

## Proyecto I

Expresiones regulares E1, E2 y E3 que corresponden, respectivamente, al reconocimiento de la palabra clave “as”, la palabra clave “array” y de un identificador de variables.

E1: as

E2: array

E3:  $(a + b + \dots + z + A + B + \dots + Z)(a + b + \dots + z + A + B + \dots + Z + 0 + 1 + \dots + 9 + \_)*$

Vemos que la expresión regular más compleja es la del identificador. Vamos a identificar algunos conjuntos de interés.

$\Sigma = \text{ASCII}$

Sea  $\bar{A} \subseteq \Sigma$ , tal que  $\bar{A} = \{a, b, c, \dots, z, A, B, C, \dots, Z\}$ , es decir, solo contiene las letras.

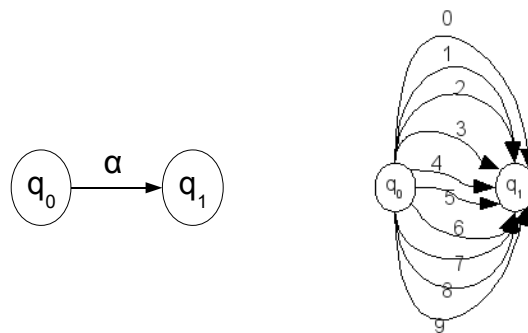
Sea  $\bar{E} \subseteq \Sigma$ , tal que  $\bar{E} = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ , es decir, solo contiene los números.

En el futuro, podremos referirnos a  $\bar{A} \subseteq \Sigma$  como el conjunto LETRA, y a  $\bar{E} \subseteq \Sigma$  como el conjunto NUMERO.

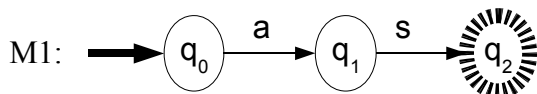
Se puede ver que si  $\alpha \in \bar{A}$ , entonces  $\alpha \in \Sigma$ . Mismo análisis para cualquier  $\beta \in \bar{E}$ .

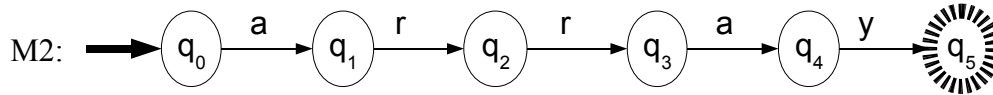
Por simplicidad, usaremos una notación simplificada en los autómatas; si usamos una letra literal, que pertenece a  $\Sigma$ , nos estaremos refiriendo precisamente a esa letra, mientras que si usamos un carácter especial, nos estaremos refiriendo a un conjunto de letras, dependiendo del conjunto del que hablemos. Por ejemplo:

Si  $\alpha \in \text{NUMERO}$ , entonces los dos autómatas de abajo son equivalentes:

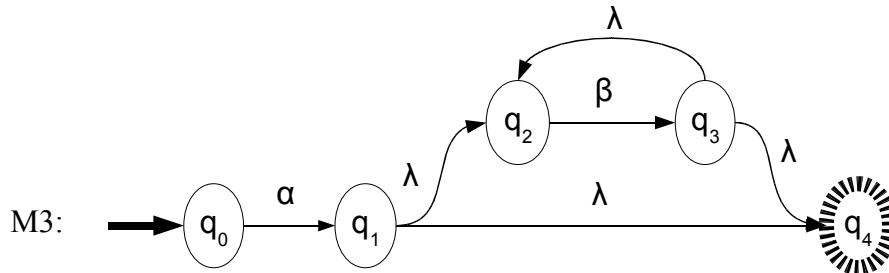


A continuación presentaremos los diagramas de transición de los autómatas finitos, M1, M2 y M3, que reconocen, respectivamente los lenguajes E1, E2 y E3. Trabajaremos sobre:

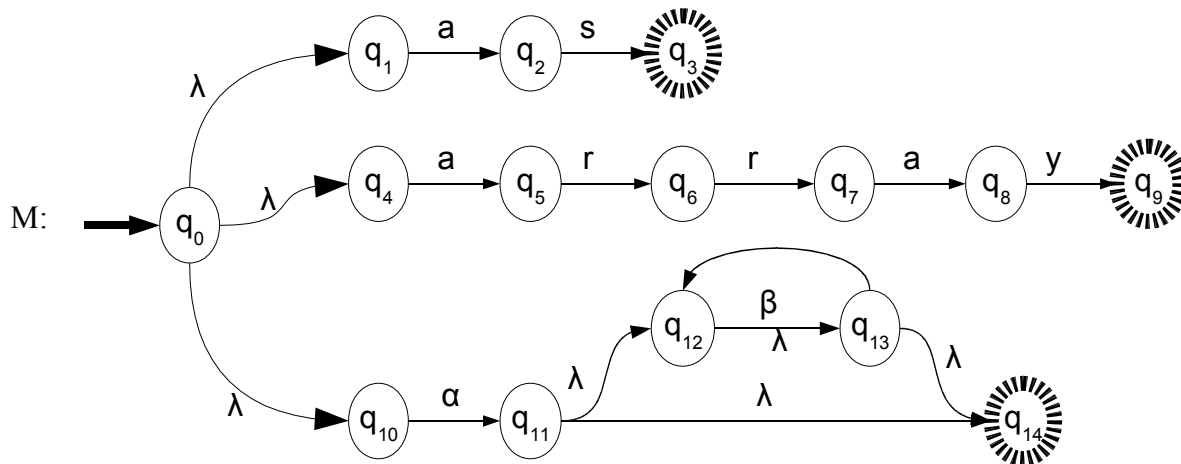




Para  $\alpha \in \text{LETRA}$ ,  $\beta \in (\text{LETRA} \cup \text{NUMERO} \cup \{\_ \})$



Estos tres autómatas son no-determinísticos, obtenidos como resultado de aplicar el algoritmo de conversión de ER a AFND. Si quisiéramos construir un autómata M que reconociera los lenguajes de M1, M2 y M3, es decir,  $L(M) = L(M1) \cup L(M2) \cup L(M3)$ , bastaría con agregar  $\lambda$ -transiciones de la siguiente forma:



A la hora de implementar un analizador lexicográfico, es importante que el autómata M presentado reconozca a que lenguaje pertenece la palabra identificada. Hay que notar que, “as” y “array” son palabras reservadas, y también son identificadores validos, por lo que “as” pertenece a  $L(M1)$  y  $L(M3)$ , mientras que “array” pertenece a  $L(M2)$  y  $L(M3)$ .

Este conflicto se puede resolver aplicando alguna distinción y/o prioridad a los estados finales de la maquina M. Dicho en otras palabras, **los estados finales son diferentes entre si**. Para esta primera maquina M podríamos resolver la prioridad como sigue:

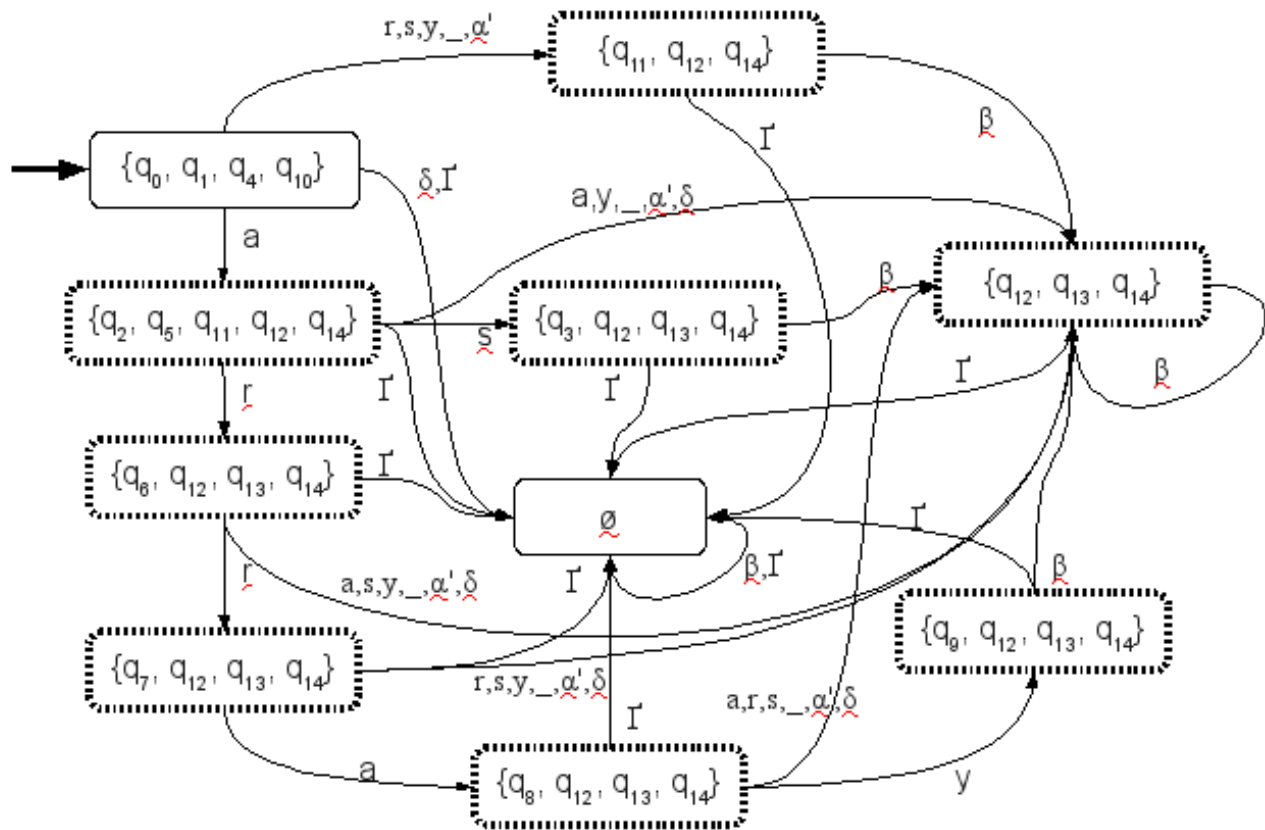
Para dos estados finales  $q_i$  y  $q_p$  cualesquiera,  $q_i$  es más relevante que  $q_p$  si y solo si  $i < p$ . Si al reconocer una palabra se llega a dos estados finales distintos, se tomara únicamente el **estado final más relevante**. Los lenguajes reconocidos serían:

$q_3$  indica que la palabra pertenece a  $L(E1)$

$q_9$  indica que la palabra pertenece a  $L(E2)$

$q_{14}$  indica que la palabra pertenece a  $L(E3)$

A continuación presentaremos el autómata finito determinístico mínimo para  $M$ , pero primero debemos convertirlo en AFD, usando el algoritmo repasado en la practica, resultando:



Con:

- $\alpha \in \text{LETRA}$
- $\delta \in \text{NUMERO}$
- $\beta \in (\text{LETRA} \cup \text{NUMERO} \cup \{\_ \})$
- $\alpha' \in (\text{LETRA} - \{a, r, s, y\})$ , es decir,  $\alpha'$  puede ser cualquier letra, excepto la "a", "r", "s" y "y".
- $\Gamma \in \Sigma - (\text{LETRA} \cup \text{NUMERO} \cup \{\_ \})$ , es decir,  $\Gamma$  es cualquier caracter de  $\Sigma$  que no sea letra o numero o el '\_' (underscore).

Los estados finales son todos aquellos que tienen el borde punteado (todos, excepto el estado inicial y el estado vacío). En este autómata, llámese  $M'$ , debemos mantener la prioridad de estados finales de la misma forma, quedando así:

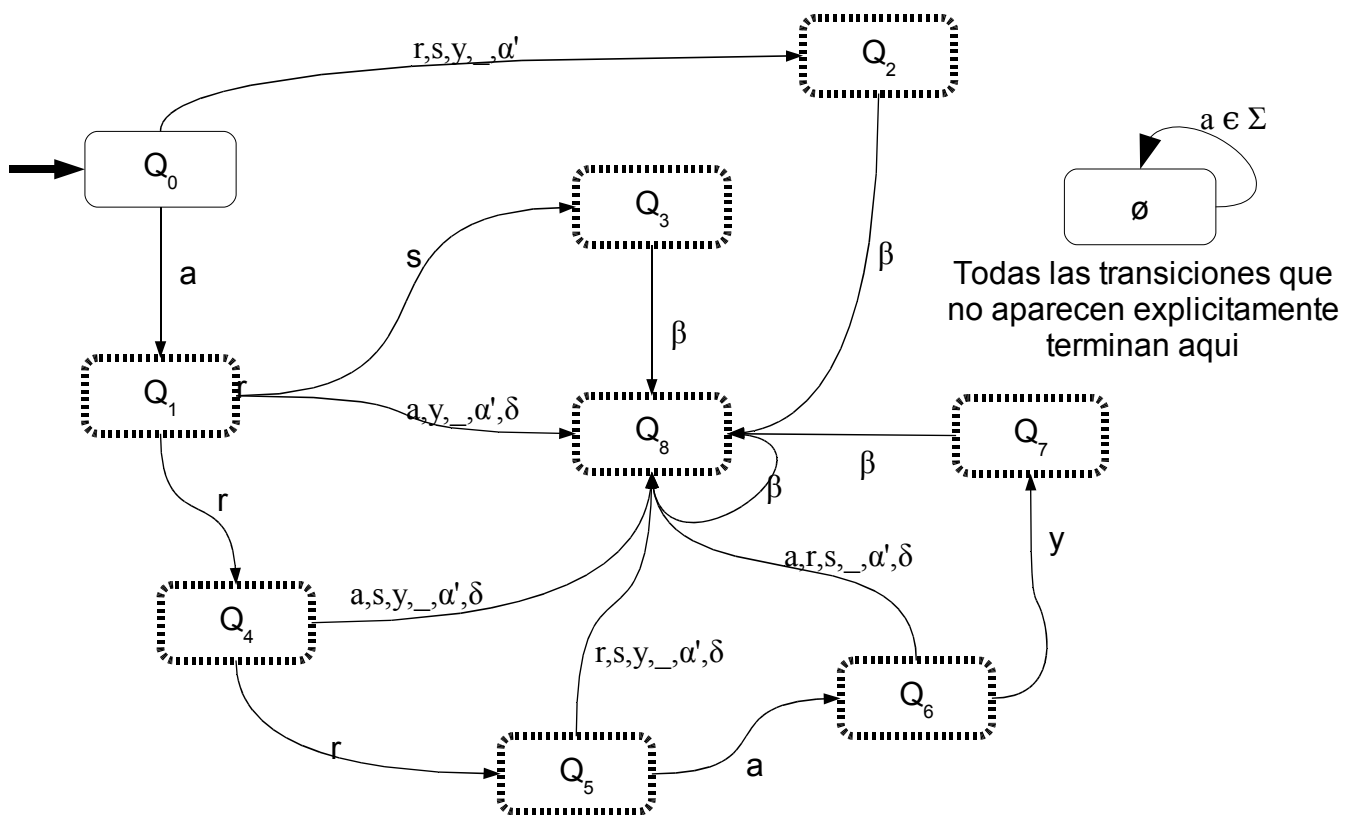
$\{Q_3, q_{12}, q_{13}, q_{14}\}$  Tiene precedencia 3, y reconoce el lenguaje  $L(E1)$ .

$\{Q_9, q_{12}, q_{13}, q_{14}\}$  Tiene precedencia 2, y reconoce el lenguaje  $L(E2)$ .

Todos los estados finales restantes tiene precedencia 1, y reconoce el lenguaje  $L(E3)$ .

Siempre que se tengan dos posibles estados finales, tomaremos aquel que tenga mayor precedencia, y así determinaremos a que lenguaje pertenece la palabra.

Para mayor claridad, vamos a reescribir al AFD  $M'$ , pero sin las transiciones con el estado vacío (asumiremos que si no se dice algo, caemos en el estado vacío), y renombrando los estados:



Con:

- $\alpha \in \text{LETRA}$
- $\delta \in \text{NUMERO}$
- $\beta \in (\text{LETRA} \cup \text{NUMERO} \cup \{\_\})$
- $\alpha' \in (\text{LETRA} - \{a, r, s, y\})$ , es decir,  $\alpha'$  puede ser cualquier letra, excepto la “a”, “r”, “s” y “y”.

En este nuevo grafo (recordemos que es un AFD, previamente definido, solo que no dibujamos las transiciones que caen en el estado vacío), llamémoslo  $M'_Q$ , nos permite visualizar mejor lo que está pasando, y así podemos aplicar el algoritmo de minimización más cómodamente. Vamos a mostrar las clases de equivalencia resultantes al aplicar el algoritmo.

$$\Xi_0 = \{ \{Q_0, \emptyset\}, \{Q_3\}, \{Q_7\}, \{Q_1, Q_2, Q_4, Q_5, Q_6, Q_8\} \}$$

Recordemos que los estados finales son **distinguibles**, lo que significa que no todos son equivalentes:  $Q_3$  reconoce el lenguaje  $L(E1)$ ,  $Q_7$  reconoce el lenguaje  $L(E2)$  y  $Q_1, Q_2, Q_4, Q_5, Q_6, Q_8$  reconocen el lenguaje  $L(E3)$ .

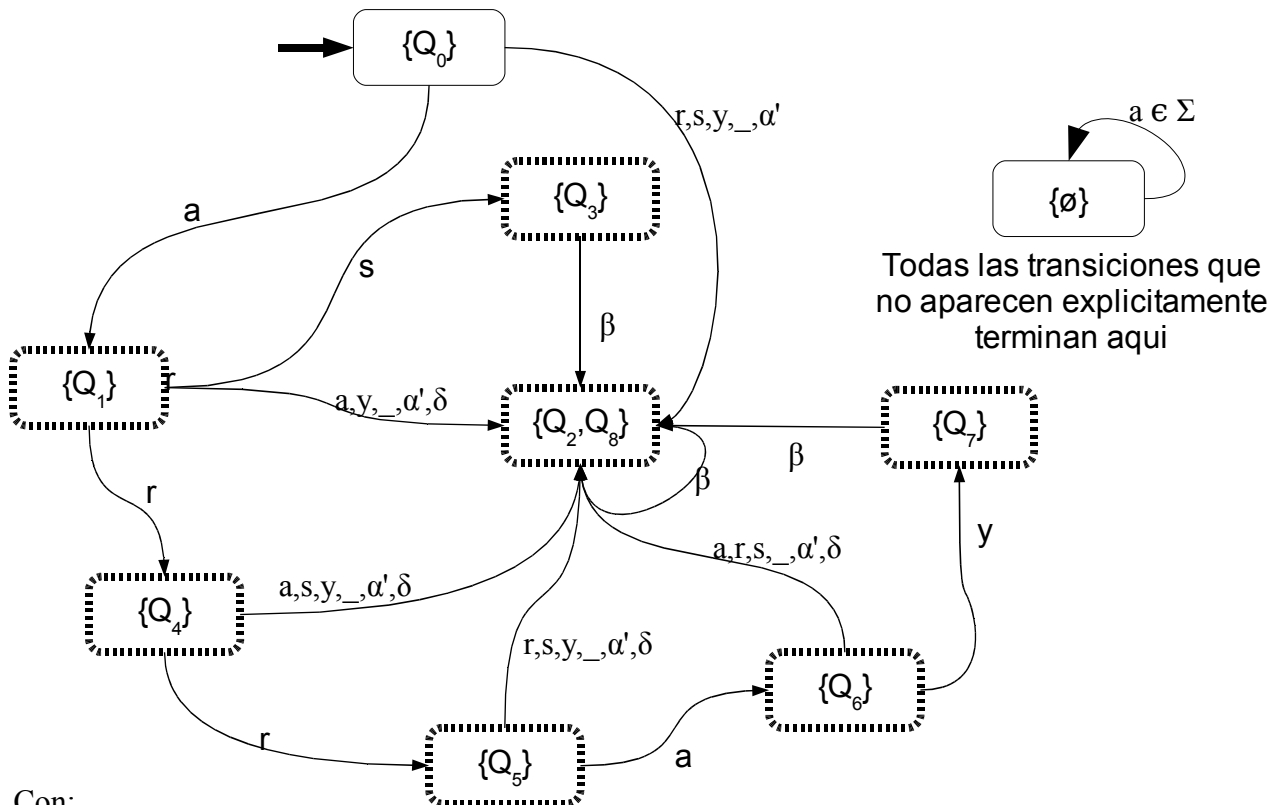
$$\Xi_1 = \{ \{\emptyset\}, \{Q_0\}, \{Q_1\}, \{Q_3\}, \{Q_6\}, \{Q_7\}, \{Q_2, Q_4, Q_5, Q_8\} \}$$

$$\Xi_2 = \{ \{\emptyset\}, \{Q_0\}, \{Q_1\}, \{Q_3\}, \{Q_5\}, \{Q_6\}, \{Q_7\}, \{Q_2, Q_4, Q_8\} \}$$

$$\Xi_3 = \{ \{\emptyset\}, \{Q_0\}, \{Q_1\}, \{Q_3\}, \{Q_4\}, \{Q_5\}, \{Q_6\}, \{Q_7\}, \{Q_2, Q_8\} \}$$

$$\Xi_4 = \{ \{\emptyset\}, \{Q_0\}, \{Q_1\}, \{Q_3\}, \{Q_4\}, \{Q_5\}, \{Q_6\}, \{Q_7\}, \{Q_2, Q_8\} \} = \Xi_3$$

Finalmente, vemos reducidos solo dos estados, quedando el autómata de la siguiente manera:



Con:

- $\alpha \in \text{LETRA}$
- $\delta \in \text{NUMERO}$
- $\beta \in (\text{LETRA} \cup \text{NUMERO} \cup \{\_ \})$
- $\alpha' \in (\text{LETRA} - \{a, r, s, y\})$ , es decir,  $\alpha'$  puede ser cualquier letra, excepto la “a”, “r”, “s” y “y”.

Recordemos, una vez más, que se trata del AFD mínimo, solo que, por simplicidad visual, no colocamos las transiciones que terminan en el estado  $\{\emptyset\}$ .

Para terminar, estos son los estados y los lenguajes asociados a ellos:

$\{Q_3\}$  Reconoce el lenguaje  $L(E1)$ .

$\{Q_7\}$  Reconoce el lenguaje  $L(E2)$ .

$\{Q_1\}$ ,  $\{Q_4\}$ ,  $\{Q_5\}$ ,  $\{Q_6\}$ ,  $\{Q_2, Q_8\}$  Reconocen el lenguaje  $L(E3)$ .