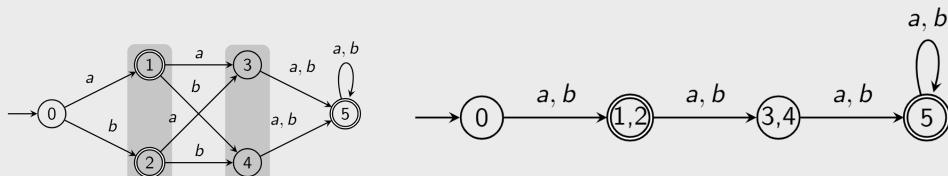
Ejemplo 2.5

Considere el siguiente autómata:



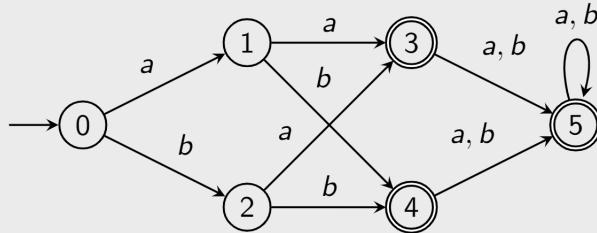
Podemos:

- ♦ Colapsar estados 1 y 2.
- ♦ Colapsar estados 3 y 4.



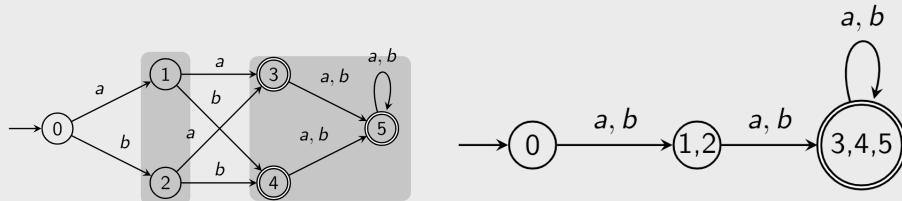
Ejemplo 2.6

Considere el siguiente autómata:



Podemos:

- ♦ Colapsar estados 1 y 2.
- ♦ Colapsar estados 3, 4 y 5.

**2.2.1. Colapsar estados**

Definición. Sea $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ un DFA.

Se define la **función de transición extendida** $\hat{\delta} : Q \times \Sigma^* \rightarrow Q$ inductivamente como:

$$\boxed{\begin{aligned}\hat{\delta}(q, \epsilon) &\stackrel{\text{def}}{=} q \\ \hat{\delta}(q, w \cdot a) &\stackrel{\text{def}}{=} \delta(\hat{\delta}(q, w), a)\end{aligned}}$$

Definición. Decimos que p y q son **indistinguibles** ($p \approx_{\mathcal{A}} q$) si:

$$p \approx_{\mathcal{A}} q \quad \text{ssi} \quad (\hat{\delta}(p, w) \in F \Leftrightarrow \hat{\delta}(q, w) \in F), \text{ para todo } w \in \Sigma^*.$$

Decimos que p y q son **distingüibles** si NO son indistinguibles ($p \not\approx_{\mathcal{A}} q$).

Recordatorio relaciones de equivalencia. Una relación \approx_R sobre un conjunto X se dice de **equivalencia** si es:

- ♦ **Refleja:** $\forall p \in X. p \approx_R p$
- ♦ **Simétrica:** $\forall p, q \in X. \text{ si } p \approx_R q \text{ entonces } q \approx_R p$.
- ♦ **Transitiva:** $\forall p, q, r \in X. \text{ si } p \approx_R q \text{ y } q \approx_R r, \text{ entonces } p \approx_R r$.

Para un elemento $p \in X$ se define su **clase de equivalencia** según \approx_R como:

$$[p]_{\approx_R} = \{q \mid q \approx_R p\}$$

Una función $f : X \rightarrow X$ se dice **bien definida** sobre \approx_R si:

$$p \approx_R q \text{ entonces } f(p) \approx_R f(q)$$

Propiedades de \approx_A . Tenemos que:

- ♦ \approx_A es una **relación de equivalencia**, es decir, es refleja, simétrica y transitiva.
- ♦ Cada estado $p \in Q$ está en exactamente una clase de equivalencia:

$$[p]_{\approx_A} = \{q \mid q \approx_A p\}$$

- ♦ Para todo $a \in \Sigma$ la función $\delta(\cdot, a) : Q \rightarrow Q$ está **bien definida** sobre \approx_A :

$$p \approx_A q \text{ entonces } \delta(p, a) \approx_A \delta(q, a)$$

El autómata cuociente. Para un DFA $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ se define el DFA:

$\mathcal{A}/\approx = (Q_\approx, \Sigma, \delta_\approx, q_\approx, F_\approx)$

- ♦ $Q_\approx = \{[p]_{\approx_A} \mid p \in Q\}$
- ♦ $\delta_\approx([p]_{\approx_A}, a) = [\delta(p, a)]_{\approx_A}$
- ♦ $q_\approx = [q_0]_{\approx_A}$
- ♦ $F_\approx = \{[p]_{\approx_A} \mid p \in F\}$

Teorema 7

Para todo autómata finito determinista \mathcal{A} se cumple que:

$$\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}/\approx)$$

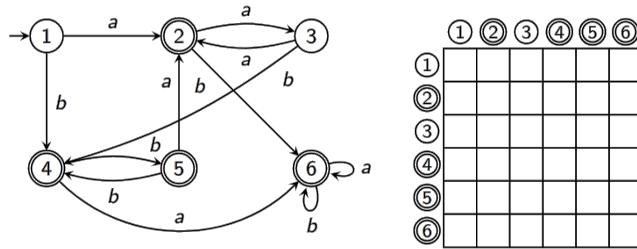
Demostración. PENDIENTE. Pero se invita al lector a hacerla <3.

2.2.2. Algoritmo de minimización

El objetivo es buscar los pares de estados que son **distingibles**:

1. Construya una tabla con los pares $\{p, q\}$ inicialmente sin marcar.
2. Marque $\{p, q\}$ si $p \in F$ y $q \notin F$ o viceversa.
3. Repita este paso hasta que no hayan más cambios:
 - ♦ Si $\{p, q\}$ no están marcados y $\{\delta(p, a), \delta(q, a)\}$ están marcados para algún $a \in \Sigma$, entonces marque $\{p, q\}$.
4. Al terminar, $p \not\approx_A q$ si, y sólo si, la entrada $\{p, q\}$ está marcada.

Veamos como funciona el algoritmo. Considere el siguiente autómata y una tabla que relacione todos los pares de estados:



1. Construya una tabla con los pares $\{p, q\}$ inicialmente sin marcar.
2. Marque $\{p, q\}$ si $p \in F$ y $q \notin F$ o viceversa.

	1	2	3	4	5	6
1						
2						
3						
4						
5						
6						

	1	2	3	4	5	6
1	✓					
2		✓				
3			✓	✓	✓	✓
4						
5						
6						

3. Repita este paso hasta que no hayan más cambios:

- ♦ Si $\{p, q\}$ no están marcados y $\{\delta(p, a), \delta(q, a)\}$ estan marcados para algún $a \in \Sigma$, entonces marque $\{p, q\}$.

	1	2	3	4	5	6
1	✓			✓	✓	✓
2		✓		a (3,6)	a (3,2)	a (3,6)
3			✓	✓	✓	✓
4						
5						
6						

	1	2	3	4	5	6
1	✓			✓	✓	✓
2		✓		✓	✓	✓
3			✓	✓	✓	✓
4						
5						
6						

	1	2	3	4	5	6
1	✓			✓	✓	✓
2		✓		✓	✓	✓
3			✓	✓	✓	✓
4					a (6,2)	
5						a (2,6)
6						

	1	2	3	4	5	6
1	✓			✓	✓	✓
2		✓	✓	✓	✓	✓
3			✓	✓	✓	✓
4				✓		
5					✓	
6						✓

	1	2	3	4	5	6
1	✓			✓	✓	✓
2		✓		✓	✓	✓
3			✓	✓	✓	✓
4				✓		b (5,6)
5					✓	
6						

	1	2	3	4	5	6
1	✓			✓	✓	✓
2		✓		✓	✓	✓
3			✓	✓	✓	✓
4				✓		✓
5					✓	
6						✓

4. Al terminar, $p \not\approx_{\mathcal{A}}$ si, y sólo si, la entrada $\{p, q\}$ está marcada.

Así, vemos que los pares indistinguibles son todas las entradas **NO** marcadas.

2.3. Teorema de Myhill-Nerode

La sección anterior deja muchas incógnitas:

1. ¿Cómo sabemos si el autómata del algoritmo es un **mínimo**?
2. Dado L , ¿existe un **único** autómata mínimo?
3. Dado un \mathcal{A} , ¿es posible **construir** un autómata mínimo equivalente?

En esta sección, demostraremos que:

- ♦ El autómata con el mínimo de estados es **único**.
- ♦ El algoritmo de minimización **siempre** construye el autómata mínimo.

Estrategia. Para demostrar lo dicho anteriormente, seguimos los siguientes pasos:

1. Desde un DFA \mathcal{A} , definiremos una relación de equivalencia (RE) $\equiv_{\mathcal{A}}$ entre palabras en Σ^* .
2. Desde una $RE \equiv$ entre palabras, construiremos un DFA \mathcal{A}_{\equiv} .
3. A partir de un lenguaje L , definiremos una $RE \equiv_L$.
4. \mathcal{A}_{\equiv_L} define el autómata con la **menor cantidad de estados**.
5. \mathcal{A}_{\equiv_L} es equivalente al resultado de nuestro **algoritmo de minimización**.

2.3.1. Relaciones de Myhill-Nerode

Sea $L \subseteq \Sigma^*$ cualquier lenguaje.

Definición. Una relación de equivalencia \equiv en Σ^* es de **Myhill-Nerode** para L si:

1. \equiv es una **congruencia por la derecha**.
2. \equiv **refina** L .
3. El número de clases de equivalencia de \equiv es **finita**.

A partir de una relación \equiv de Myhill-Nerode podemos construir un DFA \mathcal{A}_{\equiv} .

$$\boxed{\begin{array}{l} \mathcal{A} \Rightarrow \equiv_{\mathcal{A}} \\ \equiv \Rightarrow \mathcal{A}_{\equiv} \end{array}}$$

Construcción del DFA \mathcal{A}_{\equiv} . Dada una relación de Myhill-Nerode \equiv para $L \subseteq \Sigma^*$, definimos el autómata:

$$\boxed{\mathcal{A}_{\equiv} = (Q_{\equiv}, \Sigma, \delta_{\equiv}, q_{\equiv}, F_{\equiv})}$$

- ♦ $Q_{\equiv} = \{[w]_{\equiv} \mid w \in \Sigma^*\}$
- ♦ $q_{\equiv} = [\epsilon]_{\equiv}$
- ♦ $F_{\equiv} = \{[w]_{\equiv} \mid w \in L\}$
- ♦ $\delta_{\equiv} ([w]_{\equiv}, a) = [wa]_{\equiv}$

Teorema 8

Cada cualquier $L \subseteq \Sigma^*$, tenemos que

$$\mathcal{L}(\mathcal{A}_\equiv) = L$$

Podemos establecer que $\mathcal{A} \Rightarrow \equiv_{\mathcal{A}}$ y $\equiv \Rightarrow \mathcal{A}_\equiv$ son procesos inversos, conclusión que se ilustra en el siguiente teorema.

Teorema 9

1. Si \mathcal{A} es un DFA que acepta L y si construimos:

$$\mathcal{A} \Rightarrow \equiv_{\mathcal{A}} \Rightarrow \mathcal{A}_{\equiv_{\mathcal{A}}}$$

entonces \mathcal{A} es **isomorfo** (“equivalente”) a $\mathcal{A}_{\equiv_{\mathcal{A}}}$.

2. Si \equiv es una relación de Myhill-Nerode para L y si construimos:

$$\equiv \Rightarrow \mathcal{A}_\equiv \Rightarrow \equiv_{\mathcal{A}_\equiv}$$

entonces la relación \equiv es **equivalente** a $\equiv_{\mathcal{A}_\equiv}$.

La demostración de ambos teoremas queda propuesto como ejercicio al lector.

2.3.2. Camino al teorema

Antes de enunciar el Teorema de Myhill-Nerode, debemos aún mencionar algunas definiciones.

Definición. Dado un lenguaje $L \subseteq \Sigma^*$, se define la relación de equivalencia \equiv_L como:

$$u \equiv_L v \quad \text{ssi} \quad (u \cdot w \in L \Leftrightarrow v \cdot w \in L) \quad \forall w \in \Sigma^*$$

Ejemplo 2.7

Sea $L = (ab)^*$. Algunas clases de equivalencia para L son:

- ♦ $[\epsilon]_{\equiv_L} = \{\epsilon, ab, abab, ababab, \dots\}$.
- ♦ $[a]_{\equiv_L} = \{a, aba, ababa, abababa, \dots\}$.
- ♦ $[b]_{\equiv_L} = \{b, bb, ba, abb, \dots\}$

Propiedades. \equiv_L se caracteriza por:

1. Ser una **congruencia por la derecha**:

$$u \equiv_L v \text{ entonces } u \cdot w \equiv_L v \cdot w \quad \forall w \in \Sigma^*$$

2. Refinar a L :

$$u \equiv_L v \text{ entonces } (u \in L \Leftrightarrow v \in L)$$

3. Si \equiv es una congruencia por la derecha y refina L , entonces \equiv **refina** a \equiv_L :

$$u \equiv v \text{ entonces } u \equiv_L v$$

Con todo lo anterior, estamos en condiciones de enunciar el teorema.

Teorema 10

Sea $L \subseteq \Sigma^*$. Las siguientes propiedades son equivalentes:

1. L es **regular**.
2. Existe una **relación de Myhill-Nerode** para L .
3. La relación \equiv_L tiene una cantidad **finita** de clases de equivalencia.

Demostración teorema 10. Del punto 1 al 2, tenemos que si L es regular, entonces:

- ♦ existe un autómata finito \mathcal{A} tal que $L = \mathcal{L}(\mathcal{A})$.
- ♦ $\equiv_{\mathcal{A}}$ es una relación de Myhill-Nerode para L .

Del punto 2 al 3, sea \equiv una relación de Myhill-Nerode para L , entonces:

- ♦ \equiv tiene una cantidad finita de clases de equivalencia.
- ♦ \equiv_L tiene una cantidad finita de clases de equivalencia.

Del punto 3 al 1, si \equiv_L tiene una cantidad **finita** de clases de equivalencia, entonces:

- ♦ \equiv_L es una relación de Myhill-Nerode para L .
- ♦ \mathcal{A}_{\equiv_L} es un autómata finito para L . ■

Conclusiones del teorema. Tenemos que:

1. $\equiv_L \Rightarrow \mathcal{A}_{\equiv_L}$ produce el autómata con la menor cantidad de estados.
2. Todo autómata \mathcal{A} tal que $\equiv_{\mathcal{A}} = \equiv_L$ son **isomorfos** (“equivalentes”).
3. El **algoritmo de minimización** produce un autómata isomorfo \mathcal{A}_{\equiv_L} .

Demostración punto 3. Sea $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ un autómata que acepta L ya **minimizado**:

$$\begin{aligned}
 u \equiv_L v &\Leftrightarrow (u \cdot w \in L \Leftrightarrow v \cdot w \in L) \quad \forall w \in \Sigma^* \\
 &\Leftrightarrow \left(\hat{\delta}(q_0, u \cdot w) \in F \Leftrightarrow \hat{\delta}(q_0, v \cdot w) \in F \right) \quad \forall w \in \Sigma^* \\
 &\Leftrightarrow \left(\hat{\delta}(\hat{\delta}(q_0, u), w) \in F \Leftrightarrow \hat{\delta}(\hat{\delta}(q_0, v), w) \in F \right) \quad \forall w \in \Sigma^* \\
 &\Leftrightarrow \hat{\delta}(q_0, u) \approx_{\mathcal{A}} \hat{\delta}(q_0, v) \\
 &\Leftrightarrow \hat{\delta}(q_0, u) = \hat{\delta}(q_0, v) \\
 &\Leftrightarrow u \equiv_{\mathcal{A}} v
 \end{aligned}$$

2.4. Autómatas en dos direcciones

2.4.1. Definición de un 2DFA

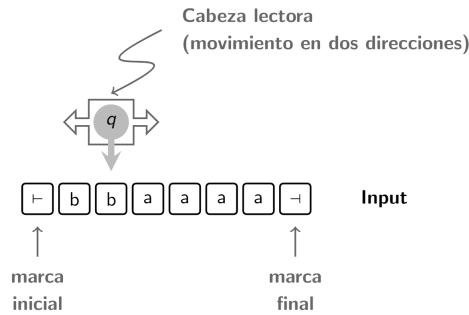


Figura 5: Representación de un 2DFA

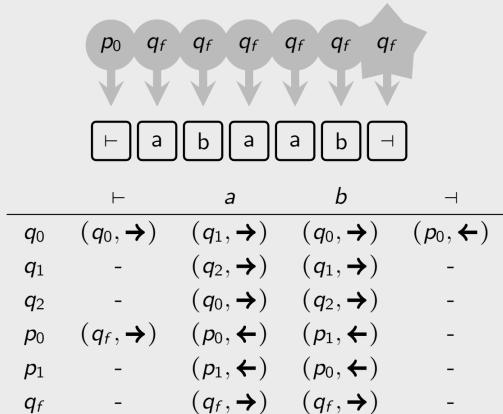
Definición. Un autómata finito determinista en 2 direcciones (2DFA) es una estructura:

$$\boxed{\mathcal{A} = (Q, \Sigma, \vdash, \dashv, \delta, q_0, q_f)}$$

- ♦ Q es un conjunto finito de estados.
- ♦ Σ es el alfabeto del input.
- ♦ \vdash y \dashv son las marcas (símbolos) iniciales y finales.
- ♦ $\delta : Q \times (\Sigma \cup \{\vdash, \dashv\}) \rightarrow Q \times \{\leftarrow, \rightarrow\}$ es la **función parcial de transición**.
- ♦ q_0 es el estado inicial.
- ♦ q_f es el estado final.

Ejemplo 2.8

$$L = \{ w \in \Sigma^* \mid (\#a(w) = 0 \bmod 3) \text{ y } (\#b(w) = 0 \bmod 2) \}$$



Configuración. Sea \mathcal{A} un 2DFA y $w = a_1 \dots a_n \in \Sigma^*$ un input. Definimos $a_0 = \vdash$ y $a_{n+1} = \dashv$ tal que el input se define como:

$$a_0 a_1 \dots a_n a_{n+1} = \vdash \cdot w \cdot \dashv$$

Una **configuración** de \mathcal{A} sobre w viene dado por un par:

$$(q, i) \in Q \times \{0, \dots, n + 1\}$$

- ♦ q es el **estado actual** del autómata.
- ♦ i es la **posición actual** de la cabeza lectora.

Se define la relación de **siguiente configuración** $\xrightarrow{\mathcal{A}}$ de \mathcal{A} sobre w como:

$$(p, i) \xrightarrow{\mathcal{A}} (q, j)$$

tal que:

- ♦ Si $\delta(p, a_i) = (q, \rightarrow)$, entonces $(p, i) \xrightarrow{\mathcal{A}} (q, i + 1)$.
- ♦ Si $\delta(p, a_i) = (q, \leftarrow)$, entonces $(p, i) \xrightarrow{\mathcal{A}} (q, i - 1)$.

Ejecución. Una ejecución (o *run*) ρ de \mathcal{A} sobre w es una secuencia de configuraciones:

$$\rho : (p_0, i_0) \rightarrow (p_1, i_1) \rightarrow \dots \rightarrow (p_m, i_m)$$

donde $p_0 = q_0$ y $i_0 = 0$. Además, $(p_j, i_j) \xrightarrow{\mathcal{A}} (p_{j+1}, i_{j+1}) \quad \forall j \in [0, m - 1]$.

Una ejecución ρ de \mathcal{A} sobre w es de **aceptación** si:

$$p_m = q_f \quad y \quad i_m = n + 1$$

Aceptación. Decimos que \mathcal{A} **acepta** w si hay una ejecución de \mathcal{A} sobre w que es de **aceptación**. Así, el lenguaje **aceptado** por \mathcal{A} se define como:

$$\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ acepta } w\}$$

Ojo: Notemos que un 2DFA puede pasar por error o NO parar nunca.

2.4.2. 2DFA vs DFA

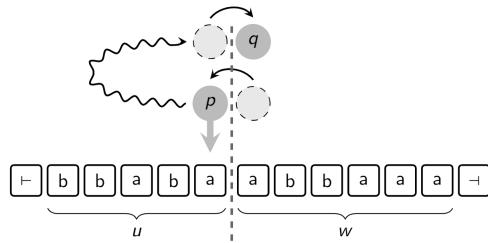
Para todo lenguaje regular L existe un 2DFA \mathcal{A} :

$$L = \mathcal{L}(\mathcal{A})$$

En otras palabras, $\text{DFA} \subseteq \text{2DFA}$.

¿Son los 2DFA más poderosos que los DFA? En esta sección, demostraremos que no, es decir, pueden representar al mismo conjunto de lenguajes.

¿Cuánta información puede almacenar un 2DFA? Veamos la siguiente figura:



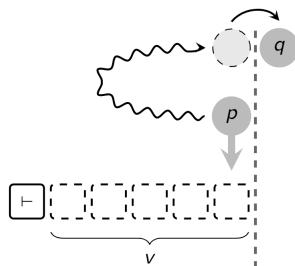
“Cada vez que \mathcal{A} cruce de w a u en el estado p , \mathcal{A} cruzará de regreso en el estado q ”.

Para cada $u \in \Sigma^*$, definimos la función $T_u : Q \cup \{\bullet\} \rightarrow Q \cup \{\perp\}$ tal que:

- ♦ $T_u(p) = q$ si, y sólo si, desde $(p, |u|)$ cruza en la configuración $(q, |u| + 1)$.
- ♦ $T_u(p) = \perp$ si, y sólo si, desde $(p, |u|)$ nunca cruza de u .
- ♦ $T_u(\bullet) = q$ si, y sólo si, desde $(q_0, 0)$ cruza por primera vez con $(q, |u| + 1)$.
- ♦ $T_u(\bullet) = q$ si, y sólo si, desde $(q_0, 0)$ cruza por primera vez con $(q, |u| + 1)$.
- ♦ $T_u(\bullet) = \perp$ si, y sólo si, desde $(q_0, 0)$ nunca cruza de u .

Suponga que ahora tenemos una palabra v tal que:

$$T_v = T_u$$



Luego, respecto a la función T_u podemos decir que:

1. Definimos la relación \equiv_T entre palabras en Σ^* tal que:

$$u \equiv_T v \text{ si, y solo si, } T_u = T_v$$

es una **relación de equivalencia**.

2. \equiv_T es una **congruencia por la derecha**:

$$u \equiv_T v \Rightarrow \forall w \cdot u \cdot w \equiv_T v \cdot w$$

3. \equiv_T **refina** a $\mathcal{L}(\mathcal{A})$:

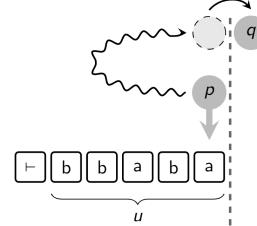
$$u \equiv_T v \Rightarrow (u \in L \Leftrightarrow v \in L)$$

4. La relación \equiv_T tiene una **cantidad finita** de clases de equivalencia:

$$T : Q \cup \{\bullet\} \rightarrow Q \cup \{\perp\}$$

Así, $\equiv_{\mathcal{L}(\mathcal{A})}$ también tiene una cantidad **finita** de **clases de equivalencia**.

Este comportamiento solo depende de u , y no de w :



Entonces, v es **indistinguible** de u según \mathcal{A} . En otras palabras:

$$u \cdot w \in \mathcal{L}(\mathcal{A}) \Leftrightarrow v \cdot w \in \mathcal{L}(\mathcal{A}) \quad \text{para todo } w \in \Sigma^*$$

Lo anterior tiene conexión con las relaciones de Myhill-Nerode, que vimos en una sección anterior.

Teorema 11

Para todo 2DFA \mathcal{A} existe un DFA \mathcal{A}' tal que:

$$\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}')$$

En otras palabras, $2\text{DFA} \equiv \text{DFA}$.

El DFA del teorema anterior se construye usando el Teorema de Myhill-Nerode a partir de las funciones T_u . La demostración del teorema y la construcción queda como ejercicio propuesto para el lector.

3. Algoritmos para lenguajes regulares

3.1. Evaluación de expresiones regulares

Expresiones regulares en la práctica. El lenguaje para definir expresiones regulares en la práctica se conoce como **RegEx** (o RexExp). Las sintaxis más usadas para definir **RegEx** son:

1. POSIX (Portable Operating System Interface for uniX).
2. Perl (PCRE = Perl Compatible Regular Expressions).

	RegEx	Teoría
Carácter	a	a
Escape	\+	—
Cualquiera	.	Σ
Clase	[abc]	$(a + b + c)$
Clase consecutivo	[a-zA-Z]	$(a + \dots + z + A + \dots + Z)$
Clase exclusivo	[^0-9]	$(+ +_{a \in \Sigma - \{0, \dots, 9\}} a)$
Alternación	cat dog	cat + dog
0 o 1	R?	$R^?$
1 o más	R+	R^+
0 o más	R*	R^*
entre n y m	R{n,m}	$R^n(R + \epsilon)^{m-n}$
Backreference	(R)\1	?
...

¿Cómo evaluamos una expresión regular? Tenemos el siguiente problema:

PROBLEMA: Evaluación de expresiones regulares
INPUT: una expresión regular R
 un documento w
OUTPUT: TRUE si, y sólo si, $w \in \mathcal{L}(R)$

Donde el tamaño del input está dado por:

- ♦ $|R| :=$ número de letras y operadores.
- ♦ $|w| :=$ largo del documento.

Buscamos un algoritmo polinomial en $|R|$ y $|w|$. Algo a notar es que $|R| \ll |w|$:

- ♦ $|R|$ puede ser del orden **decenas** operadores (~ 1 KB).
- ♦ $|w|$ puede ser del orden de **miles a millones** de símbolos (~ 100 MB).

Análisis de tiempo diferenciado. Tenemos dos tipos:

- ♦ **Combined complexity:** expresión y documento son parte del input.
- ♦ **Data complexity:** solo documento es parte del input (tamaño de expresión es considerada una constante).

Buscamos algoritmos que sean **lineales en data complexity** y ojalá **polinomiales en combined complexity**.

Volviendo a nuestro problema:

PROBLEMA:	Evaluación de expresiones regulares
INPUT:	una expresión regular R
	un documento w
OUTPUT:	TRUE si, y sólo si, $w \in \mathcal{L}(R)$

Los pasos que podemos seguir para resolverlo son:

1. Convertimos R a un ϵ -NFA \mathcal{A}_R .
2. Verificamos si $w \in \mathcal{L}(\mathcal{A}_R)$.

Ahora, el tamaño del input es:

- ♦ $|w| :=$ largo del documento.
- ♦ $|\mathcal{A}| := |Q| + |\Delta|$.

Hay varias soluciones para resolver nuestro problema, que iremos enumerado a continuación.

Solución 1: Backtracking. Esta solución es la que usa la mayoría de los motores de RegEx en la práctica. Se puede hacer una prueba en <https://regexpr.com/>.

Solución 2: DFA. Para un DFA $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ y una palabra $w = a_1 \dots a_n$.

```
Function eval-DFA( $\mathcal{A}$ ,  $w$ ):
     $q := q_0$ 
    for  $i = 1$  to  $n$  do
         $| q := \delta(q, a_i)$ 
    return check ( $q \in F$ )
```

Solución 3: NFA determinización. Para un NFA $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ y una palabra $w = a_1 \dots a_n$.

```
Function eval-NFA( $\mathcal{A}$ ,  $w$ ):
     $\mathcal{A}^{\text{det}} := \text{NFAToDFA}(\mathcal{A})$ 
    return eval-DFA( $\mathcal{A}^{\text{det}}$ ,  $w$ )
```

¿Es necesario construir la determinización completa? La respuesta es que no. Recordemos que es la determinización. Para un autómata no-determinista $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$, definimos el autómata determinista (**determinización** de \mathcal{A}):

$$\boxed{\mathcal{A}^{\text{det}} = (2^Q, \Sigma, \delta^{\text{det}}, q_0^{\text{det}}, F^{\text{det}})}$$

- ♦ $2^Q = \{S \mid S \subseteq Q\}$ es el conjunto potencia de Q .

♦ $q_0^{\det} = I$.

♦ $\delta^{\det} : 2^Q \times \Sigma \rightarrow 2^Q$ tal que:

$$\delta^{\det}(S, a) = \{q \in Q \mid \exists p \in S. (p, a, q) \in \Delta\}$$

♦ $F^{\det} = \{S \in 2^Q \mid S \cap F \neq \emptyset\}$, es decir, todos los conjuntos que tengan al menos un estado final.

Solución 4: NFA *on-the-fly*. Fijémonos en la función de transición $\delta^{\det} : 2^Q \times \Sigma \rightarrow 2^Q$ tal que:

$$\delta^{\det}(S, a) = \{q \in Q \mid \exists p \in S. (p, a, q) \in \Delta\}$$

La estrategia *on-the-fly* corresponde a:

1. Mantenemos un conjunto S de estados actuales.
2. Por cada nueva letra a , calculamos el conjunto $\delta^{\det}(S, a)$.

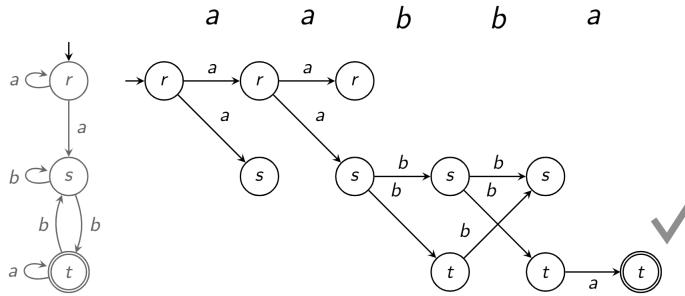


Figura 6: Ejemplo de NFA *on-the-fly*

Definimos el algoritmo para un NFA $\mathcal{A} = (Q, \Sigma, \Delta, I, F)$ y una palabra $w = a_1 \dots a_n$ como:

```

Function eval-NFAonthefly( $\mathcal{A}, w$ ):
   $S := I$ 
  for  $i = 1$  to  $n$  do
     $S_{\text{old}} := S$ 
     $S := \emptyset$ 
    foreach  $p \in S_{\text{old}}$  do
       $| \quad S := S \cup \{q \mid (p, a_i, q) \in \Delta\}$ 
  return check ( $S \cap F \neq \emptyset$ )
  
```

Resumen y complejidades. Para un autómata \mathcal{A} y una palabra $w = a_1 \dots a_n$ tenemos que:

	Tiempo
Backtracking	$\mathcal{O}(\mathcal{A} ^{ w })$
DFA	$\mathcal{O}(\mathcal{A} + w)$
NFA	$\mathcal{O}(2^{ Q } + w)$
ϵ -NFA <i>on-the-fly</i>	$\mathcal{O}(\mathcal{A} \cdot w)$

3.2. Transductores

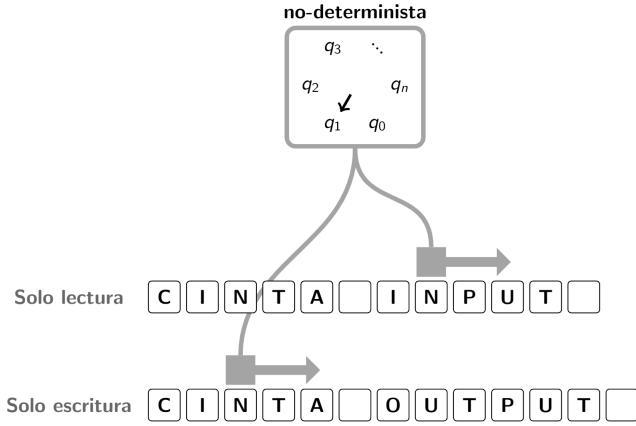


Figura 7: Representación de un transductor

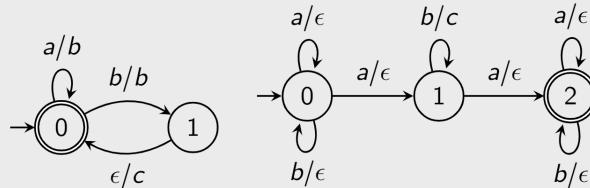
Definición. Un transductor (en inglés, *transducer*) es una tupla:

$$\boxed{\mathcal{T} = (Q, \Sigma, \Omega, \Delta, I, F)}$$

- ♦ Q es un conjunto finito de estados.
- ♦ Σ es el alfabeto del input.
- ♦ Ω es el alfabeto del **output**.
- ♦ $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Omega \cup \{\epsilon\}) \times Q$ es la **relación de transición**.
- ♦ $I \subseteq Q$ es un conjunto de estados iniciales.
- ♦ $F \subseteq Q$ es el conjunto de estados finales.

Ejemplo 3.1

Algunos transductores:



Configuración. Sea \mathcal{T} un transductor. Definimos:

- ♦ Un par $(q, u, v) \in Q \times \Sigma^* \times \Omega^*$ es una **configuración** de \mathcal{T} .
- ♦ Una configuración (q, u, ϵ) es **inicial** si $q \in F$.
- ♦ Una configuración (q, ϵ, v) es **final** si $q \in F$.

"Intuitivamente, una configuración (q, au, vb) representa que \mathcal{T} se encuentra en el estado q procesando la palabra au y leyendo a , y hasta ahora grabó la palabra vb y el último símbolo impreso es b ".

Ejecución. Se define la relación $\vdash_{\mathcal{T}} \subseteq (Q \times \Sigma^* \times \Omega^*) \times (Q \times \Sigma^* \times \Omega^*)$ de **siguiente-paso** entre configuraciones de \mathcal{T} :

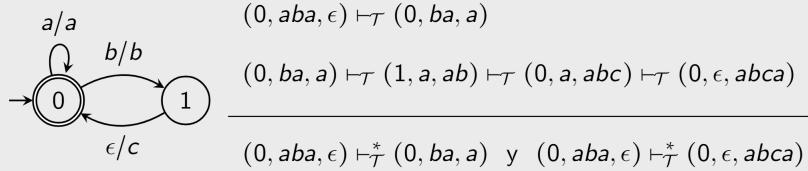
$$(p, u_1, v_1) \vdash_{\mathcal{T}} (q, u_2, v_2)$$

si, y sólo si, existe $(p, a, b, q) \in \Delta$ tal que $u_1 = a \cdot u_2$ y $v_2 = v_1 \cdot b$.

Se define $\vdash_{\mathcal{T}}^*$ como la clausura **refleja** y **transitiva** de $\vdash_{\mathcal{T}}$:

$$\boxed{\text{para toda configuración } (q, u, v) : (q, u, v) \vdash_{\mathcal{T}}^* (q, u, v) \\ \text{si } (q_1, u_1, v_1) \vdash_{\mathcal{T}}^* (q_2, u_2, v_2) \text{ y } (q_2, u_2, v_2) \vdash_{\mathcal{T}} (q_3, u_3, v_3) : (q_1, u_1, v_1) \vdash_{\mathcal{T}}^* (q_3, u_3, v_3)}$$

Ejemplo 3.2



Función definida por un transductor. Dado un transductor \mathcal{T} , decimos que:

- ♦ \mathcal{T} entrega v con **input** u si existe una configuración inicial (q_0, u, ϵ) y una configuración final (q_f, ϵ, v) tal que:

$$(q_0, u, \epsilon) \vdash_{\mathcal{T}}^* (q_f, \epsilon, v)$$

- ♦ Se define la función $\llbracket \mathcal{T} \rrbracket : \Sigma^* \rightarrow 2^{\Omega^*}$:

$$\boxed{\llbracket \mathcal{T} \rrbracket(u) = \{v \in \Omega^* \mid \mathcal{T} \text{ entrega } v \text{ con input } u\}}$$

- ♦ Se dice que $f : \Sigma^* \rightarrow 2^{\Omega^*}$ es una **función racional** si existe un transductor \mathcal{T} tal que $f = \llbracket \mathcal{T} \rrbracket$.

- ♦ Un transductor define una función de palabras a conjuntos de palabras.

Dos interpretaciones para un transductor. Podemos ver a $\llbracket \mathcal{T} \rrbracket$ de dos formas:

1. \mathcal{T} define la **función** $\llbracket \mathcal{T} \rrbracket : \Sigma^* \rightarrow 2^{\Omega^*}$:

$$\llbracket \mathcal{T} \rrbracket(u) = \{v \in \Omega^* \mid \mathcal{T} \text{ entrega } v \text{ con input } u\}$$

2. \mathcal{T} define la **relación** $\llbracket \mathcal{T} \rrbracket \subseteq \Sigma^* \times \Omega^*$:

$$(u, v) \in \llbracket \mathcal{T} \rrbracket \quad \text{si, y solo si, } \mathcal{T} \text{ entrega } v \text{ con input } u$$

Desde ahora, hablaremos de función o relación **indistintamente** y hablaremos de las **relaciones racionales** (definidas por un transductor).

Lenguaje de input y output. Para una relación $R \subseteq \Sigma^* \times \Omega^*$ se define:

- ♦ $\pi_1(R) = \{u \in \Sigma^* \mid \exists v \in \Omega^*. (u, v) \in R\}$.
- ♦ $\pi_2(R) = \{v \in \Omega^* \mid \exists u \in \Sigma^*. (u, v) \in R\}$.

Teorema 12

Si \mathcal{T} es un transductor, entonces $\pi_1([\mathcal{T}])$ y $\pi_2([\mathcal{T}])$ son lenguajes regulares sobre Σ y Ω , respectivamente.

Idea demostración teorema 12. Para $\mathcal{T} = (Q, \Sigma, \Omega, \Delta, I, F)$, defina $\mathcal{A}_1 = (Q, \Sigma, \Delta_1, I, F)$ tal que

$$(p, a, q) \in \Delta_1 \quad \text{si, y solo si, } \quad \exists b \in \Omega \cup \{\epsilon\}. (p, a, b, q) \in \Delta$$

y demuestre que $\mathcal{L}(\mathcal{A}_1) = \pi_1([\mathcal{T}])$.

Teorema 13

Sea \mathcal{T}_1 y \mathcal{T}_2 dos transductores con Σ y Ω alfabetos de input y output. Las siguientes son **relaciones racionales**:

1.

$$[\mathcal{T}_1] \cup [\mathcal{T}_2] = \{(u, v) \in \Sigma^* \times \Omega^* \mid (u, v) \in [\mathcal{T}_1] \vee (u, v) \in [\mathcal{T}_2]\}$$

$$2. [\mathcal{T}_1] \cdot [\mathcal{T}_2] = \{(u_1 u_2, v_1 v_2) \in \Sigma^* \times \Omega^* \mid (u_1, v_1) \in [\mathcal{T}_1] \wedge (u_2, v_2) \in [\mathcal{T}_2]\}$$

$$3. [\mathcal{T}_1]^* = \bigcup_{k=0}^{\infty} [\mathcal{T}_1]^k$$

La demostración de este teorema queda como ejercicio propuesto al lector.

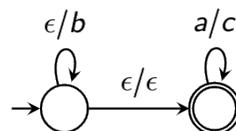
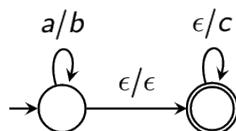
Teorema 14

Existen transductores \mathcal{T}_1 y \mathcal{T}_2 sobre Σ y Ω tal que:

$$[\mathcal{T}_1] \cap [\mathcal{T}_2] = \{(u, v) \in \Sigma^* \times \Omega^* \mid (u, v) \in [\mathcal{T}_1] \wedge (u, v) \in [\mathcal{T}_2]\}$$

NO es una relación racional.

Demostración teorema 14. Considere los siguientes transductores:



$$[\mathcal{T}_1] = \{(a^n, b^n c^m) \mid n \geq 0, m \geq 0\} \quad [\mathcal{T}_2] = \{(a^n, b^m c^n) \mid n \geq 0, m \geq 0\}$$

Vemos que $[\mathcal{T}_1] \cap [\mathcal{T}_2] = \{(a^n, b^n c^n) \mid n \geq 0\}$, pero el lenguaje que define el output no es regular, y por lo tanto $[\mathcal{T}_1] \cap [\mathcal{T}_2]$ no es racional. ■

Definición. Decimos que un transductor \mathcal{T} define una **función parcial** si:

para todo $u \in \Sigma^*$ se tiene que $|[\mathcal{T}](u)| \leq 1$.

Definición. Decimos que $\mathcal{T} = (Q, \Sigma, \Omega, \Delta, I, F)$ es **determinista** si cumple que:

1. \mathcal{T} define una función $\llbracket \mathcal{T} \rrbracket : \Sigma^* \rightarrow \Omega^*$.
2. para todo $(p, a_1, b_1, q_1) \in \Delta$ y $(p, a_2, b_2, q_2) \in \Delta$, si $a_1 = a_2$, entonces

$$b_1 = b_2 \quad \text{y} \quad q_1 = q_2$$

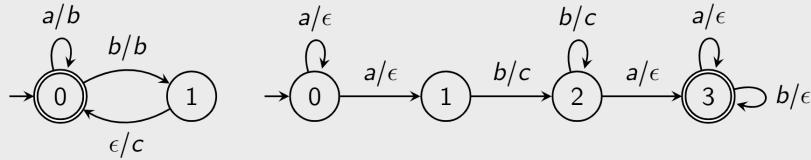
Es decir, si tenemos dos transiciones desde un estado p a otros dos estados q_1 y q_2 , y ambas transiciones leen lo mismo: $a_1 = a_2$, entonces lo que imprime el transductor y los estados deben ser iguales: $b_1 = b_2$ y $q_1 = q_2$. Básicamente, buscamos que sean la misma transición.

3. si $(p, \epsilon, b, q) \in \Delta$, entonces para todo $(p, a', b', q') \in \Delta$, se tiene que

$$(a', b', c') = (\epsilon, b, q)$$

Es decir, si tenemos una ϵ -transición, entonces es la única transición que sale desde p .

Ejemplo 3.3



3.3. Análisis léxico

PENDIENTE.

3.4. Algoritmo de Knuth-Morris-Prat

Pattern matching. Veamos el siguiente problema. Dado un **patrón** $w = w_1 \dots w_m$ y un documento $d = d_1 \dots d_n$, encontrar todas las posiciones donde aparece w en d , o sea, enumerar:

$$\{(i, j) \mid w = d_i d_{i+1} \dots d_j\}$$

Podríamos implementar la siguiente solución ingenua (y poco eficiente):

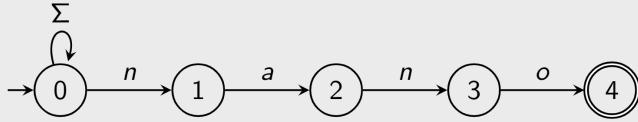
```

for  $i = 0$  to  $n - m$  do
|    $j := 1$ 
|   while  $j \leq m \wedge w_j = d_{i+j}$  do
|   |    $j := j + 1$ 
|   if  $j > m$  then
|   |   output  $(i + 1, i + m - 1)$ 
```

3.4.1. Autómata de un patrón

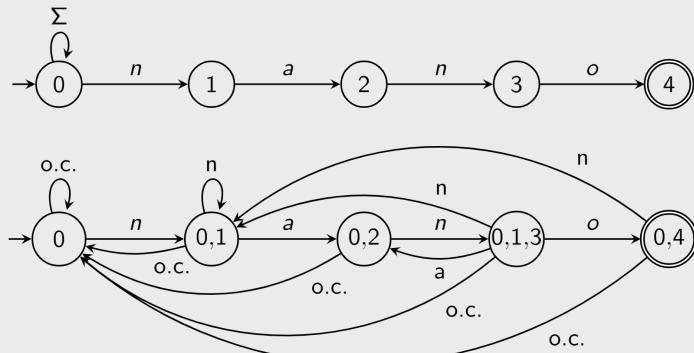
Definición. Dado una palabra $w = w_1 \dots w_m$, sea el NFA $\mathcal{A}_w = (Q, \Sigma, \Delta, I, F)$ tal que:

- ♦ $Q = \{0, 1, \dots, m\}$
- ♦ $\Delta = \{(0, a, 0) \mid a \in \Sigma\} \cup \{(i, w_{i+1}, i+1) \mid i < m\}$
- ♦ $I = \{0\}$ y $F = \{m\}$.

Ejemplo 3.4: Palabra $w = \text{nano}$ 

Podemos usar \mathcal{A}_w para encontrar todas las apariciones de w en w haciendo su **determinización**.

Definición. Sea $\mathcal{A}_w^{\det} = (Q^{\det}, \Sigma, \delta^{\det}, \{0\}, F^{\det})$ la determinización de \mathcal{A}_w tal que Q^{\det} contiene **solo los estados alcanzables** desde $\{0\}$.

Ejemplo 3.5: Palabra $w = \text{nano}$ y $d = \text{un nano no nana}$ 

q_0	q_0	q_1	q_0	q_1	q_2	$q_3 \checkmark$	q_4	q_0	q_1	q_0	q_1	q_2	q_3	q_2
u	n		n	a	n	o		n	o		n	a	n	a
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

Teorema 15

Para todo $S \in Q^{\det}$ y $i \in \{0, 1, \dots, m\}$ se cumple que

$$i \in S \quad \text{si, y solo si, } w_1 \dots w_i \text{ es un sufijo de } w_1 \dots w_{\max(S)}.$$

Corolarios. Tenemos que:

♦ Para todo $S_1, S_2 \in Q^{\det}$, si $\max(S_1) = \max(S_2)$, entonces $S_1 = S_2$.

♦ \mathcal{A}_w^{\det} tiene $|w| + 1$ estados y a lo más $\mathcal{O}(|w|^2)$ transiciones.

Por lo tanto, encontrar todos los substrings de w en d toma tiempo $\mathcal{O}(|d| + |w|^2)$.

Demostración teorema 15. Sea $S \in Q^{\det}$ un conjunto de estados cualquiera alcanzable desde $\{0\}$. Entonces, existe una palabra $u = a_1 \dots a_k$ tal que $\hat{\delta}^{\det}(\{0\}, u) = S$. Por la demostración de que $\mathcal{L}(\mathcal{A}^{\det}) = \mathcal{L}(\mathcal{A})$ para todo NFA \mathcal{A} , sabemos que $j \in S$ si, y sólo si, existe una ejecución \mathcal{A}_w sobre u :

$$0 = q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} \dots \xrightarrow{a_k} q_k = j$$

Por la definición de \mathcal{A}_w esta ejecución es de la forma:

$$0 \xrightarrow{a_1} 0 \xrightarrow{a_2} \dots \xrightarrow{a_{k-j}} 0 \underbrace{\xrightarrow{a_{k-j+1}} 1 \xrightarrow{a_{k-j+2}} 2 \dots \xrightarrow{a_k} j}_{w_1 \dots w_j}$$

Por lo tanto, $w_1 w_2 \dots w_j$ es sufijo de $a_1 \dots a_k$. Usaremos este último hecho para demostrar **ambas** direcciones del teorema, que formalizamos a continuación.

Propiedad. Para toda $u = a_1 \dots a_k$ tal que $\hat{\delta}^{\text{det}}(\{0\}, u) = S$, y para todo $j \leq m$:

$$j \in S \quad \text{si, y solo si, } w_1 \dots w_j \text{ es sufijo de } a_1 \dots a_k$$

Demostración (\Rightarrow). Como S es alcanzable desde $\{0\}$, entonces existe $u = a_1 \dots a_k$ tal que $\hat{\delta}^{\text{det}}(\{0\}, u) = S$. Como $\max(S) \in S$, entonces $W_1 \dots W_{\max(S)}$ es sufijo de $a_1 \dots a_k$.

Suponga que $i \in S$. Entonces $w_1 \dots w_i$ es sufijo de $a_1 \dots a_k$. Como $i \leq \max(S)$, entonces:

$$a_1 a_2 \dots a_{k-\max(S)} \overbrace{a_{k-\max(S)+1} \dots a_{k-i}}^{w_1 \dots w_{\max(S)}} \underbrace{a_{k-i+1} \dots a_k}_{w_{i+1} \dots w_k}$$

Por lo tanto, $w_1 \dots w_i$ es sufijo de $w_1 \dots w_{\max(S)}$.

Demostración (\Leftarrow). Como S es alcanzable desde $\{0\}$, entonces existe $u = a_1 \dots a_k$ tal que $\hat{\delta}^{\text{det}}(\{0\}, u) = S$. Como $\max(S) \in S$, entonces $w_1 \dots w_{\max(S)}$ es sufijo de $a_1 \dots a_k$.

Suponga que $w_1 \dots w_i$ es sufijo de $w_1 \dots w_{\max(S)}$. Como $w_1 \dots w_i$ es sufijo de $w_1 \dots w_{\max(S)}$ y $w_1 \dots w_{\max(S)}$ es sufijo de u , entonces $w_1 \dots w_i$ es sufijo de $u = a_1 \dots a_k$.

Por la “propiedad”, concluimos que $i \in S$. ■

3.4.2. Autómata finito con k -lookahead

Definición . Se definen los siguientes conjuntos de palabras sobre un alfabeto finito Σ :

- ♦ $\Sigma_\bullet = \Sigma^* \times \Sigma^*$
- ♦ $\Sigma_\bullet^k = \{(u, v) \in \Sigma_\bullet \mid |uv| = k\}$

Notación. En vez de $(u, v) \in \Sigma_\bullet$, escribiremos $u.v \in \Sigma_\bullet$.

Ejemplo 3.6

Si $\Sigma = \{a, b\}$, entonces:

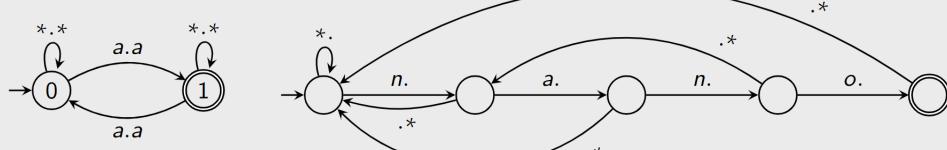
- ♦ $ab.ba \in \Sigma_\bullet$ y $.aba \in \Sigma_\bullet$
- ♦ $ab.ba \in \Sigma_\bullet^4$ y $.aba \in \Sigma_\bullet^3$

Definición. Un autómata finito determinista con k -lookahead es una tupla:

$$\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$$

- ♦ Q es un conjunto finito de estados.
- ♦ Σ es el alfabeto del input.
- ♦ q_0 es el estado inicial.
- ♦ $F \subseteq Q$ es el conjunto de estados iniciales.
- ♦ $\delta : Q \times (\Sigma \cup \{\$\})_\bullet^k \rightarrow Q$ es una función parcial, tal que:

$$\text{para todo } p \in Q \text{ y } w \in (\Sigma \cup \{\$\})_\bullet^k : |\{u.v \mid \delta(p, u \cdot v) = q \text{ y } uv = w\}| \leq 1$$

Ejemplo 3.7

Ejecución. Sea \mathcal{A} un DFA con k -lookahead. Tenemos que:

- ♦ Un par $(q, w) \in Q \times (\Sigma \cup \{\$\})^*$ es una **configuración** de \mathcal{A} .
- ♦ Una configuración $(q_0, w\$^k)$ es **inicial**.
- ♦ Una configuración $(q, \$^k)$ es **final** si $q \in F$.

El sufijo $\k nos sirve para marcar el final del input (y simplificar la definición de lookahead al leer el final de la palabra).

Se define la relación $\vdash_{\mathcal{A}}$ de **siguiente-paso** entre configuraciones de \mathcal{A} :

$$(p_1, w_1) \vdash_{\mathcal{A}} (p_2, w_2)$$

si, y sólo si, $\delta(p_1, u.v) = p_2$ y existe $w \in \Sigma^*$ tal que $w_1 = uw$ y $w_2 = vw$.

Se define $\vdash_{\mathcal{A}}^*$ como la clausura **refleja** y transitiva de $\vdash_{\mathcal{A}}$:

$$\boxed{\begin{array}{l} \text{para toda configuración } (p, w) : (p, w) \vdash_{\mathcal{A}}^* (p, w) \\ \text{si } (p_1, w_1) \vdash_{\mathcal{A}}^* \text{ y } (p_2, w_2) \vdash_{\mathcal{A}} (p_3, w_3) : (p_1, w_1) \vdash_{\mathcal{A}}^* (p_3, w_3) \end{array}}$$

Decimos que $(p, u) \vdash_{\mathcal{A}}^* (q, v)$ si uno puede ir de (p, u) a (q, v) en **0 o más pasos**.

Aceptación. Decimos que \mathcal{A} **acepta** w si existe una configuración inicial $(q_0, w\$^k)$ y una configuración final $(q_f, \$^k)$ tal que:

$$(q_0, w\$^k) \vdash_{\mathcal{A}}^* (q_f, \$^k)$$

El **lenguaje aceptado** por \mathcal{A} se define como:

$$\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \mathcal{A} \text{ acepta } w\}$$

Vemos que son las mismas definiciones para un ϵ -NFA.

Teorema 16

Para todo DFA con k -lookahead \mathcal{A} se tiene que $\mathcal{L}(\mathcal{A})$ es un **lenguaje regular**.

La demostración de este teorema queda como ejercicio propuesto al lector.

Definición. Llamaremos un **lazy autómata** a un DFA con 1-lookahead.

3.4.3. Algoritmo KMP

Construcción de un lazy automata. Sea $w = w_1 \dots w_m$ y $\mathcal{A}_w^{\det} = (Q^{\det}, \Sigma, \delta^{\det}, \{0\}, F^{\det})$ la determinización de \mathcal{A}_w .

Definición. Para $i \in [0, m]$, sea S_i el **único estado** en Q^{det} tal que $i = \max(S_i)$.

Propiedad. Para todo $a \in \{w_1, \dots, w_m\}$ y $i \in [0, m - 1]$:

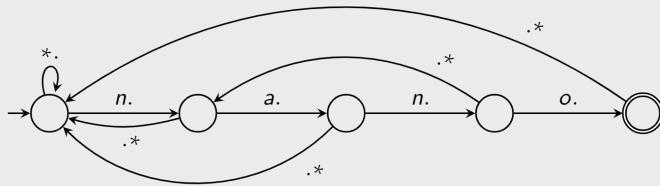
1. $S_i - \{i\} \in Q^{\text{det}}$
2. $a = w_{i+1}$, entonces $\delta^{\text{det}}(S_i, a) = S_{i+1}$.
3. $a \neq w_{i+1}$, entonces $\delta^{\text{det}}(S_i, a) = \delta^{\text{det}}(S_i - \{i\}, a)$.

La demostración de esta propiedad queda como ejercicio propuesto al lector.

En base a la propiedad anterior, se define el lazy autómata $\mathcal{A}_w^{\text{lazy}} = (Q^{\text{det}}, \Sigma, \delta^{\text{lazy}}, \{0\}, F^{\text{det}})$ tal que:

- ♦ para todo $a \neq w_1$: $\delta^{\text{lazy}}(\{0\}, a) = \{0\}$.
- ♦ para todo $a \in \{w_1, \dots, w_m\}$ y $i \in [0, m - 1]$:
 - si $a = w_{i+1}$, entonces $\delta^{\text{lazy}}(S_i, a) = S_{i+1}$
 - si $a \neq w_{i+1}$ y $i \neq 0$, entonces $\delta^{\text{lazy}}(S_i, a) = S_i - \{i\}$.

Ejemplo 3.8



Teorema 17

Para todo w se cumple que $\mathcal{L}(\mathcal{A}_w^{\text{det}}) = \mathcal{L}(\mathcal{A}_w^{\text{lazy}})$.

La demostración queda como ejercicio propuesto para el lector (usando la propiedad).

Complejidad KMP. Vemos que:

- ♦ El número de pasos que $\mathcal{A}_w^{\text{lazy}}$ **consume** letras = $|d|$
- ♦ El número de pasos que $\mathcal{A}_w^{\text{lazy}}$ **retrocede** $\leq |d|$
- ♦ El número de **pasos totales** de $\mathcal{A}_w^{\text{lazy}}$ $\leq 2 \cdot |d|$

Por lo tanto, la cantidad de pasos es **lineal** en $\mathcal{O}(|d|)$.

Algoritmo KMP. Dado una palabra w y un documento d :

1. Construimos $\mathcal{A}_w^{\text{lazy}}$ desde \mathcal{A}_w .
2. Ejecutamos $\mathcal{A}_w^{\text{lazy}}$ sobre d .

El paso 1 toma $\mathcal{O}(|w|)$ y el paso 2 toma $\mathcal{O}(|d|)$, por lo tanto, el tiempo del algoritmo es $\mathcal{O}(|w| + |d|)$.

Queda como ejercicio para el lector demostrar que construir $\mathcal{A}_w^{\text{lazy}}$ toma tiempo $\mathcal{O}(|w|)$.

4. Lenguajes libres de contexto

4.1. Gramáticas libres de contexto

4.1.1. Gramáticas

Definición. Una **gramática libre de contexto** (CFG) es una tupla:

$$\boxed{\mathcal{G} = (V, \Sigma, P, S)}$$

- ♦ V es un conjunto finito de **variables** o **no-terminales**.
- ♦ Σ es el alfabeto finito (o **terminales**) tal que $\Sigma \cap V = \emptyset$.
- ♦ $P \subseteq V \times (V \cup \Sigma)^*$ es un subconjunto finito de **reglas** o **producciones**.
- ♦ $S \in V$ es la **variable inicial**.

Ejemplo 4.1

Considere la gramática $\mathcal{G} = (V, \Sigma, P, S)$ tal que:

- ♦ $V = \{X, Y\}$
- ♦ $\Sigma = \{a, b\}$
- ♦ $\{(X, aXb), (X, Y), (Y, \epsilon)\}$
- ♦ $S = X$

$$\begin{aligned}\mathcal{G} : \quad X &\rightarrow aXb \\ &\quad X \rightarrow Y \\ &\quad Y \rightarrow \epsilon\end{aligned}$$

Notación. En este texto:

- ♦ Para las **variables** en una gramática usaremos letras mayúsculas: X, Y, Z, A, B, C, \dots
- ♦ Para los **terminales** en una gramática usaremos letras minúsculas: a, b, c, \dots
- ♦ Para palabras en $(V \cup \Sigma)^*$ usaremos símbolos: $\alpha, \beta, \gamma, \dots$
- ♦ Para una producción $(A, \alpha) \in P$ la escribimos como: $A \rightarrow \alpha$

Simplificación. Si tenemos un conjunto de reglas de la forma:

$$\begin{aligned}X &\rightarrow \alpha_1 \\ X &\rightarrow \alpha_2 \\ &\dots \\ X &\rightarrow \alpha_n\end{aligned}$$

entonces escribimos estas reglas **sucintamente** como

$$X \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

Recordando que $\alpha_1, \alpha_2, \dots, \alpha_n \in (V \cup \Sigma)^*$.

Ejemplo 4.2

La gramática del ejemplo anterior:

$$\begin{aligned}\mathcal{G} : \quad X &\rightarrow aXb \\ &X \rightarrow Y \\ &Y \rightarrow \epsilon\end{aligned}$$

Podemos escribirla en notación **sucinta** como:

$$\begin{aligned}\mathcal{G} : \quad X &\rightarrow aXb \mid Y \\ &Y \rightarrow \epsilon\end{aligned}$$

Producciones. Sea \mathcal{G} una CFG. Definimos la relación $\Rightarrow \subseteq (V \cup \Sigma)^* \times (V \cup \Sigma)^*$ de **producción** tal que:

$$\alpha \cdot X \cdot \beta \Rightarrow \alpha \cdot \gamma \cdot \beta \quad \text{si, y solo si, } (X \rightarrow \gamma) \in P$$

para todo $X \in V$ y $\alpha, \beta, \gamma \in (V \cup \Sigma)^*$.

Si $\alpha X \beta \Rightarrow \alpha \gamma \beta$ entonces decimos que:

- ♦ $\alpha X \beta$ **produce** $\alpha \gamma \beta$ o
- ♦ $\alpha \gamma \beta$ **es producible** desde $\alpha X \beta$.
- ♦ $\alpha X \beta \Rightarrow \alpha \gamma \beta$ es **reemplazar** γ en X en la palabra $\alpha X \beta$.

Derivaciones. Sea \mathcal{G} una CFG. Dadas dos palabras $\alpha, \beta \in (V \cup \Sigma)^*$ decimos que α **deriva** β :

$$\alpha \xrightarrow{*} \beta$$

si existe $\alpha_1, \alpha_2, \dots, \alpha_n \in (V \cup \Sigma)^*$ tal que: $\alpha \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \beta$, con \Rightarrow^* la **clausura refleja y transitiva** de \Rightarrow , esto es:

1. $\alpha \xrightarrow{*} \alpha$
2. $\alpha \xrightarrow{*} \beta$ si, y sólo si, existe γ tal que $\alpha \xrightarrow{*} \gamma$ y $\gamma \Rightarrow \beta$.

para todo $\alpha, \beta \in (V \cup \Sigma)^*$. Notemos que \Rightarrow y \Rightarrow^* son relaciones entre palabras en $(V \cup \Sigma)^*$.

Lenguaje. Sea \mathcal{G} una CFG. El **lenguaje** de una gramática \mathcal{G} se define como:

$$\mathcal{L}(\mathcal{G}) = \left\{ w \in \Sigma^* \mid S \xrightarrow{*} w \right\}$$

$\mathcal{L}(\mathcal{G})$ son todas las palabras en Σ^* que se pueden derivar desde S .

Ejemplo 4.3

Sea \mathcal{G} una CFG tal que:

$$\begin{aligned}\mathcal{G} : \quad X &\rightarrow aXb \mid Y \\ Y &\rightarrow \epsilon\end{aligned}$$

- ♦ Como $X \xrightarrow{*} aaabbb$, entonces $aaabbb \in \mathcal{L}(\mathcal{G})$.
- ♦ En general, uno puede demostrar por **inducción** que:

$$\mathcal{L}(\mathcal{G}) = \{a^n b^n \mid n \geq 0\}$$

Lenguaje libre de contexto. Diremos que $L \subseteq \Sigma^*$ es un **lenguaje libre de contexto** si, y sólo si, existe una gramática libre de contexto \mathcal{G} tal que:

$$L = \mathcal{L}(\mathcal{G})$$

Ejemplo 4.4

Los siguientes son lenguajes libres de contexto:

- ♦ $L = \{a^n b^n \mid n \geq 0\}$
- ♦ Par = $\{w \in \{a, b\}^* \mid w \text{ tiene largo par}\}$
- ♦ Pal = $\{w \in \{a, b\}^* \mid w = w^{\text{rev}}\}$

4.1.2. Árboles y derivaciones

4.2. Simplificación de gramáticas

4.3. Forma normal de Chomsky

4.4. Lema de bombeo para lenguajes libres de contexto

4.5. Algoritmo CKY

5. Algoritmos para lenguajes libres de contexto

5.1. Autómatas apiladores

5.1.1. Versión normal

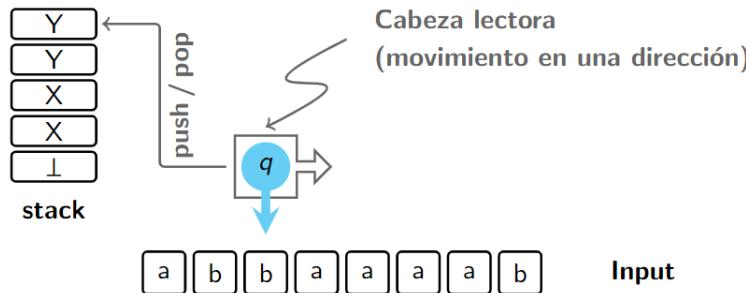


Figura 8: Idea de un autómata apilador

Definición. Un autómata apilador (*PushDown Automata*, PDA) es una estructura:

$$\mathcal{P} = (Q, \Sigma, \Gamma, \Delta, q_0, \perp, F)$$

- ♦ Q es un conjunto finito de **estados**.
- ♦ Σ es el alfabeto del **input**.
- ♦ $q_0 \in Q$ es el estado **inicial**.
- ♦ F es el conjunto de estados **finales**.
- ♦ Γ es el alfabeto de **stack**.
- ♦ $\perp \in \Gamma$ es el símbolo **inicial del stack** (fondo).
- ♦ $\Delta \subseteq (Q \times (\Sigma \cup \{\epsilon\}) \times \Gamma) \times (Q \times \Gamma^*)$ es una relación finita de transición.

Intuitivamente, la transición:

$$\left((p, a, A), (q, B_1 B_2 \cdots B_k) \right) \in \Delta$$

si el autómata apilador está:

- ♦ en el estado p , leyendo a , y en el tope del stack hay una A ,
- entonces:
- ♦ cambia al estado q , y modifíco el tope A por $B_1 B_2 \cdots B_k$.

Intuitivamente, la transición **en vacío**:

$$\left((p, \epsilon, A), (q, B_1 B_2 \cdots B_k) \right) \in \Delta$$

si el autómata apilador está:

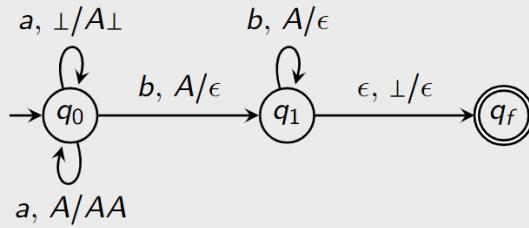
- ♦ en el estado p , *sin lectura de una letra*, y en el tope del stack hay una A ,
- entonces:
- ♦ cambia al estado q , y modifíco el tope A por $B_1 B_2 \cdots B_k$.

Ejemplo 5.1

$$\mathcal{P} = (Q, \Sigma, \Gamma, \Delta, q_0, \perp, \{q_f\})$$

- ♦ $Q = \{q_0, q_1, q_f\}$, $\Sigma = \{a, b\}$, $\Gamma = \{A, \perp\}$ y Δ :

$$\begin{array}{ll}
 (q_0, a, \perp, q_0, A \perp) & q_0 \perp \xrightarrow{a} q_0 A \perp \\
 (q_0, a, A, q_0, AA) & q_0 A \xrightarrow{a} q_0 AA \\
 (q_0, b, A, q_1, \epsilon) & q_0 A \xrightarrow{b} q_1 \\
 (q_1, b, A, q_1, \epsilon) & q_1 A \xrightarrow{b} q_1 \\
 (q_1, \epsilon, \perp, q_f, \epsilon) & q_1 \perp \xrightarrow{\epsilon} q_f
 \end{array}$$



Notación. Dada una palabra $A_1 A_2 \dots A_k \in \Gamma^+$ decimos que:

- ♦ $A_1 A_2 \dots A_k$ es un stack (contenido),
- ♦ A_1 es el **tope** del stack y
- ♦ $A_2 \dots A_k$ es la **cola** del stack.

Definición. Una **configuración** de \mathcal{P} es una tupla $(q \cdot \gamma, w) \in (Q \cdot \Gamma^*, \Sigma^*)$ tal que:

- ♦ q es el estado actual.
- ♦ γ es el contenido del stack.
- ♦ w es el contenido del input.

Decimos que una configuración:

$$(q \cdot \gamma, w) \in (Q \cdot \Gamma^*, \Sigma^*)$$

- ♦ es **inicial** si $q \cdot \gamma = q_0 \cdot \perp$.
- ♦ es **final** si $q \cdot \gamma = q_f \cdot \epsilon$ con $q_f \in F$ y $w = \epsilon$.

Definición. Se define la relación $\vdash_{\mathcal{P}}$ de **siguiente-paso** entre configuraciones de \mathcal{P} :

$$(q_1 \cdot \gamma_1, w_1) \vdash_{\mathcal{P}} (q_2 \cdot \gamma_2, w_2)$$

si, y sólo si, existe una transición $(q_1, a, A, q_2, \alpha) \in \Delta$ y $\gamma \in \Gamma^*$ tal que:

- ♦ $w_1 = a \cdot w_2$

$$\diamond \gamma_1 = A \cdot \gamma$$

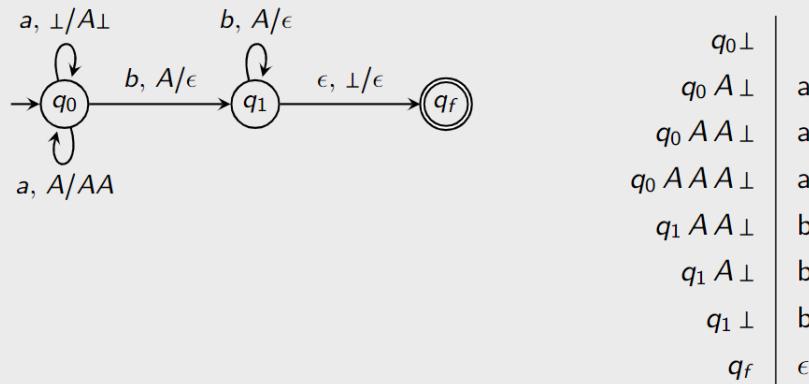
$$\diamond \gamma_2 = \alpha \cdot \gamma$$

Se define $\vdash_{\mathcal{P}}^*$ como la clausura **refleja** y **transitiva** de $\vdash_{\mathcal{P}}$. En otras palabras:

$(q_1\gamma_1, w_1) \vdash_{\mathcal{P}}^* (q_2\gamma_2, w_2)$ si uno puede ir de $(q_1\gamma_1, w_1)$ a $(q_2\gamma_2, w_2)$ en 0 o más pasos.

Ejemplo 5.2

Para la palabra $w = aaabbb$, tenemos la ejecución:



Definiciones. \mathcal{P} **acepta** w si, y sólo si, $(q_0 \perp, w) \vdash_{\mathcal{P}}^* (q_f, \epsilon)$ para algún $q_f \in F$.
El **lenguaje aceptado** por \mathcal{P} se define como:

$$\mathcal{L}(\mathcal{P}) = \{w \in \Sigma^* \mid \mathcal{P} \text{ acepta } w\}$$

Ejemplo 5.3

El lenguaje aceptado por el PDA utilizado en los ejemplos anteriores es $\mathcal{L}(\mathcal{P}) = \{a^n b^n \mid n \geq 0\}$.

5.1.2. Versión alternativa

Esta definición de autómata apilador es poco común pero trae algunas ventajas:

- ♦ Es un modelo que ayuda a entender mejor los algoritmos de evaluación para gramáticas.
- ♦ Es un modelo menos estándar pero mucho más sencillo.
- ♦ Al prof Cristian le gustó y lo encontró interesante.

Definición. Un **PDA alternativo** es una estructura:

$$\mathcal{D} = (Q, \Sigma, \Delta, q_0, F)$$

- ♦ Q es un conjunto finito de **estados**.

- ♦ Σ es el alfabeto del **input**.
- ♦ $q_0 \in Q$ es el estado **inicial**.
- ♦ F es el conjunto de estados **finales**.
- ♦ $\Delta \subseteq Q^+ \times (\Sigma \cup \{\epsilon\}) \times Q^*$ es una **relación finita de transición**.

Intuitivamente, la transición:

$$\boxed{(A_1 \dots A_i, a, B_1 \dots B_j) \in \Delta}$$

si el autómata apilador tiene:

- ♦ $A_1 \dots A_i$ en el tope del stack y leyendo a ,

entonces:

- ♦ cambia el tope $A_1 \dots A_i$ por $B_1 \dots B_j$.

En este tipo de autómata apilador, **no hay diferencia** entre estados y alfabeto del stack.

Definición. Una **configuración** de \mathcal{D} es una tupla

$$\boxed{(q_1 \dots q_k, w) \in (Q^+, \Sigma^*)}$$

tal que:

- ♦ $q_1 \dots q_k$ es el contenido del stack con q_1 el tope del stack.
- ♦ w es el contenido del input.

Decimos que una configuración:

- ♦ (q_0, w) es **inicial**.
- ♦ (Q_f, ϵ) es **final** si $q_f \in F$.

Definición. Se define la relación $\vdash_{\mathcal{D}}$ de **siguiente-paso** entre configuraciones de \mathcal{D} :

$$\boxed{(\gamma_1, w_1) \quad \vdash_{\mathcal{D}} \quad (\gamma_2, w_2)}$$

si, y sólo si, existe una transición $(\alpha, a, \beta) \in \Delta$ y $\gamma \in \Gamma^*$ tal que:

- ♦ $w_1 = a \cdot w_2$
- ♦ $\gamma_1 = \alpha \cdot \gamma$
- ♦ $\gamma_2 = \beta \cdot \gamma$

Se define $\vdash_{\mathcal{D}}^*$ como la clausura **refleja** y **transitiva** de $\vdash_{\mathcal{D}}$.

Definiciones. \mathcal{D} **acepta** w si, y sólo si, $(q_0, w) \vdash_{\mathcal{D}}^* (q_f, \epsilon)$ para algún $q_f \in F$. Además, el **lenguaje aceptado** por \mathcal{D} se define como:

$$\boxed{\mathcal{L}(\mathcal{D}) = \{w \in \Sigma^* \mid \mathcal{D} \text{ acepta } w\}}$$

Ejemplo 5.4

$$\mathcal{D} = (Q, \{a, b\}, \Delta, q_0, F)$$

- ♦ $Q = \{\perp, q_0, q_1, q_f\}$ y Δ :

$(\perp, a, q_0 \perp)$	$\perp \xrightarrow{a} q_0 \perp$	\perp	$q_0 \perp$ a $q_0 q_0 \perp$ a $q_0 q_0 q_0 \perp$ a $q_1 q_0 q_0 \perp$ b $q_1 q_0 \perp$ b $q_1 \perp$ b q_f ϵ
$(q_0, a, q_0 q_0)$	$q_0 \xrightarrow{a} q_0 q_0$	$q_0 \perp$	
(q_0, b, q_1)	$q_0 \xrightarrow{a} q_1$	$q_0 q_0 \perp$	
$(q_1 q_0, b, q_1)$	$q_1 q_0 \xrightarrow{b} q_1$	$q_0 q_0 q_0 \perp$	
$(q_1 \perp, \epsilon, q_f)$	$q_1 \perp \xrightarrow{b} q_f$	$q_1 q_0 q_0 \perp$	

$$\mathcal{L}(\mathcal{D}) = \{a^n b^n \mid n \geq 1\}$$

Teorema 18

Para todo autómata apilador \mathcal{P} existe un autómata apilador alternativo \mathcal{D} , y viceversa, tal que:

$$\mathcal{L}(\mathcal{P}) = \mathcal{L}(\mathcal{D})$$

El teorema anterior nos dice que podemos usar ambos modelos de manera **equivalente**.

5.2. Autómatas apiladores vs gramáticas libres de contexto

¿En qué se parecen CFG a PDA?

Gramáticas libre de contexto	Autómatas apiladores
■ Terminales (Σ)	■ Alfabeto input (Σ)
■ Variables (V)	■ Alfabeto stack (Γ)
■ Producciones	■ Transiciones
■ $A \rightarrow \gamma$	■ $pA \xrightarrow{c} q\gamma$
■ $A \rightarrow a$	■ $pA \xrightarrow{c} q$
■ Derivaciones	■ Ejecuciones
■ Árbol de derivación	■ Stack

Figura 9: Gramáticas vs Autómatas apiladores

Teorema 19

Todo lenguaje libre de contexto puede ser descrito equivalentemente por:

- ♦ Una gramática libre de contexto (**CFG**).
- ♦ Un autómata apilador (**PDA**).

5.2.1. Desde CFG a PDA

Partimos enunciado un teorema:

Teorema 20

Para toda gramática libre de contexto \mathcal{G} , existe un **autómata apilador alternativo** \mathcal{D} , tal que:

$$\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{D})$$

Construcción \mathcal{D} desde \mathcal{G} . Sea $\mathcal{G} = (V, \Sigma, P, S)$ una CFG. Construimos un PDA alternativo \mathcal{D} que acepta $\mathcal{L}(\mathcal{G})$:

$$\mathcal{D} = \left(V \cup \Sigma \cup \{q_0, q_f\}, \Sigma, \Delta, q_0, \{q_f\} \right)$$

La relación de transición Δ se define como:

$$\begin{aligned} \Delta &= \{(q_0, \epsilon, S \cdot q_f)\} \quad \cup \\ &\quad \{(X, \epsilon, \gamma) \mid X \rightarrow \gamma \in P\} \quad \cup \quad (\textbf{Expandir}) \\ &\quad \{(a, a, \epsilon) \mid a \in \Sigma\} \quad \quad \quad (\textbf{Reducir}) \end{aligned}$$

Demostración $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{D})$. Debemos demostrar dos direcciones: $\mathcal{L}(\mathcal{G}) \subseteq \mathcal{L}(\mathcal{D})$ y $\mathcal{L}(\mathcal{D}) \subseteq \mathcal{L}(\mathcal{G})$.

Demostración $\mathcal{L}(\mathcal{G}) \subseteq \mathcal{L}(\mathcal{D})$. Para cada $w \in \mathcal{L}(\mathcal{G})$ debemos encontrar una ejecución de aceptación de \mathcal{D} sobre w . ¿Cómo encontramos esta ejecución? La idea es que para cada árbol de derivación \mathcal{T} de \mathcal{G} sobre w , construimos una ejecución de \mathcal{D} sobre w que recorre el árbol \mathcal{T} en profundidad (DFS). Por tanto, debemos usar **inducción** sobre la altura del árbol \mathcal{T} .

Hipótesis de inducción. Para todo árbol de derivación \mathcal{T} de \mathcal{G} con **altura** h tal que:

♦ la raíz de \mathcal{T} es X , y

♦ \mathcal{T} produce la palabra w

entonces $(X \cdot \gamma, w) \vdash_{\mathcal{D}}^* (\gamma, \epsilon)$ para todo $\gamma \in Q^+$.

Caso base: $h = 1$. Si \mathcal{T} tiene altura 1, entonces:

♦ \mathcal{T} produce la palabra $w = a$ para algún $a \in \Sigma$ y

♦ \mathcal{T} consiste de un nodo X y un hijo a con $X \rightarrow a$.

Entonces para todo $\gamma \in Q^+$:

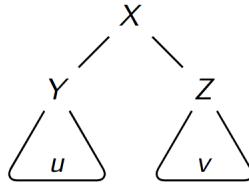
$$(X \cdot \gamma, a) \vdash_{\mathcal{D}} (a \cdot \gamma, a) \vdash_{\mathcal{D}} (\gamma, \epsilon)$$

es una ejecución de \mathcal{D} sobre a .

Caso inductivo: $h = n$. Suponemos que el árbol de derivación \mathcal{T} de \mathcal{G} tiene **altura** n tal que:

- ♦ la raíz de \mathcal{T} es X , y
- ♦ \mathcal{T} produce la palabra w .

Sin pérdida de generalidad, suponga que \mathcal{T} es de la forma:



donde $w = u \cdot v$ y $X \rightarrow YZ$. Por HI, se tiene que para todo $\gamma_1, \gamma_2 \in Q^+$:

$$\begin{aligned} (Y \cdot \gamma_1, u) &\vdash_{\mathcal{D}}^* (\gamma_1, \epsilon) \\ (Z \cdot \gamma_2, v) &\vdash_{\mathcal{D}}^* (\gamma_2, \epsilon) \end{aligned}$$

Para $\gamma \in Q^+$ construimos la siguiente ejecución de \mathcal{D} sobre $w = uv$:

$$(X \cdot \gamma, uv) \vdash_{\mathcal{D}} (YZ \cdot \gamma, uv) \vdash_{\mathcal{D}}^* (Z \cdot \gamma, v) \vdash_{\mathcal{D}}^* (\gamma, \epsilon)$$

■

La demostración de $\mathcal{L}(\mathcal{D}) \subseteq \mathcal{L}(\mathcal{G})$ se deja como ejercicio propuesto al lector.

5.2.2. Desde PDA a CFG

Partimos enunciando el siguiente teorema:

Teorema 21

Para todo autómata apilador \mathcal{P} , existe una gramática libre de contexto \mathcal{G} tal que:

$$\mathcal{L}(\mathcal{P}) = \mathcal{L}(\mathcal{G})$$

Demostración $\mathcal{L}(\mathcal{P}) = \mathcal{L}(\mathcal{G})$. Sea $\mathcal{P} = (Q, \Sigma, \Gamma, \Delta, q_0, \perp, F)$ un PDA (normal). Los pasos a seguir son:

1. Convertir \mathcal{P} a un PDA \mathcal{P}' con **UN SOLO ESTADO**.
2. Convertir \mathcal{P}' a una gramática libre de contexto \mathcal{G} .

Paso 1. Sea $\mathcal{P} = (Q, \Sigma, \Gamma, \Delta, q_0, \perp, F)$ un PDA. Podemos analizar:

- ♦ ¿Por qué NO necesitamos la información de los estados?
- ♦ ¿Cómo guardamos la información de los estados en el stack?

Esto conlleva a la siguiente pregunta: *Si el PDA está en el estado p y en el tope del stack hay una A , ¿a cuál estado llegaré al remover A del stack?*

La solución a esta pregunta es que podemos **adivinar** (no-determinismo) el estado que vamos a llegar cuando removamos A del stack.

Sin pérdida de generalidad, podemos asumir que

1. Todas las transiciones son de la forma:

$$qA \xrightarrow{c} pB_1B_2 \quad \text{o} \quad qA \xrightarrow{c} p\epsilon$$

con $c \in (\Sigma \cup \{\epsilon\})$.

2. Existe $q_f \in Q$ tal que si $w \in \mathcal{L}(\mathcal{P})$ entonces:

$$(q_0\perp, w) \vdash_{\mathcal{D}}^* (q_f, \epsilon)$$

Estos dos puntos nos aseguran que siempre llegamos al **mismo estado** q_f . Luego, construimos el autómata apilador \mathcal{P}' con **un solo estado**:

$$\mathcal{P}' = (\{q\}, \Sigma, \Gamma', \Delta', \{q\}, \perp', \{q\})$$

- ♦ $\Gamma' = Q \times \Gamma \times Q$.

“($p, A, q) \in \Gamma'$ si desde p leyendo A en el tope del stack llegamos a q al hacer pop de A ”.

- ♦ $\perp' = (q_0, \perp, q_f)$.

“El autómata parte en q_0 y al hacer pop de \perp llegará a q_f ”.

- ♦ Si $pA \xrightarrow{c} p'B_1B_2 \in \Delta$ con $c \in (\Sigma \cup \{\epsilon\})$, entonces **para todo** $p_1, p_2 \in Q$:

$$q(p, A, p_2) \xrightarrow{c} q(p', B_1, p_1) (p_1, B_2, p_2) \in \Delta'$$

- ♦ Si $pA \xrightarrow{c} p' \in \Delta$ con $c \in (\Sigma \cup \{\epsilon\})$, entonces:

$$q(p, A, p') \xrightarrow{c} q \in \Delta'$$

Hipótesis de inducción (en el número de pasos n). Para todo $p, p' \in Q$, $A \in \Gamma$ y $w \in \Sigma^*$ se cumple que:

$$(pA, w) \vdash_{\mathcal{P}}^n (p', \epsilon) \quad \text{si, y solo si,} \quad (q(p, A, p'), w) \vdash_{\mathcal{P}'}^n (q, \epsilon)$$

donde $\vdash_{\mathcal{P}}^n$ es la relación de **siguiente-paso** de \mathcal{P} n -veces.

Si demostramos esta hipótesis, habremos demostrado que $\mathcal{L}(\mathcal{P}) = \mathcal{L}(\mathcal{P}')$. ¿Por qué?

Caso base: $n = 1$. Para todo $p, p' \in Q$, y $A \in \Gamma$ se cumple que:

$$(pA, c) \vdash_{\mathcal{P}} (p', \epsilon) \quad \text{si, y solo si,} \quad (q(p, A, p'), c) \vdash_{\mathcal{P}'} (q, \epsilon)$$

para todo $c \in (\Sigma \cup \{\epsilon\})$.

Caso inductivo. **Sin pérdida de generalidad**, suponga que $pA \xrightarrow{a} p_1A_1A_2$ y $w = auv$, entonces

$$\begin{aligned} (pA, \underbrace{auv}_w) \vdash_{\mathcal{P}}^n (p', \epsilon) \quad &\text{ssi } (pA, auv) \vdash_{\mathcal{P}} (p_1A_1A_2, uv) \vdash_{\mathcal{P}}^i (p_2A_2, v) \vdash_{\mathcal{P}}^j (p', \epsilon) \\ &\text{ssi } (p_1A_1, u) \vdash_{\mathcal{P}}^i (p_2, \epsilon) \quad \text{y} \quad (p_2A_2, v) \vdash_{\mathcal{P}}^j (p', \epsilon) \\ &\text{ssi } (q(p_1, A_1, p_2), u) \vdash_{\mathcal{P}'}^i (q, \epsilon) \quad \text{y} \quad (q(p_2, A_2, p'), v) \vdash_{\mathcal{P}'}^j (q, \epsilon) \\ &\text{ssi } (q(p, A, p'), auv) \vdash_{\mathcal{P}} (q(p_1, A_1, p_2)(p_2, A_2, q)), uv) \vdash_{\mathcal{P}}^{i+j} (q, \epsilon) \end{aligned}$$

■

Paso 2. Sea $\mathcal{P} = (\{q\}, \Sigma, \Gamma, \Delta, q, \perp, \{q\})$ un PDA con **UN solo estado**. Construimos la gramática:

$$\mathcal{G} = (V, \Sigma, P, \perp)$$

- ♦ $V = \Gamma$.
- ♦ Si $qA \xrightarrow{\epsilon} q\alpha \in \Delta$ entonces $A \rightarrow \alpha \in P$
- ♦ Si $qA \xrightarrow{a} q\alpha \in \Delta$ entonces $A \rightarrow a\alpha \in P$

La demostración de este paso queda como ejercicio propuesto al lector.

5.3. Parsing: cómputo de First y Follow

Recordatorio. La **sintaxis** de un lenguaje es un conjunto de reglas que describen los programas válidos que tienen significado. Por otro lado, la **semántica** de un lenguaje define el significado de un programa correcto según la sintaxis.

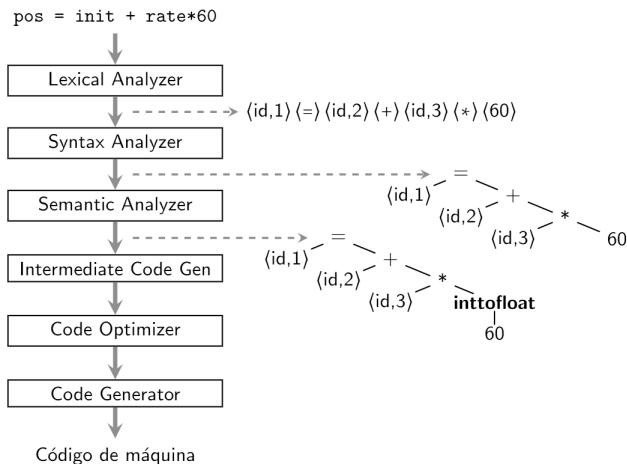


Figura 10: La estructura de un compilador

Lo que se busca es un proceso de **verificación de sintaxis** de un programa, y que entregue la estructura del mismo (árbol de parsing). Consta de tres etapas:

1. Análisis léxico (**Lexer**).
2. Análisis sintático (**Parser**).
3. Análisis semántico.

En una sección anterior vimos el **Lexer**. Ahora, veremos como hacer el **Parser**.

Informalmente: “*dado una secuencia de tokens w' y una gramática \mathcal{G} , construir un árbol de derivación (parsing) de \mathcal{G} para w'* ”.

Con el **árbol de derivación** habremos verificado la sintaxis y obtenido la estructura.

Ejemplo 5.5: Parsing de gramática

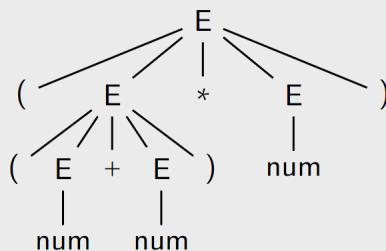
$$E \rightarrow (E + E) \mid (E * E) \mid \text{num}$$

Para un input $w = ((43 + 56) * 27)$:

- ♦ Convertimos w en una secuencia de **tokens**:

$$w' = ((\text{num} + \text{num}) * \text{num})$$

- ♦ Construimos un árbol de **parsing** para w' :



Problema de parsing. Dado una palabra w y dado una gramática \mathcal{G} , generar un árbol de parsing \mathcal{T} de \mathcal{G} para w . ¿Ya sabemos resolver este problema? El algoritmo CKY nos permite hacer esto, pero:

- ♦ es impracticable para grandes inputs.
- ♦ múltiples pasadas sobre el input.

Deseamos hacer parsing en **tiempo lineal** en el tamaño del input. ¿Quién nos puede rescatar ante tal problema? Efectivamente, los autómatas apiladores.

Recordemos que, para una gramática $\mathcal{G} = (V, \Sigma, P, S)$ **podemos construir un PDA alternativo** \mathcal{D} que acepta $\mathcal{L}(\mathcal{G})$:

$$\mathcal{D} = \left(V \cup \Sigma \cup \{q_0, q_f\}, \Sigma, \Delta, q_0, \{q_f\} \right)$$

La relación de transición Δ se define como:

$$\begin{aligned} \Delta &= \{(q_0, \epsilon, S \cdot q_f)\} \quad \cup \\ &\quad \{(X, \epsilon, \gamma) \mid X \rightarrow \gamma \in P\} \quad \cup \quad (\text{Expandir}) \\ &\quad \{(a, a, \epsilon) \mid a \in \Sigma\} \quad \quad \quad (\text{Reducir}) \end{aligned}$$

Con esto, nos encontramos con otro **problema**: hay muchas alternativas para **expandir**. ¿Cómo elegir entonces la siguiente producción para expandir? Por ejemplo, si tenemos la regla $X \rightarrow \alpha \mid \beta$, ¿cómo elegir entre α o β ?

Queremos elegir la **próxima producción** $X \rightarrow \gamma$ de tal manera que, si existe una derivación para el input, entonces $X \rightarrow \gamma$ es parte de esa derivación:

$$\text{si } S \xrightarrow[\text{lm}]{*} uX\gamma' \xrightarrow[\text{lm}]{*} uv, \text{ entonces } \gamma\gamma' \xrightarrow[\text{lm}]{*} v$$

Necesitamos **mirar las siguientes letras** en v y ver si pueden ser producidas por α o β . Para esto, ocuparemos los conceptos de **first** y **follow**.

5.3.1. Prefijos

Definición. Sea Σ un alfabeto finito. Para un $k \geq 0$, se define

$$\begin{aligned}\Sigma^{\leq k} &= \bigcup_{i=0}^k \Sigma^i \\ \Sigma_{\#}^{\leq k} &= \Sigma^{\leq k} \cup (\Sigma^{\leq k-1} \cdot \{\#\})\end{aligned}$$

Ejemplo 5.6

Para $\Sigma = \{a, b\}$:

- ♦ $\Sigma^{\leq 2} = \{\epsilon, a, b, aa, ab, ba, bb\}$
- ♦ $\Sigma_{\#}^{\leq 2} = \{\epsilon, a, b, aa, ab, ba, bb\} \cup \{\#, a\#, b\#\}$

El símbolo $\#$ representará un EOF (End Of File), marcando el **fin de una palabra**.

Definición. Para una palabra $w = a_1 a_2 \dots a_n \in \Sigma^*$ se define el k -prefijo de w como:

$$w|_k = \begin{cases} a_1 \dots a_n & \text{si } n \leq k \\ a_1 \dots a_k & \text{si } k < n \end{cases}$$

Definimos la k -concatenación \odot_k entre strings $u, v \in \Sigma$ como:

$$u \odot_k v = (u \cdot v)|_k$$

Ejemplo 5.7

Sea $\Sigma = \{a, b\}$, entonces:

- ♦ $(abaa)|_2 = ab$ $(ab)|_2 = ab$ $(a)|_2 = a$ $(\epsilon)|_2 = \epsilon$
- ♦ $a \odot_2 baa = (abaa)|_2 = ab$
- ♦ $bba \odot_2 a = (bbaa)|_2 = bb$
- ♦ $b \odot_2 \epsilon = (b)|_2 = b$

Extendemos estas operaciones para lenguajes $L, L_1, L_2 \subseteq \Sigma^*$ como:

$$\begin{aligned}L|_k &= \{w|_k \mid w \in L\} \\ L_1 \odot_k L_2 &= \{w_1 \odot_k w_2 \mid w_1 \in L_1 \text{ y } w_2 \in L_2\}\end{aligned}$$

Ejemplo 5.8

- ♦ $((ab)^*)|_3 = \{\epsilon, ab, aba\}$
- ♦ $(a)^* \odot_3 (ab)^* = \{\epsilon, a, aa, aaa, ab, aba, aab\}$

Podemos decir que los operadores $|_k$ y \odot_k “miran” hasta un prefijo k .

Propiedades. Para todo $k \geq 1$ y $L_1, L_2, L_3 \subseteq \Sigma^*$:

1. $L_1 \odot_k (L_2 \odot_k L_3) = (L_1 \odot_k L_2) \odot_k L_3$
2. $L_1 \odot_k \{\epsilon\} = \{\epsilon\} \odot_k L_1 = L_1|_k$
3. $(L_1 L_2)|_k = L_1|_k \odot_k L_2|_k$
4. $L_1 \odot_k (L_2 \cup L_3) = (L_1 \odot_k L_2) \cup (L_1 \odot_k L_3)$

La demostración de estas propiedades queda como ejercicio propuesto al lector.

5.3.2. First y Follow

Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática libre de contexto y $k \geq 1$.

Definición. Se define la función $\text{first}_k : (V \cup \Sigma)^* \rightarrow 2^{\Sigma^{\leq k}}$ tal que, para $\gamma \in (V \cup \Sigma)^*$:

$$\boxed{\text{first}_k(\gamma) = \{u|_k \mid \gamma \xrightarrow{*} u\}}$$

Ejemplo 5.9

$$E \rightarrow (E + E) \mid (E * E) \mid n$$

- ♦ $\text{first}_1(E) = \{(, n\}$
- ♦ $\text{first}_2(E) = \{n, (n, ()\}$
- ♦ $\text{first}_3(E) = \{n, (n+, (n*, ((n, (((\}$

Definición. Se define la función $\text{follow}_k : V \rightarrow 2^{\Sigma^{\leq k}}$ como:

$$\boxed{\text{follow}_k(X) = \{w \mid S \xrightarrow{*} \alpha X \beta \text{ y } w \in \text{first}_k(\beta \#)\}}$$

Ejemplo 5.10

$$E \rightarrow (E + E) \mid (E * E) \mid n$$

- ♦ $\text{follow}_1(E) = \{\#, +, *,)\}$
- ♦ $\text{follow}_2(E) = \{\#,)\#,),),)+,)*, +, *, (, +n, *n\}$

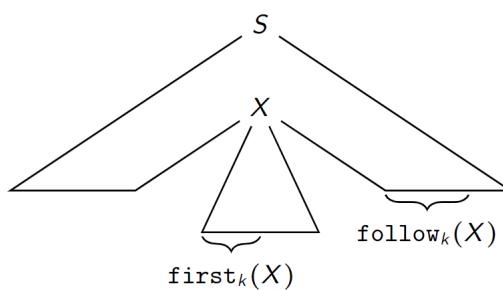


Figura 11: Representación de first y follow

5.3.3. Calcular First

Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática libre de contexto y $k \geq 1$.

Proposición. Para $X_1, \dots, X_n \in (V \cup \Sigma)$:

$$\text{first}_k(X_1 \dots X_n) = \text{first}_k(X_1) \odot_k \dots \odot_k \text{first}_k(X_n)$$

Demostración. Defina $\mathcal{L}(X) = \{w \mid X \xrightarrow{*} w\}$ y $\mathcal{L}(\gamma) = \{w \mid \gamma \xrightarrow{*} w\}$. Notar que $\text{first}_k(\gamma) = \mathcal{L}(\gamma)|_k$, por lo tanto, tenemos que

$$\begin{aligned} \text{first}_k(X_1 \dots X_n) &= \mathcal{L}(X_1 \dots X_n)|_k \\ &= (\mathcal{L}(X_1) \cdot \mathcal{L}(X_2) \cdot \dots \cdot \mathcal{L}(X_n))|_k \\ &= \mathcal{L}(X_1)_k \odot_k \mathcal{L}(X_2)_k \odot_k \dots \odot_k \mathcal{L}(X_n)|_k \\ &= \text{first}_k(X_1) \odot_k \dots \odot_k \text{first}_k(X_n) \end{aligned}$$

■

En particular, tenemos que:

$$\text{first}_k(X) = \bigcup_{X \rightarrow X_1 \dots X_n \in P} \text{first}_k(X_1) \odot_k \dots \odot_k \text{first}_k(X_n)$$

Definimos el siguiente **programa recursivo** para todo $X \in (V \cup \Sigma)$:

$$\begin{aligned} \text{first}_k^0(X) &:= \bigcup_{X \rightarrow w \in P} w|_k \\ \text{first}_k^i(X) &:= \bigcup_{X \rightarrow X_1 \dots X_n \in P} \text{first}_k^{i-1}(X_1) \odot_k \dots \odot_k \text{first}_k^{i-1}(X_n) \end{aligned}$$

Es fácil ver que:

- ♦ $\text{first}_k^{i-1}(X) \subseteq \text{first}_k^i(X)$ para todo $i > 1$.
- ♦ Como $\text{first}_k(X) \subseteq \Sigma^{\leq k}$, entonces para algún $i \leq k \cdot |\Sigma|^k \cdot |V|$ tendremos:

$$\text{first}_k^j(X) = \text{first}_k^{j+1}(X) \text{ para todo } j \geq i$$

Teorema 22

Sea i^* el menor número tal que $\text{first}_k^{i^*}(X) = \text{first}_k^{i^*+1}(X)$ para todo $X \in V$. Entonces, para todo $X \in V$:

$$\text{first}_k^{i^*}(X) = \text{first}_k(X)$$

La demostración del teorema anterior queda como ejercicio propuesto para el lector. Una idea para la dirección \subseteq , es demostrar por inducción que $\text{first}_k^i(X) \subseteq \text{first}_k(X)$. Para la dirección \supseteq , una idea es demostrar por inducción que si $X \xrightarrow{*} w$, entonces $w|_k \in \text{first}_k^i(X)$ para algún i .

Algoritmo. A continuación se presenta un algoritmo para calcular first_k :

- ♦ **Input:** Gramática $\mathcal{G} = (V, \Sigma, P, S)$ y $k \geq 1$.
- ♦ **Output:** Todos los conjuntos $\text{first}_k(X)$ para todo $X \in (V \cup \Sigma)$.

Function CalcularFirst(\mathcal{G}, k):

```

foreach  $a \in \Sigma$  do
|  $\text{first}_k^0(a) := \{a\}$ 
foreach  $X \in V$  do
|  $\text{first}_k^0(X) := \cup_{X \rightarrow w \in P} w|_k$ 
 $i := 0$ 
repeat
|  $i := i + 1$ 
| foreach  $a \in \Sigma$  do
| |  $\text{first}_k^i(a) := \{a\}$ 
| foreach  $X \in V$  do
| |  $\text{first}_k^i(X) := \bigcup_{X \rightarrow X_1 \dots X_n \in P} \text{first}_k^{i-1}(X_1) \odot_k \dots \odot_k \text{first}_k^{i-1}(X_n)$ 
| until  $\text{first}_k^i(X) = \text{first}_k^{i-1}(X)$  para todo  $X \in (V \cup \Sigma)$ 
return  $\{\text{first}_k(X)\}_{X \in (V \cup \Sigma)}$ 
```

5.3.4. Calcular Follow

Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática libre de contexto y $k \geq 1$. Si consideramos $X \neq S$:

$$\begin{aligned}
\text{follow}_k(X) &= \bigcup_{S \xrightarrow{*} \alpha X \beta} \text{first}_k(\beta \#) \\
&= \bigcup_{S \xrightarrow{*} \alpha Y \beta \Rightarrow \alpha \alpha' X \beta' \beta} \text{first}_k(\beta' \beta \#) \\
&= \bigcup_{Y \rightarrow \alpha' X \beta'} \bigcup_{S \xrightarrow{*} \alpha Y \beta} \text{first}_k(\beta' \beta \#) \\
&= \bigcup_{Y \rightarrow \alpha' X \beta'} \bigcup_{S \xrightarrow{*} \alpha Y \beta} \text{first}_k(\beta') \odot_k \text{first}_k(\beta \#) \\
&= \bigcup_{Y \rightarrow \alpha' X \beta'} \text{first}_k(\beta') \odot_k \bigcup_{S \xrightarrow{*} \alpha Y \beta} \text{first}_k(\beta \#) \\
&= \bigcup_{Y \rightarrow \alpha' X \beta'} \text{first}_k(\beta') \odot_k \text{follow}_k(Y)
\end{aligned}$$

Si consideramos $X = S$:

$$\begin{aligned}
\text{follow}_k(S) &= \{\#\} \cup \bigcup_{S \xrightarrow{+} \alpha S \beta} \text{first}_k(\beta \#) \\
&= \{\#\} \cup \bigcup_{S \xrightarrow{*} \alpha Y \beta \Rightarrow \alpha \alpha' S \beta' \beta} \text{first}_k(\beta' \beta \#) \\
&= \{\#\} \cup \bigcup_{Y \rightarrow \alpha' S \beta'} \text{first}_k(\beta') \odot_k \text{follow}_k(Y)
\end{aligned}$$

Dado lo anterior, podemos definir el siguiente teorema.

Teorema 23

Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática libre de contexto y $k \geq 1$. Entonces:

$$\text{Para } X \neq S : \quad \text{follow}_k(X) = \bigcup_{Y \rightarrow \alpha X \beta} \text{first}_k(\beta) \odot_k \text{follow}_k(Y)$$

$$\text{Para } X = S : \quad \text{follow}_k(S) = \{\#\} \cup \bigcup_{Y \rightarrow \alpha S \beta} \text{first}_k(\beta) \odot_k \text{follow}_k(Y)$$

Definimos el siguiente **programa recursivo** para todo $X \in V$:

$$\text{Para } X \neq S : \quad \text{follow}_k^0(X) := \emptyset$$

$$\text{Para } X = S : \quad \text{follow}_k^0(S) := \{\#\}$$

$$\text{Para } X \neq S : \quad \text{follow}_k^i(X) := \bigcup_{Y \rightarrow \alpha X \beta} \text{first}_k(\beta) \odot_k \text{follow}_k^{i-1}(Y)$$

$$\text{Para } X = S : \quad \text{follow}_k^i(S) := \{\#\} \cup \bigcup_{Y \rightarrow \alpha S \beta} \text{first}_k(\beta) \odot_k \text{follow}_k^{i-1}(Y)$$

Similar al caso de first_k , es fácil ver que:

- ♦ $\text{follow}_k^{i-1}(X) \subseteq \text{follow}_k^i(X)$ para todo $i > 1$.
- ♦ Como $\text{follow}_k(X) \subseteq \Sigma^{\leq k}$, entonces para algún $i \leq k \cdot |\Sigma|^k \cdot |V|$:

$$\text{follow}_k^j(X) = \text{follow}_k^{j+1}(X) \text{ para todo } j \geq i$$

Teorema 24

Sea i^* el menor número tal que $\text{follow}_k^{i^*}(X) = \text{follow}_k^{i^*+1}(X)$ para todo $X \in V$. Entonces, para todo $X \in V$:

$$\text{follow}_k^{i^*}(X) = \text{follow}_k(X)$$

La demostración de este teorema se deja como ejercicio propuesto al lector.

Con todo lo anterior, podemos calcular $\text{follow}_k(X)$ con un algoritmo similar que $\text{first}_k(X)$. Respecto a la eficiencia de este tipo de algoritmos:

- ♦ Toman $\mathcal{O}(k \cdot |\Sigma|^k \cdot |V|)$ en el peor caso.
- ♦ Si $k = 1$, el número de repeticiones será $\mathcal{O}(|\Sigma| \cdot |V|)$ y el tiempo del algoritmo será polinomial en $|\mathcal{G}|$ en el peor caso. Incluso, se puede hacer en tiempo $\mathcal{O}(|V| \cdot |P|)$ en total.

5.4. Gramáticas LL

Volvamos a la idea de buscar un algoritmo que haga parsing en **tiempo lineal**. Para esto, construimos un autómata apilador alternativo \mathcal{D} al cual le expandimos sus producciones. ¿El problema? No sabemos cómo elegir qué producciones expandir. Debido a lo anterior, introducimos los conceptos de **first** y **follow**. Así que, si tenemos una producción de la forma $X \rightarrow \alpha \mid \beta$, ¿cómo elegir entre α o β ?

Estrategia (intuición). La idea es la siguiente:

1. Mirar k símbolos del resto del input v (**k -lookahead**).
2. Usar $v|_k$ y decidir cuál regla $X \rightarrow \gamma$ elegimos para expandir.

La caracterización de las gramáticas que cumplen las propiedades anteriores se denominan **Gramáticas LL(k)**, donde

- ♦ **Primera L:** leer el input de izquierda a derecha (**Left-right**).
 - ♦ **Segunda L:** producir una derivación por la izquierda (**Leftmost**).
 - ♦ **Parámetro k :** el número de letras en adelante que utiliza (**lookahead**).

5.4.1. Definición Gramáticas LL

Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática libre de contexto y $k \geq 1$.

Definición. Decimos que \mathcal{G} es una gramática LL(k) si para todas las derivaciones:

- ♦ $S \xrightarrow[\text{lm}]{*} uY\beta \xrightarrow[\text{lm}]{*} u\gamma_1\beta \xrightarrow[\text{lm}]{*} uv_1$
 - ♦ $S \xrightarrow[\text{lm}]{*} uY\beta \xrightarrow[\text{lm}]{*} u\gamma_2\beta \xrightarrow[\text{lm}]{*} uv_2 \text{ y}$
 - ♦ $v_1|_k = v_2|_k$

entonces se cumple que $\gamma_1 = \gamma_2$.

Notar que la elección de $Y \rightarrow \gamma$ depende de $Y, v|_k$ y u .

Ejemplo 5.11: Gramáticas LL(1)

$$\begin{array}{ccccccc}
 & \mathcal{G}_1 : & S & \rightarrow & (S) & | & n \\
 \\
 S & \xrightarrow[\text{Im}]{\star} & \overbrace{\cdots(\overbrace{S})\cdots}^{\substack{u \\ Y \\ \beta}} & \Rightarrow & (\cdots(\gamma_1)\cdots) & \xrightarrow[\text{Im}]{\star} & \overbrace{\cdots(\overbrace{v'_1})\cdots}^{\substack{u \\ v_1}} \\
 & & \xrightarrow[\text{Im}]{\star} & & (\cdots(\gamma_2)\cdots) & \xrightarrow[\text{Im}]{\star} & \overbrace{\cdots(\overbrace{v'_2})\cdots}^{\substack{u \\ v_2}}
 \end{array}$$

- ♦ Si $v_1|_1 = v_2|_1 = n$, entonces $\gamma_1 = \gamma_2 = n$.
 - ♦ Si $v_1|_1 = v_2|_1 = ($, entonces $\gamma_1 = \gamma_2 = (S)$

En ambos casos, tenemos que $\gamma_1 = \gamma_2$ y \mathcal{G}_1 es una gramática LL(1).

$$\begin{array}{rcccl}
 \mathcal{G}_2 : & S & \rightarrow & (X) & | \quad n \\
 & X & \rightarrow & SX & | \quad \epsilon
 \end{array}$$

$$\begin{array}{ccccccc}
 S & \xrightarrow[\text{Im}]{\star} & uX\beta & \xrightarrow[\text{Im}]{\star} & u\gamma_1\beta & \xrightarrow[\text{Im}]{\star} & uv_1 \\
 & & & & & & \\
 & & & & \xrightarrow[\text{Im}]{\star} & u\gamma_2\beta & \xrightarrow[\text{Im}]{\star} uv_2
 \end{array}$$

- ♦ Si $v_1|_1 = v_2|_1 = ($ o ‘n’, entonces $\gamma_1 = \gamma_2 = SX$.
 - ♦ Si $v_1|_1 = v_2|_1 =),$ entonces $\gamma_1 = \gamma_2 = \epsilon$.

Por lo tanto, tenemos que $\gamma_1 = \gamma_2$ y \mathcal{G}_2 es también una gramática LL(1).

Ejemplo 5.12: Gramática NO LL(1) pero si LL(2)

$$\begin{array}{rcl}
 \mathcal{G}_3 : & S & \rightarrow (X) \quad | \quad n + S \quad | \quad n \\
 & X & \rightarrow SX \quad | \quad \epsilon
 \end{array}$$

$$S \xrightarrow[\text{Im}]{\star} \underbrace{(S X)}_{\substack{u \\ \gamma_1 \\ \beta}} \Rightarrow \underbrace{(n + S X)}_{\substack{\gamma_1 \\ u \\ n + S}} \xrightarrow[\text{Im}]{\star} \underbrace{(n + n)}_{\substack{u \\ v_1 \\ n + n}} \quad S \xrightarrow[\text{Im}]{\star} u S \beta \Rightarrow u \gamma_1 \beta \xrightarrow[\text{Im}]{\star} u v_1 \beta$$

$$\xrightarrow[\text{Im}]{\star} \underbrace{(n X)}_{\substack{n \\ \gamma_2}} \Rightarrow \underbrace{(n)}_{\substack{u \\ v_2}}$$

Como $v_1|_1 = v_2|_1 = n$ pero $\gamma_1 \neq \gamma_2$, entonces \mathcal{G}_3 NO es una gramática LL(1).

- ♦ Si $v_1|_2 = v_2|_2 = n+$, entonces $\gamma_1 = \gamma_2 = n + S$.
 - ♦ Si $v_1|_1 = v_2|_1 = na$, con $a \neq +$, entonces $\gamma_1 = \gamma_2 = n$.

Por lo tanto, tenemos que $\gamma_1 = \gamma_2$ y entonces \mathcal{G}_3 es LL(2).

Ejemplo 5.13: Gramática NO LL(k)

$$\begin{array}{l} \mathcal{G}_4 : \quad S \quad \rightarrow \quad (X) \quad | \quad (X)^e \quad | \quad n + S \quad | \quad n \\ \qquad X \quad \rightarrow \quad SX \quad | \quad \epsilon \\ \\ S \quad \stackrel{\star}{\Rightarrow}_{\text{Im}} \quad \overbrace{(\quad S \quad X)}^{\substack{u \\ Y \\ \beta}} \quad \Rightarrow \quad \overbrace{((S)^e \quad X)}^{\gamma_1} \quad \stackrel{\star}{\Rightarrow} \quad \overbrace{((\quad \underbrace{\cdots (n) \cdots }^u \quad)^e)}^{v_1} \\ \qquad \qquad \qquad \stackrel{\star}{\Rightarrow}_{\text{Im}} \quad \overbrace{((S) \quad X)}^{\gamma_2} \quad \stackrel{\star}{\Rightarrow} \quad \overbrace{((\quad \underbrace{\cdots (n) \cdots }_u \quad)^e)}^{v_2} \end{array}$$

Como $v_1|_k = v_2|_k = (\cdots ($ pero $\gamma_1 \neq \gamma_2$, entonces \mathcal{G}_4 NO es una gramática LL(k) para todo k .

Ejemplo 5.14: Gramática NO LL(k) transformada en LL(2)

La gramática G_4 del ejemplo anterior se puede transformar para que sea LL(2) de la siguiente manera:

$$\begin{array}{lll} \mathcal{G}'_4 : & S & \rightarrow (XY \quad | \quad n+S \quad | \quad n \\ & X & \rightarrow SX \quad | \quad \epsilon \\ & Y & \rightarrow) \quad | \quad)^{\wedge} e \end{array}$$

Queda como ejercicio para el lector demostrar que \mathcal{G}'_4 es LL(2).

Ejemplo 5.15: Lenguaje NO LL(k)

$$\begin{array}{ll} \mathcal{G}_5 : & S \rightarrow X \mid Y \\ & X \rightarrow aXb \mid 0 \\ & Y \rightarrow aYbb \mid 1 \end{array} \quad \mathcal{L}(\mathcal{G}_5) = \{a^n0b^n \mid n \geq 0\} \cup \{a^n1b^{2n} \mid n \geq 0\}$$

$$\begin{array}{lcl} S & \xrightarrow[\text{lm}]{*} & S \xrightarrow[\text{lm}]{*} X \xrightarrow[\text{lm}]{*} a^k0b^k \\ & & \Rightarrow Y \xrightarrow[\text{lm}]{*} a^k1b^{2k} \end{array}$$

Para todo $k \geq 1$, se tiene que \mathcal{G}_5 NO es una gramática LL(k).

Es posible demostrar que, para toda gramática \mathcal{G} con $\mathcal{L}(\mathcal{G}_5) = \mathcal{L}(\mathcal{G})$, \mathcal{G} NO es una gramática LL(k) para todo $k \geq 1$.

5.4.2. Caracterización LL

Para esta parte es importante manejar las definiciones de prefijos vistas en la sección 5.3.1.

Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática libre de contexto **reducida** y $k \geq 1$. En base a esto definimos el siguiente teorema:

Teorema 25

\mathcal{G} es una gramática LL(k) si, y sólo si, para todas dos reglas distintas $Y \rightarrow \gamma_1, Y \rightarrow \gamma_2 \in P$ y para todo $S \xrightarrow[\text{lm}]{*} uY\beta$, se tiene que:

$$\text{first}_k(\gamma_1\beta) \cap \text{first}_k(\gamma_2\beta) = \emptyset$$

Demostración. (\Rightarrow) Por contrapositivo, supongamos que $v \in \text{first}_k(\gamma_1\beta) \cap \text{first}_k(\gamma_2\beta)$. Como \mathcal{G} es reducida (sin variables inútiles), entonces

$$\begin{array}{lcl} S & \xrightarrow[\text{lm}]{*} & uY\beta \xrightarrow[\text{lm}]{*} u\gamma_1\beta \xrightarrow[\text{lm}]{*} uvv_1 \\ & & \Rightarrow u\gamma_2\beta \xrightarrow[\text{lm}]{*} uvv_2 \end{array}$$

para algún $v_1, v_2 \in \Sigma^*$. Como $\gamma_1 \neq \gamma_2$, entonces \mathcal{G} NO es LL(k).

(\Leftarrow) Por contrapositivo (de nuevo), supongamos que \mathcal{G} no es LL(k). Como \mathcal{G} no es LL(k), entonces tenemos derivaciones de la forma:

$$\begin{array}{lcl} S & \xrightarrow[\text{lm}]{*} & uY\beta \xrightarrow[\text{lm}]{*} u\gamma_1\beta \xrightarrow[\text{lm}]{*} uv_1 \\ & & \Rightarrow u\gamma_2\beta \xrightarrow[\text{lm}]{*} uv_2 \end{array}$$

Vemos que $v_1|_k = v_2|_k = v$, pero $\gamma_1 \neq \gamma_2$. Por lo tanto, $v \in \text{first}_k(\gamma_1\beta) \cap \text{first}_k(\gamma_2\beta)$. ■

¿Cómo usamos la caracterización del teorema para demostrar que una gramática es LL(k)? Buscaremos condiciones más simples para verificar si una gramática es LL(k).

Definición. \mathcal{G} es una gramática LL(k) **fuerte** si para todas dos reglas distintas $Y \rightarrow \gamma_1, Y \rightarrow \gamma_2 \in P$ se tiene que:

$$\text{first}_k(\gamma_1) \odot_k \text{follow}_k(Y) \cap \text{first}_k(\gamma_2) \odot_k \text{follow}_k(Y) = \emptyset$$

Ejemplo 5.16: Si \mathcal{G} es LL(k) fuerte, entonces \mathcal{G} es LL(k)

Una gramática \mathcal{G} que sea LL(k) fuerte siempre es LL(k), ya que si definimos dos conjuntos dados por el teorema de LL(k) (F_1) y la definición de LL(k) fuerte (F_2), dados por:

$$\begin{aligned} F_1 &= \text{first}_k(\gamma_1\beta) \cap \text{first}_k(\gamma_2\beta) = \text{first}_k(\gamma_1) \odot_k \text{first}_k(\beta) \cap \text{first}_k(\gamma_2) \odot_k \text{first}_k(\beta) \\ F_2 &= \text{first}_k(\gamma_1) \odot_k \text{first}_k(Y) \cap \text{first}_k(\gamma_2) \odot_k \text{first}_k(Y) = \emptyset \end{aligned}$$

Entonces, tenemos que $F_1 \subseteq F_2$.

Ejemplo 5.17: Si \mathcal{G} es LL(k), ¿es LL(k) fuerte?

La respuesta directa es que no. Con un contrajemplo, tomemos la gramática \mathcal{G} definida por

$$\begin{aligned} \mathcal{G} : \quad S &\rightarrow aXaa \mid bXba \\ X &\rightarrow b \mid \epsilon \end{aligned}$$

Recordatorio: \mathcal{G} es LL(k) si para todas dos reglas distintas $Y \rightarrow \gamma_1, Y \rightarrow \gamma_2 \in P$ y para todo $S \xrightarrow[\text{lm}]{} uY\beta$, se tiene que

$$\text{first}_k(\gamma_1\beta) \cap \text{first}_k(\gamma_2\beta) = \emptyset$$

- ♦ Si $S \xrightarrow[\text{lm}]{} aXaa$, entonces $\text{first}_2(baa) \cap \text{first}_2(aa) = \emptyset$.
- ♦ Si $S \xrightarrow[\text{lm}]{} bXba$, entonces $\text{first}_2(baa) \cap \text{first}_2(ba) = \emptyset$

Por lo tanto, \mathcal{G} es LL(2).

Recordatorio: \mathcal{G} es una gramática LL(k) fuerte si para todas dos reglas distintas $Y \rightarrow \gamma_1, Y \rightarrow \gamma_2 \in P$ se tiene que:

$$\text{first}_k(\gamma_1) \odot_k \text{follow}_k(Y) \cap \text{first}_k(\gamma_2) \odot_k \text{follow}_k(Y) = \emptyset$$

Si vemos $X \rightarrow b$ y $X \rightarrow \epsilon$:

$$\begin{aligned} \text{first}_2(b) \odot_2 \text{follow}_2(X) &\cap \text{first}_2(\epsilon) \odot_2 \text{follow}_2(X) \\ &= \{b\} \odot_2 \{aa, ba\} \cap \{\epsilon\} \odot_2 \{aa, ba\} \\ &= \{ba, bb\} \cap \{aa, ba\} \\ &= \{ba\} \quad \text{y por ende } \mathcal{G} \text{ no es LL(2) fuerte.} \end{aligned}$$

¿Qué pasa con el caso LL(1)? Supongamos que \mathcal{G} es LL(1) y $Y \rightarrow \gamma_1, Y \rightarrow \gamma_2 \in P$ son reglas distintas.

1. Si $\epsilon \notin \text{first}_1(\gamma_1)$ y $\epsilon \notin \text{first}_1(\gamma_2)$, entonces, por la caracterización de LL(1):

$$\begin{aligned} \emptyset &= \text{first}_1(\gamma_1\beta) \cap \text{first}_1(\gamma_2\beta) \\ &= \text{first}_1(\gamma_1) \cap \text{first}_1(\gamma_2) \\ &= \text{first}_1(\gamma_1) \odot_1 \text{follow}_1(Y) \cap \text{first}_1(\gamma_2) \odot_1 \text{follow}_1(Y) \end{aligned}$$

2. Si $\epsilon \in \text{first}_1(\gamma_1)$ y $\epsilon \notin \text{first}_1(\gamma_2)$, entonces, por la caracterización de LL(1):

$$\begin{aligned} \emptyset &= \text{first}_1(\gamma_1\beta) \cap \text{first}_1(\gamma_2\beta) \\ &= \text{first}_1(\gamma_1\beta) \cap \text{first}_1(\gamma_2) \\ &= \text{first}_1(\gamma_1\beta) \cap \text{first}_1(\gamma_2\beta') \end{aligned}$$

para todo $\beta' \in (V \cup \Sigma)^*$. Por lo tanto:

$$\begin{aligned} & \text{first}_1(\gamma_1) \odot_1 \text{follow}_1(Y) \cap \text{first}_1(\gamma_2) \odot_1 \text{follow}_1(Y) \\ &= \bigcup_{\substack{S \xrightarrow[\text{Im}]{*} uY\beta}} \text{first}_1(\gamma_1\beta) \cap \bigcup_{\substack{S \xrightarrow[\text{Im}]{*} uY\beta'}} \text{first}_1(\gamma_2\beta') = \emptyset \end{aligned}$$

Por lo tanto, establecemos el siguiente teorema.

Teorema 26

Una gramática \mathcal{G} es LL(1) si, y sólo si, \mathcal{G} es LL(1) **fuerte**, esto es, para todas dos reglas distintas $Y \rightarrow \gamma_1, Y \rightarrow \gamma_2 \in P$:

$$\text{first}_1(\gamma_1) \odot_1 \text{follow}_1(Y) \cap \text{first}_1(\gamma_2) \odot_1 \text{follow}_1(Y) = \emptyset$$

La condición del teorema anterior se puede verificar en **tiempo polinomial** en \mathcal{G} .

5.5. Parsing con gramáticas LL(k)

5.5.1. Algunas consideraciones

Considere la siguiente gramática \mathcal{G} :

$$\begin{aligned} S &\rightarrow Xa \mid Xb \\ X &\rightarrow c \end{aligned}$$

¿Es esta gramática del tipo LL(1)? Podemos ver que $\text{first}_1(\gamma_1\beta) = \{c\}$ y $\text{first}_1(\gamma_2\beta) = \{c\}$, con $\gamma_1 = Xa$, $\gamma_2 = Xb$ y $\beta = \epsilon$. Por lo tanto su intersección no es vacía y entonces \mathcal{G} no es LL(1). ¿Podemos establecer una solución para este problema?

Factorización. En general, si tenemos una regla:

$$X \rightarrow \gamma\alpha_1 \mid \gamma\alpha_2$$

siempre podemos “**factorizar**” la regla manteniendo la semántica, como:

$$\begin{aligned} X &\rightarrow \gamma X' \\ X' &\rightarrow \alpha_1 \mid \alpha_2 \end{aligned}$$

Considere ahora la siguiente gramática \mathcal{G} :

$$E \rightarrow E * E \mid n$$

¿Es esta gramática del tipo LL(1)? ¿LL(k)? Pues no es ninguna. El problema con esta gramática es su recursividad, en específico, por la izquierda.

Definición. Una gramática \mathcal{G} se dice **recursiva por la izquierda** si existe $X \in V$ tal que:

$$X \stackrel{\pm}{\Rightarrow} X\gamma \quad \text{para algún } \gamma \in (V \cup \Sigma)^*$$

Teorema 27

Si $\mathcal{G} = (V, \Sigma, P, S)$ es una gramática reducida y recursiva por la izquierda, entonces \mathcal{G} NO es LL(k) para todo $k \geq 1$.

Demostración. Por simplicidad, suponga que $X \rightarrow X\beta \in P$ y $X \rightarrow w \in P$.

Como \mathcal{G} es reducida, entonces existe una derivación $S \xrightarrow[\text{lm}]{}^* uX\gamma$:

$$S \xrightarrow[\text{lm}]{}^* uX\gamma \xrightarrow[\text{lm}]{}^{\text{n-veces}} \xrightarrow[\text{lm}]{}^* uX\beta^n\gamma$$

Por **contradicción**, suponga que \mathcal{G} es LL(k). Por lo tanto:

$$\text{first}_k(X\beta^{n+1}\gamma) \cap \text{first}_k(w\beta^n\gamma) = \emptyset$$

Suponga que $\beta \xrightarrow{*} v \in \Sigma^*$ y $\gamma \xrightarrow{*} v' \in \Sigma^*$. Con $n = k$, tendremos que

$$(wv^kv')|_k \in \text{first}_k(X\beta^{k+1}\gamma) \cap \text{first}_k(w\beta^k\gamma) \rightarrow \leftarrow (\text{¡contradicción! el conjunto no es vacío})$$

■

Hablemos de recursión **inmediata** por la izquierda. Suponga que existe $X \in V$ tal que:

$$X \rightarrow X\alpha_1 | \dots | X\alpha_m | \beta_1 | \dots | \beta_n$$

¿Cómo podemos **eliminar** la recursión inmediata por la izquierda? Consideramos la misma gramática pero **cambiando** las reglas de X por:

$$\begin{aligned} X &\rightarrow \beta_1X' | \dots | \beta_nX' \\ X' &\rightarrow \alpha_1X' | \dots | \alpha_mX' | \epsilon \end{aligned}$$

Ejemplo 5.18: Eliminando recursión inmediata

$\begin{aligned} S &\rightarrow Xa b \\ X &\rightarrow Xc d \end{aligned}$	$\begin{aligned} S &\rightarrow Xa b \\ X &\rightarrow dX' \\ X' &\rightarrow cX' \epsilon \end{aligned}$
--	---

Teorema 28

Sea \mathcal{G} una gramática tal que existe $X \in V$:

$$X \rightarrow X\alpha_1 | \dots | X\alpha_m | \beta_1 | \dots | \beta_n$$

Sea \mathcal{G}' la misma gramática \mathcal{G} pero cambiando las reglas de X por:

$$\begin{aligned} X &\rightarrow \beta_1X' | \dots | \beta_nX' \\ X' &\rightarrow \alpha_1X' | \dots | \alpha_mX' | \epsilon \end{aligned}$$

Entonces $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}')$

Demostración. Una derivación por la izquierda de X en \mathcal{G} :

$$X \xrightarrow[\text{lm}]{} X\alpha_{i_1} \xrightarrow[\text{lm}]{} X\alpha_{i_2}\alpha_{i_1} \xrightarrow[\text{lm}]{} \dots \xrightarrow[\text{lm}]{} X\alpha_{i_p}\alpha_{i_{p-1}}\dots\alpha_{i_1} \xrightarrow[\text{lm}]{} \beta_j\alpha_{i_p}\alpha_{i_{p-1}}\dots\alpha_{i_1}$$

Una derivación por la derecha de X en \mathcal{G}' **equivalente**:

$$X \xrightarrow[\text{rm}]{} \beta_jX' \xrightarrow[\text{rm}]{} \beta_j\alpha_{i_p}X' \xrightarrow[\text{rm}]{} \dots \xrightarrow[\text{rm}]{} \beta_j\alpha_{i_p}\dots\alpha_{i_2}\alpha_{i_1}X' \xrightarrow[\text{rm}]{} \beta_j\alpha_{i_p}\alpha_{i_{p-1}}\dots\alpha_{i_1}$$

■

Ahora, ¿qué pasa si la recursión por la izquierda es **no-inmediata**? Considere la siguiente gramática **recursiva por la izquierda**:

$$\begin{aligned} S &\rightarrow Xa \mid b \\ X &\rightarrow Yc \\ Y &\rightarrow Xd \mid e \end{aligned}$$

¿Cómo eliminamos la recursión por la izquierda no-inmediata?

Estrategia. Dado $V = \{X_1, \dots, X_n\}$, removemos la recursión inductivamente en n , tal que, en cada paso i de la inducción, se cumplira que para todo $i, j \leq n$:

$$\text{si } X_i \rightarrow X_j \alpha, \text{ entonces } i < j$$

input : Gramática $\mathcal{G} = (V, \Sigma, P, S)$ y $V = \{X_1, \dots, X_n\}$

output: Gramática \mathcal{G} sin recursión por la izquierda

Function EliminarRecursion(\mathcal{G}):

```

 $P' := P$ 
for  $i = 1$  to  $n$  do
  for  $j = 1$  to  $i - 1$  do
    foreach  $X_i \rightarrow X_j \gamma \in P'$  do
      foreach  $X_j \rightarrow \alpha \in P'$  do
         $P' := P' \cup \{X_i \rightarrow \alpha \gamma\}$ 
     $P' := P' - \{X_i \rightarrow X_j \gamma\}$ 
    Remover recursión inmediata para  $X_i$  en  $P'$  (si existe)
   $V' := \{X_1, \dots, X_n\} \cup \{X'_1, \dots, X'_n\}$ 
return  $(V', \Sigma, P', S)$ 

```

Queda como ejercicio propuesto al lector demostrar la correctitud del algoritmo.

Ejemplo 5.19: Eliminando recursión

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid n \end{aligned}$$

Eliminando la **recursión inmediata** de E :

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid n \end{aligned}$$

Eliminando la **recusión inmediata** de T :

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \epsilon \\ F &\rightarrow (E) \mid n \end{aligned}$$

Conclusión. Es **posible eliminar** la recursividad por la izquierda, pero esto **NO asegura** que el resultado sea una gramática LL(k) para algún k .

5.5.2. Parsing de LL(k)

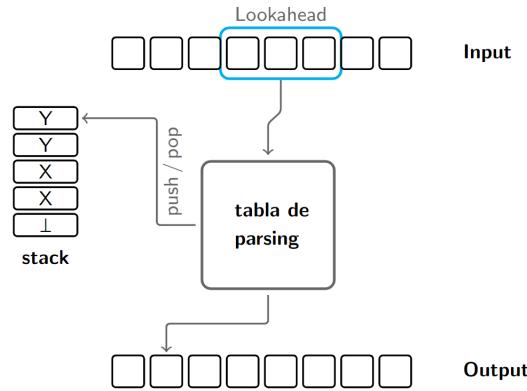


Figura 12: Máquina de parsing

Definición. Sea Σ un alfabeto finito. Se definen los siguientes conjuntos de palabras:

- ♦ $\dot{\Sigma} = \Sigma^* \times \Sigma^*$
- ♦ $\dot{\Sigma}^{\leq k} = \{(u, v) \in \dot{\Sigma} \mid |uv| \leq k\}$
- ♦ $\dot{\Sigma}_{\#}^{\leq k} = \{(u, v) \in \dot{\Sigma} \mid |uv| \leq k\} \cup \{(u, v\#) \mid (u, v) \in \dot{\Sigma} \mid |uv| < k\}$

Notación. En vez de usar $(u, v) \in \dot{\Sigma}_{\#}^{\leq k}$, escribiremos $u.v \in \dot{\Sigma}_{\#}^{\leq k}$. El par $\epsilon.\epsilon$ lo denotaremos solamente por ϵ .

Definición. Un transductor apilador con k -lookahead (k -PDT) es una tupla:

$$\boxed{\mathcal{T} = (Q, \Sigma, \Omega, \Delta, q_0, F)}$$

- ♦ Q es un conjunto finito de estados.
- ♦ Σ es el alfabeto de input.
- ♦ Ω es el alfabeto de output.
- ♦ $\Delta \subseteq Q^+ \times \dot{\Sigma}_{\#}^{\leq k} \times (\Omega \cup \{\epsilon\}) \times Q^*$ es la relación de transición.
- ♦ $q_0 \in Q$ es un conjunto de estados iniciales.
- ♦ $F \subseteq Q$ es el conjunto de estados finales.

Definición. Una **configuración** de \mathcal{T} es una tupla:

$$(q_1 \dots q_k, w, o) \in (Q^+, \Sigma^* \cdot \{\#\}, \Omega^*)$$

- ♦ $q_1 \dots q_k$ es el contenido del stack con q_1 el tope del stack.
- ♦ w es el contenido del input.
- ♦ o es el contenido del output.

Decimos que una configuración:

- ♦ $(q_0, w\#, \epsilon)$ es **inicial** y
- ♦ $(q_f, \#, o)$ es **final** si $q_f \in F$.

Definición. Se define la relación $\vdash_{\mathcal{T}}$ de **siguiente-paso** entre configuraciones de \mathcal{T} :

$$(\gamma_1, w_1, o_1) \vdash_{\mathcal{T}} (\gamma_2, w_2, o_2)$$

si, y sólo si, existe $(\alpha, u.v, a, \beta) \in \Delta$, $\gamma \in \Gamma^*$ y $w \in \Sigma^* \cdot \{\#\}$ tal que:

- ♦ **Stack:** $\gamma_1 = \alpha \cdot \gamma$ y $\gamma_2 = \beta \cdot \gamma$
- ♦ **Look-ahead:** $w_1 = u \cdot v \cdot w$ y $w_2 = v \cdot w$
- ♦ **Output:** $o_2 = o_1 \cdot a$

Se define $\vdash_{\mathcal{T}}^*$ como la clausura **refleja** y **transitiva** de $\vdash_{\mathcal{T}}$.

Definición. \mathcal{T} entrega o con **input** w si existe una configuración inicial $(q_0, w \cdot \#, \epsilon)$ y una configuración final $(q_f, \#, o)$ tal que:

$$(q_0, w \cdot \#, \epsilon) \vdash_{\mathcal{T}}^* (q_f, \#, o)$$

Se define la función $\llbracket \mathcal{T} \rrbracket : \Sigma^* \rightarrow 2^{\Omega^*}$:

$$\llbracket \mathcal{T} \rrbracket(w) = \{o \in \Omega^* \mid \mathcal{T} \text{ entrega } o \text{ con input } w\}$$

Definición. \mathcal{T} es **determinista** si para todo $(\alpha_1, u_1.v_1, a_1, \beta_1), (\alpha_2, u_2.v_2, a_2, \beta_2) \in \Delta$ con $(\alpha_1, u_1.v_1, a_1, \beta_1) \neq (\alpha_2, u_2.v_2, a_2, \beta_2)$ se cumple que

$$\alpha_1 \text{ NO es prefijo de } \alpha_2 \quad \text{o} \quad u_1v_1 \text{ NO es prefijo de } u_2v_2.$$

“Para cualquier configuración (γ, w, o) existe a lo más una configuración (γ', w', o') tal que $(\gamma, w, o) \vdash_{\mathcal{T}}^ (\gamma', w', o')$ ”*

La **ventaja** de un k -PDT determinista es que nos aseguramos de que siempre obtenemos un solo output para cada input (el no-determinismo nos podría generar muchos outputs distintos).

Construcción del parser. Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática LL(k) fuerte. Se define el k -PDT para \mathcal{G} :

$$\mathcal{T}[\mathcal{G}] = \left(V \cup \Sigma \cup \{q_0, q_f\}, \Sigma, \underbrace{P}_{\Omega}, \Delta, q_0, \{q_f\} \right)$$

La relación de transición Δ se define como:

- Iniciar:** $(q_0, \epsilon, \epsilon, S \cdot q_f)$
- Reducir:** $(a, a, \epsilon, \epsilon)$ para cada $a \in \Sigma$
- Expandir:** $(X, .u, p, \gamma)$
para cada $p := (X \rightarrow \gamma) \in P$ tal que $u \in \text{first}_k(\gamma) \odot_k \text{follow}_k(X)$

Propiedades. $\mathcal{T}[\mathcal{G}]$ tiene las siguientes propiedades:

1. $\mathcal{T}[\mathcal{G}]$ es un k -PDT **determinista** si, y sólo si, \mathcal{G} es LL(k) fuerte.
2. Si $w \notin \mathcal{L}(\mathcal{G})$ entonces $\llbracket \mathcal{T} \rrbracket(w) = \emptyset$.
3. Si $w \in \mathcal{L}(\mathcal{G})$ entonces $\llbracket \mathcal{T} \rrbracket(w) = \{r_1 \dots r_m\}$ es una derivación por la izquierda de \mathcal{G} sobre w .

Algoritmo. Para una gramática LL(k) \mathcal{G} y una palabra $w \in \Sigma^*$:

1. Construya el k -PDT determinista $\mathcal{T}[\mathcal{G}]$ a partir de \mathcal{G} .
2. Ejecute $\mathcal{T}[[G]]$ sobre w .

Como $\mathcal{T}[\mathcal{G}]$ es determinista, entonces el algoritmo toma **tiempo lineal** en w .

Tabla predictiva para LL(k) fuerte. Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática LL(k) fuerte. Para cada $u \in \Sigma^k \cup \Sigma^{<k} \cdot \{\#\}$, se define $M[X, u] \in (V \cup \Sigma)^* \cup \{\text{ERROR}\}$:

$$M[X, u] = \begin{cases} \gamma & \text{si } X \rightarrow \gamma \in P \text{ y } u \in \text{first}_k(\gamma) \odot_k \text{follow}_k(X) \\ \text{ERROR} & \text{en otro caso.} \end{cases}$$

El cálculo de la tabla predictiva puede tomar **tiempo exponencial** en $|\mathcal{G}|$ y k .

Caso especial: tabla predictiva para LL(1). Sea $\mathcal{G} = (V, \Sigma, P, S)$ una gramática LL(1) fuerte. Para cada $a \in \Sigma \cup \{\#\}$, se define $M[X, a] \in (V \cup \Sigma)^* \cup \{\text{ERROR}\}$:

$$M[X, a] = \begin{cases} \gamma & \text{si } X \rightarrow \gamma \in P \text{ y } a \in \text{first}_1(\gamma) \\ \gamma & \text{si } X \rightarrow \gamma \in P, \epsilon \in \text{first}_1(\gamma) \text{ y } a \in \text{follow}_1(X) \\ \text{ERROR} & \text{en otro caso.} \end{cases}$$

Este cálculo se puede hacer en tiempo $\mathcal{O}(|V| \cdot |P|)$.

Ejemplo 5.20: Tabla predictiva

$E \rightarrow TE'$	$\text{first}_1(E) = \{(\text{id})\}$	$\text{first}_1(E') = \{+, \epsilon\}$
$E' \rightarrow +TE' \mid \epsilon$	$\text{first}_1(T) = \{(\text{id})\}$	$\text{first}_1(T') = \{*, \epsilon\}$
$T \rightarrow FT'$	$\text{first}_1(F) = \{(\text{id})\}$	
$T' \rightarrow *FT' \mid \epsilon$	$\text{follow}_1(E) = \{(), \#\}$	$\text{follow}_1(E') = \{(), \#\}$
$F \rightarrow (E) \mid \text{id}$	$\text{follow}_1(T) = \{+, (), \#\}$	$\text{follow}_1(T') = \{+, (), \#\}$
	$\text{follow}_1(F) = \{+, *, (), \#\}$	

	id	+	*	()	#
E	TE'	ERROR	ERROR	TE'	ERROR	ERROR
E'	ERROR	$+TE'$	ERROR	ERROR	ϵ	ϵ
T	FT'	ERROR	ERROR	FT'	ERROR	ERROR
T'	ERROR	ϵ	$*FT'$	ERROR	ϵ	ϵ
F	id	ERROR	ERROR	(E)	ERROR	ERROR

6. Extracción de información

6.1. Extracción

6.2. Enumeración de resultados: Autómatas con anotaciones