

# Extracción de información

Clase 25

IIC 2223

Prof. Cristian Riveros

# Extracción de información desde documentos

```
18:30 ERROR 06  
19:10 OK 00  
20:00 ERROR 19
```

*"Obtener todas las horas HH:MM"*

$R := (\backslash d \backslash d : \backslash d \backslash d)$

¿cómo podemos **automatizar** esta tarea de extraer datos?

18:30 ERROR 06 ← 19:10 OK 00 ← 20:00 ERROR 19  
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40

## Extracción de información

La tarea de automatizar la extracción de datos desde documentos no estructurados o semi-estructurados.

# Extracción de información desde documentos

18:30	ERROR	06
19:10	OK	00
20:00	ERROR	19

*"Obtener todas las horas HH:MM"*

$R := (\backslash d \backslash d : \backslash d \backslash d)$

¿cómo podemos **automatizar** esta tarea de extraer datos?

1. Una tarea como la anterior **se puede programar fácilmente**, ¿por qué queremos automatizarla?
2. Las expresiones regulares ya hacen el trabajo anterior, ¿por qué queremos **estudiar este problema de nuevo**?

En esta última parte del curso veremos **nuevas técnicas formales y algorítmicas** para entender y resolver este problema.

# Outline

Spans

Regex

Vset automata

Desde regex a VA

# Outline

**Spans**

Regex

Vset automata

Desde regex a VA

# Representación de intervalos: spans

## Definición

- Para un documento  $d = a_0 a_1 \dots a_{n-1}$  se define un **span**  $s$  de  $d$  como:

$$s = [i, j\rangle$$

tal que  $0 \leq i \leq j \leq n$ .

- Si  $s = [i, j\rangle$  es un span de  $d$  se define:

$$d[s] = d[[i, j\rangle] = a_i a_{i+1} \dots a_{j-1}$$

como el **contenido del span**  $s$  en  $d$ . Si  $i = j$ , entonces  $d[[i, i\rangle] = \epsilon$ .

## Ejemplo

18:30 ERROR 06 ↵ 19:10 OK 00 ↵ 20:00 ERROR 19

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40

$s_1$  is highlighted over the span [15, 20) containing the text "19:10".  
 $s_2$  is highlighted over the span [32, 32) which is empty.

$$s_1 = [15, 20\rangle \quad w[s_1] = 19:10 \quad s_2 = [32, 32\rangle \quad w[s_2] = \epsilon$$

# Representación de intervalos: spans

## Definición

- Para un documento  $d = a_0 a_1 \dots a_{n-1}$  se define un **span**  $s$  de  $d$  como:

$$s = [i, j\rangle$$

tal que  $0 \leq i \leq j \leq n$ .

- Si  $s = [i, j\rangle$  es un span de  $d$  se define:

$$d[s] = d[[i, j\rangle] = a_i a_{i+1} \dots a_{j-1}$$

como el **contenido del span**  $s$  en  $d$ . Si  $i = j$ , entonces  $d[[i, i\rangle] = \epsilon$ .

- Denotamos el conjunto de **todos los spans** de  $d$  como  $\text{Spans}(d)$ .

Queremos estudiar como **extraer spans desde documentos**.

# Outline

Spans

**Regex**

Vset automata

Desde regex a VA



# Expresiones regulares con variables

Sea  $\Sigma$  un alfabeto finito y  $\mathcal{X}$  un conjunto de variables.

## Definición (Sintaxis)

$R$  es una expresión regular con variables (o **regex**) sobre  $\Sigma$  si  $R$  es igual a:

1.  $a$  para alguna letra  $a \in \Sigma$ .
2.  $\epsilon$
3.  $x\{R_1\}$  donde  $R_1$  es una regex y  $x \in \mathcal{X}$ .
4.  $(R_1 + R_2)$  donde  $R_1$  y  $R_2$  son regex.
5.  $(R_1 \cdot R_2)$  donde  $R_1$  y  $R_2$  son regex.
6.  $(R_1^*)$  donde  $R_1$  es una regex.

$x\{R_1\}$  representa que el **span capturado** por  $R_1$   
lo almacenaremos en la variable  $x$

# Expresiones regulares con variables

## Ejemplos de regex

Sea  $\Sigma = \{a, b, \sqcup\}$ .

- $\Sigma^* \cdot \mathbf{x}\{abba\} \cdot \Sigma^*$
- $\Sigma^* \cdot \sqcup \cdot \mathbf{x}\{ab \cdot (a + b)^*\} \cdot \sqcup \cdot \Sigma^*$
- $\Sigma^* \cdot \sqcup \cdot \mathbf{x}\{(a + b)^*\} \cdot \sqcup \cdot \mathbf{y}\{(a + b)^*\} \cdot \sqcup \cdot \Sigma^*$
- $(\Sigma^* \cdot \sqcup + \epsilon) \cdot \mathbf{z}\{\mathbf{x}\{(a + b)^+\} \cdot \sqcup \cdot \mathbf{y}\{(a + b)^+\}\} \cdot (\epsilon + \sqcup \cdot \Sigma^*)$
- $\Sigma^* \cdot \mathbf{x}\{abba\} \cdot \sqcup \cdot \mathbf{x}\{baba\} \cdot \Sigma^*$

# Regex válidas

Sea  $\text{Var}(R)$  el conjunto de **todas las variables** en  $\mathcal{X}$  mencionadas en  $R$ .

## Definición

$R$  es una regex **válida** ssi todas las siguientes condiciones se cumplen:

1. si  $R = R_1 \cdot R_2$ , entonces  $\text{Var}(R_1) \cap \text{Var}(R_2) = \emptyset$  y  $R_1, R_2$  son válidas
2. si  $R = R_1 + R_2$ , entonces  $\text{Var}(R_1) = \text{Var}(R_2)$  y  $R_1, R_2$  son válidas
3. si  $R = R_1^*$ , entonces  $\text{Var}(R_1) = \emptyset$  y  $R_1$  es válida
4. si  $R = \mathbf{x}\{R_1\}$ , entonces  $\mathbf{x} \notin \text{Var}(R_1)$  y  $R_1$  es válida

Si  $R$  es **válida** nos permite asignar las variables de **manera correcta**.

# Regex válidas

## Definición regex válidas (resumen)

1. si  $R = R_1 \cdot R_2$ , entonces  $\text{Var}(R_1) \cap \text{Var}(R_2) = \emptyset$  y  $R_1, R_2$  son válidas
2. si  $R = R_1 + R_2$ , entonces  $\text{Var}(R_1) = \text{Var}(R_2)$  y  $R_1, R_2$  son válidas
3. si  $R = R_1^*$ , entonces  $\text{Var}(R_1) = \emptyset$  y  $R_1$  es válida
4. si  $R = \mathbf{x}\{R_1\}$ , entonces  $\mathbf{x} \notin \text{Var}(R_1)$  y  $R_1$  es válida

## Ejemplos

- $\Sigma^* \cdot \mathbf{x}\{abba\} \cdot \Sigma^*$
- $\Sigma^* \cdot \mathbf{x}\{abba\} \cdot \sqcup \cdot \mathbf{y}\{baba\} \cdot \Sigma^*$
- $\Sigma^* \cdot \sqcup \cdot \mathbf{x}\{ab \cdot (a + b^*)\} \cdot \sqcup \cdot \Sigma^*$
- $\Sigma^* \cdot \mathbf{x}\{abba\} \cdot \sqcup \cdot \mathbf{x}\{baba\} \cdot \Sigma^*$
- $\Sigma^* \cdot (\mathbf{x}\{abba\} + \mathbf{y}\{baba\}) \cdot \Sigma^*$
- $\Sigma^* \cdot \sqcup \cdot (\mathbf{x}\{ab \cdot (a + b)^*\} \cdot \sqcup)^* \cdot \Sigma^*$



# Mappings como outputs de regex

## Definiciones

- Un **mapping** de  $R$  sobre un documento  $d \in \Sigma^*$  es una función:

$$\mu : \text{Var}(R) \rightarrow \text{Spans}(d)$$

donde el dominio de  $\mu$  es  $\text{dom}(\mu) = \text{Var}(R)$ .

- Se define el mapping  $\mu = \perp$  como el **mapping vacío** donde  $\text{dom}(\perp) = \emptyset$ .

## Ejemplo



$$[x \mapsto [0, 5], y \mapsto [6, 11], z \mapsto [12, 14]]$$

$$[x \mapsto [15, 20], y \mapsto [21, 23], z \mapsto [24, 26]]$$

$$\perp = [\cdot]$$

# Mappings como outputs de regex

## Definiciones

- Un **mapping** de  $R$  sobre  $d$  es una función  $\mu : \text{Var}(R) \rightarrow \text{Spans}(d)$ .
- Se define el mapping  $\mu = \perp$  como el **mapping vacío** donde  $\text{dom}(\perp) = \emptyset$ .
- Para una variable  $x$  y span  $s$  se define el **mapping de una variable**:

$$\mu = [x \mapsto s] \quad \text{tal que} \quad \text{dom}(\mu) = \{x\} \quad \text{y} \quad \mu(x) = s$$

- Para  $k \in \mathbb{N}$  se define el **mapping**  $\mu + k$  tal que  $\text{dom}(\mu + k) = \text{dom}(\mu)$  y:

$$\text{si } \mu(x) = [i, j] \text{ entonces } [\mu + k](x) = [i + k, j + k].$$

- Para mappings  $\mu_1, \mu_2$  con  $\text{dom}(\mu_1) \cap \text{dom}(\mu_2) = \emptyset$  se define **la unión**:

$$\mu = \mu_1 \cup \mu_2 \quad \text{tal que} \quad \mu(x) = \begin{cases} \mu_1(x) & \text{si } x \in \text{dom}(\mu_1) \\ \mu_2(x) & \text{si } x \in \text{dom}(\mu_2) \end{cases}$$

# Cada regex define un conjunto de mappings

## Semántica regex

Para una regex válida  $R$  cualquiera,  
se define la función  $\llbracket R \rrbracket$  **inductivamente** sobre documentos  $d \in \Sigma^*$ :

1.  $\llbracket a \rrbracket(d) = \{\perp\}$  si  $d = a$ , y  $\emptyset$  en otro caso
2.  $\llbracket \epsilon \rrbracket(d) = \{\perp\}$  si  $d = \epsilon$ , y  $\emptyset$  en otro caso

## Ejemplos

- $\llbracket b \rrbracket(a) = \emptyset$
- $\llbracket b \rrbracket(b) = \{\perp\}$

# Cada regex define un conjunto de mappings

## Semántica regex

1.  $\llbracket a \rrbracket(d) = \{\perp\}$  si  $d = a$ , y  $\emptyset$  en otro caso
2.  $\llbracket \epsilon \rrbracket(d) = \{\perp\}$  si  $d = \epsilon$ , y  $\emptyset$  en otro caso
3.  $\llbracket \mathbf{x}\{R_1\} \rrbracket(d) = \{ \mu \cup [\mathbf{x} \mapsto s] \mid \mu \in \llbracket R_1 \rrbracket(d) \text{ y } s = [0, |d|] \}$

## Ejemplos

- $\llbracket b \rrbracket(a) = \emptyset$
- $\llbracket b \rrbracket(b) = \{\perp\}$
- $\llbracket \mathbf{x}\{b\} \rrbracket(a) = \emptyset$
- $\llbracket \mathbf{x}\{b\} \rrbracket(b) = \{ [\mathbf{x} \mapsto [0, 1]] \}$



# Cada regex define un conjunto de mappings

## Semántica regex

1.  $\llbracket a \rrbracket(d) = \{\perp\}$  si  $d = a$ , y  $\emptyset$  en otro caso
2.  $\llbracket \epsilon \rrbracket(d) = \{\perp\}$  si  $d = \epsilon$ , y  $\emptyset$  en otro caso
3.  $\llbracket \mathbf{x}\{R_1\} \rrbracket(d) = \{ \mu \cup [\mathbf{x} \mapsto s] \mid \mu \in \llbracket R_1 \rrbracket(d) \text{ y } s = [0, |d|] \}$

Como  $\mathbf{x}\{R_1\}$  es **válida**, entonces sabemos que  $\mathbf{x} \notin \text{dom}(\mu)$ .

# Cada regex define un conjunto de mappings

## Semántica regex

$$3. \llbracket \mathbf{x}\{R_1\} \rrbracket(d) = \{ \mu \cup [\mathbf{x} \mapsto s] \mid \mu \in \llbracket R_1 \rrbracket(d) \text{ y } s = [0, |d|] \}$$

$$4. \llbracket R_1 \cdot R_2 \rrbracket(d) = \left\{ \mu_1 \cup (\mu_2 + |d_1|) \mid \begin{array}{l} \text{existe } d_1, d_2 \text{ tal que } d = d_1 \cdot d_2, \\ \mu_1 \in \llbracket R_1 \rrbracket(d_1) \text{ y } \mu_2 \in \llbracket R_2 \rrbracket(d_2) \end{array} \right\}$$

## Ejemplos

- $\llbracket b \rrbracket(b) = \{\perp\}$
- $\llbracket \mathbf{x}\{b\} \rrbracket(b) = \{ [\mathbf{x} \mapsto [0, 1]] \}$
- $\llbracket \mathbf{x}\{b\}b \rrbracket(bb) = \{ [\mathbf{x} \mapsto [0, 1]] \}$
- $\llbracket b\mathbf{x}\{b\} \rrbracket(bb) = \{ [\mathbf{x} \mapsto [1, 2]] \}$
- $\llbracket \mathbf{x}\{b\}\mathbf{y}\{b\} \rrbracket(bb) = \{ [\mathbf{x} \mapsto [0, 1], \mathbf{y} \mapsto [1, 2]] \}$

# Cada regex define un conjunto de mappings

## Semántica regex

$$3. \llbracket \mathbf{x}\{R_1\} \rrbracket(d) = \{ \mu \cup [\mathbf{x} \mapsto s] \mid \mu \in \llbracket R_1 \rrbracket(d) \text{ y } s = [0, |d|] \}$$

$$4. \llbracket R_1 \cdot R_2 \rrbracket(d) = \left\{ \mu_1 \cup (\mu_2 + |d_1|) \mid \begin{array}{l} \text{existe } d_1, d_2 \text{ tal que } d = d_1 \cdot d_2, \\ \mu_1 \in \llbracket R_1 \rrbracket(d_1) \text{ y } \mu_2 \in \llbracket R_2 \rrbracket(d_2) \end{array} \right\}$$

Como  $R_1 \cdot R_2$  son válidas, entonces sabemos que  $\text{dom}(\mu_1) \cap \text{dom}(\mu_2) = \emptyset$

# Cada regex define un conjunto de mappings

## Semántica regex

$$3. \llbracket \mathbf{x}\{R_1\} \rrbracket(d) = \{ \mu \cup [\mathbf{x} \mapsto s] \mid \mu \in \llbracket R_1 \rrbracket(d) \text{ y } s = [0, |d|] \}$$

$$4. \llbracket R_1 \cdot R_2 \rrbracket(d) = \left\{ \mu_1 \cup (\mu_2 + |d_1|) \mid \begin{array}{l} \text{existe } d_1, d_2 \text{ tal que } d = d_1 \cdot d_2, \\ \mu_1 \in \llbracket R_1 \rrbracket(d_1) \text{ y } \mu_2 \in \llbracket R_2 \rrbracket(d_2) \end{array} \right\}$$

$$5. \llbracket R_1 + R_2 \rrbracket(d) = \llbracket R_1 \rrbracket(d) \cup \llbracket R_2 \rrbracket(d)$$

## Ejemplos

$$\blacksquare \llbracket \mathbf{x}\{b\}b \rrbracket(bb) = \{ [\mathbf{x} \mapsto [0, 1]] \}$$

$$\blacksquare \llbracket b\mathbf{x}\{b\} \rrbracket(bb) = \{ [\mathbf{x} \mapsto [1, 2]] \}$$

$$\blacksquare \llbracket \mathbf{x}\{b\}b + b\mathbf{x}\{b\} \rrbracket(bb) = \{ [\mathbf{x} \mapsto [0, 1]], [\mathbf{x} \mapsto [1, 2]] \}$$

# Cada regex define un conjunto de mappings

## Semántica regex

$$3. \llbracket \mathbf{x}\{R_1\} \rrbracket(d) = \{ \mu \cup [\mathbf{x} \mapsto s] \mid \mu \in \llbracket R_1 \rrbracket(d) \text{ y } s = [0, |d|] \}$$

$$4. \llbracket R_1 \cdot R_2 \rrbracket(d) = \left\{ \mu_1 \cup (\mu_2 + |d_1|) \mid \begin{array}{l} \text{existe } d_1, d_2 \text{ tal que } d = d_1 \cdot d_2, \\ \mu_1 \in \llbracket R_1 \rrbracket(d_1) \text{ y } \mu_2 \in \llbracket R_2 \rrbracket(d_2) \end{array} \right\}$$

$$5. \llbracket R_1 + R_2 \rrbracket(d) = \llbracket R_1 \rrbracket(d) \cup \llbracket R_2 \rrbracket(d)$$

$$6. \llbracket R_1^* \rrbracket(d) = \bigcup_{k=0}^{\infty} \llbracket (R_1)^k \rrbracket(d)$$

Como  $\text{Var}(R_1) = \emptyset$ , entonces  $\llbracket R_1^* \rrbracket(d) = \{\perp\}$  ssi  $d \in \mathcal{L}(R_1^*)$ .

# Cada regex define un conjunto de mappings

## Semántica regex (completa)

Para una regex válida  $R$  cualquiera,  
se define la función  $\llbracket R \rrbracket$  **inductivamente** sobre documentos  $d \in \Sigma^*$ :

1.  $\llbracket a \rrbracket(d) = \{\perp\}$  si  $d = a$ , y  $\emptyset$  en otro caso
2.  $\llbracket \epsilon \rrbracket(d) = \{\perp\}$  si  $d = \epsilon$ , y  $\emptyset$  en otro caso
3.  $\llbracket \mathbf{x}\{R_1\} \rrbracket(d) = \{ \mu \cup [\mathbf{x} \mapsto s] \mid \mu \in \llbracket R_1 \rrbracket(d) \text{ y } s = [0, |d|] \}$
4.  $\llbracket R_1 \cdot R_2 \rrbracket(d) = \left\{ \mu_1 \cup (\mu_2 + |d_1|) \mid \begin{array}{l} \text{existe } d_1, d_2 \text{ tal que } d = d_1 \cdot d_2, \\ \mu_1 \in \llbracket R_1 \rrbracket(d_1) \text{ y } \mu_2 \in \llbracket R_2 \rrbracket(d_2) \end{array} \right\}$
5.  $\llbracket R_1 + R_2 \rrbracket(d) = \llbracket R_1 \rrbracket(d) \cup \llbracket R_2 \rrbracket(d)$
6.  $\llbracket R_1^* \rrbracket(d) = \bigcup_{k=0}^{\infty} \llbracket (R_1)^k \rrbracket(d)$

# Evaluación de regex

¿cuál es el resultado de las siguientes regex?

$$d = \begin{array}{cccccccccccccccc} \text{a} & \text{a} & \text{b} & \text{b} & \text{a} & & \text{b} & \text{a} & & \text{b} & \text{a} & \text{b} & \text{a} & \text{b} & \text{a} \\ 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \end{array}$$

- $\llbracket \Sigma^* \cdot \mathbf{x}\{aba + bab\} \cdot \Sigma^* \rrbracket(d) = \{ [\mathbf{x} \mapsto [10, 13]], [\mathbf{x} \mapsto [11, 14]] \}$
- $\llbracket \Sigma^* \cdot \sqcup \cdot \mathbf{x}\{ba \cdot (a + b)^*\} \cdot \sqcup \cdot \Sigma^* \rrbracket(d) = \{ [\mathbf{x} \mapsto [7, 9]], [\mathbf{x} \mapsto [10, 14]] \}$
- $\llbracket \Sigma^* \cdot \mathbf{x}\{abba\} \cdot \sqcup \cdot \mathbf{y}\{ba\} \cdot \Sigma^* \rrbracket(d) = \{ [\mathbf{x} \mapsto [2, 6], \mathbf{y} \mapsto [7, 9]] \}$
- $\llbracket \Sigma^* \cdot \sqcup \cdot \mathbf{x}\{(a + b)^*\} \cdot \sqcup \cdot \mathbf{y}\{(a + b)^*\} \cdot \sqcup \cdot \Sigma^* \rrbracket(d) =$   
 $\{ [\mathbf{x} \mapsto [2, 6], \mathbf{y} \mapsto [7, 9]], [\mathbf{x} \mapsto [7, 9], \mathbf{y} \mapsto [10, 14]] \}$
- $\llbracket \Sigma^* \cdot \mathbf{x}\{\Sigma^*\} \cdot \Sigma^* \rrbracket(d) = \text{Spans}(d)$

# Evaluación de regex

¿cuál es el resultado de las siguientes regex?

$$d = \begin{array}{cccccccccccccccc} a & a & b & b & a & b & a & b & a & b & a & b & a & b & a \\ \hline 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \end{array}$$

- $\llbracket \Sigma^* \cdot x \{ \Sigma^* \} \cdot \Sigma^* \rrbracket(d) = \text{Spans}(d)$
- Dado una regex  $R$  con una variable, ¿de qué tamaño puede ser  $\llbracket R \rrbracket(d)$  con respecto a  $|d|$  **en el peor caso**?
- Dado una regex  $R$  con  $k$  variables, ¿de qué tamaño puede ser  $\llbracket R \rrbracket(d)$  con respecto a  $|d|$  y  $k$  **en el peor caso**?

Entonces, ¿cómo podemos evaluar una regex  $R$  **eficientemente**?



# Outline

Spans

Regex

**Vset automata**

Desde regex a VA

# Autómata con variables (vset automata)

¿qué tienen de nuevo?

1. tiene transiciones con **abre** y **cierra** de variable  $x$ :

$$p \xrightarrow{x} q \quad y \quad p \xrightarrow{x} q$$

2. cada ejecución define un **mapping** de las variables a spans.

**Vset autómatas** será nuestro primer modelo para **compilar regex**

# Autómata con variables (vset automata)

## Definición

Un vset automata (VA) es una tupla:

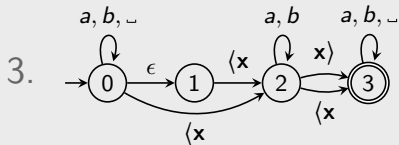
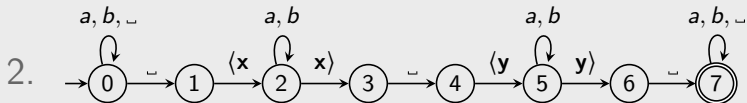
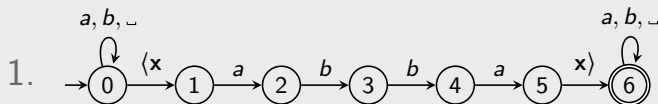
$$\mathcal{A} = (Q, \Sigma, \mathcal{X}, \Delta, I, F)$$

- $Q$  es un conjunto finito de estados.
- $\Sigma$  es el alfabeto de input.
- $\mathcal{X}$  es un conjunto finito de variables.
- $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\} \cup \{\langle \mathbf{x}, \mathbf{x} \rangle \mid \mathbf{x} \in \mathcal{X}\}) \times Q$  es la relación de transición.
- $I \subseteq Q$  es un conjunto de estados iniciales.
- $F \subseteq Q$  es el conjunto de estados finales (o aceptación).

‘ $\langle \mathbf{x} \rangle$ ’ simboliza **abrir** y ‘ $\mathbf{x}$ ’ simboliza **cerrar** la variable  $\mathbf{x}$

# Autómata con variables (vset automata)

## Ejemplos de vset autómata



# Ejecución de un vset autómatata (VA)

Sea  $\mathcal{A} = (Q, \Sigma, \mathcal{X}, \Delta, I, F)$  un VA y  $d = a_0 \dots a_{n-1} \in \Sigma^*$  un documento.

## Definiciones

- Un par  $(p, i) \in Q \times \{0, \dots, n\}$  es una **configuración** de  $\mathcal{A}$  sobre  $d$ .
- Una configuración  $(p, 0)$  es **inicial** si  $q \in I$ .
- Una configuración  $(p, |d|)$  es **final** si  $q \in F$ .

*“Intuitivamente, una configuración  $(p, i)$  representa que  $\mathcal{A}$  se encuentra en el estado  $p$  antes de leer  $a_i$ .”*

# Ejecución de un vset autómatata (VA)

Sea  $\mathcal{A} = (Q, \Sigma, \mathcal{X}, \Delta, I, F)$  un VA y  $d = a_0 \dots a_{n-1} \in \Sigma^*$  un documento.

## Definiciones

Una **ejecución** (o run)  $\rho$  de  $\mathcal{A}$  sobre  $d$  es una secuencia:

$$\rho : (p_0, i_0) \xrightarrow{o_1} (p_1, i_1) \xrightarrow{o_2} \dots \xrightarrow{o_m} (p_m, i_m)$$

tal que cumple todas las siguientes condiciones:

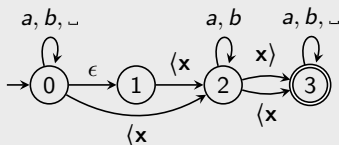
- $o_k \in \Sigma \cup \{\epsilon\} \cup \{\langle \mathbf{x}, \mathbf{x} \rangle \mid \mathbf{x} \in \mathcal{X}\}$  con  $k \leq m$
- $(p_0, i_0)$  es una configuración inicial
- para todo  $k < m$ ,  $(p_k, o_{k+1}, p_{k+1}) \in \Delta$
- para todo  $k \leq m$ , si  $o_k \in \Sigma$ , entonces  $o_k = a_{i_{k-1}}$  y  $i_k = i_{k-1} + 1$
- para todo  $k \leq m$ , si  $o_k \in \{\epsilon\} \cup \{\langle \mathbf{x}, \mathbf{x} \rangle \mid \mathbf{x} \in \mathcal{X}\}$ , entonces  $i_k = i_{k-1}$ .

Una ejecución  $\rho$  es de **aceptación** si  $(p_m, i_m)$  es de aceptación.

# Ejecución de un vset autómatata (VA)

## Ejemplos de ejecuciones

$$d = \frac{a \ b \ a}{0 \ 1 \ 2}$$



Algunas ejecuciones sobre  $d$ :

$$\rho_1 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{b} (p_0, 2) \xrightarrow{a} (p_0, 3)$$

$$\rho_2 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{\epsilon} (p_1, 1) \xrightarrow{\langle x \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{\langle x \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

$$\rho_3 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{\langle x \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{\langle x \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

$$\rho_4 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{\langle x \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{\langle x \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

¿cuáles ejecuciones son de **aceptación**? ¿cuáles **deberían** dar output?

# Ejecución válida de un vset autómatas (VA)

Sea  $\mathcal{A} = (Q, \Sigma, \mathcal{X}, \Delta, I, F)$  un VA y  $d = a_0 \dots a_{n-1} \in \Sigma^*$  un documento.

## Definiciones

Para una ejecución  $\rho$  de  $\mathcal{A}$  sobre  $d$ :

$$\rho : (p_0, i_0) \xrightarrow{o_1} (p_1, i_1) \xrightarrow{o_2} \dots \xrightarrow{o_m} (p_m, i_m)$$

decimos que  $\rho$  es **válida** si, y solo si, para todo  $x \in \mathcal{X}$ :

- existe un único  $k_1 \leq m$  tal que  $o_{k_1} = \langle x$
- existe un único  $k_2 \leq m$  tal que  $o_{k_2} = x \rangle$  y
- $k_1 < k_2$ .

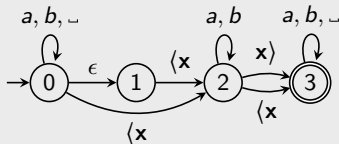
*“Intuitivamente,  $\rho$  es válida si, y solo si, todas las variables se abren y cierran correctamente.”*



# Ejecución válida de un vset autómeta (VA)

¿cuáles ejecuciones son válidas?

$$d = \frac{a \ b \ a}{\quad \quad \quad}$$



Algunas ejecuciones sobre  $d$ :

$$\rho_1 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{b} (p_0, 2) \xrightarrow{a} (p_0, 3)$$

$$\rho_2 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{\epsilon} (p_1, 1) \xrightarrow{\langle x \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{x \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

$$\rho_3 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{\langle x \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{x \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

$$\rho_4 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{\langle x \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{\langle x \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

# El mapping de una ejecución válida

Sea  $\mathcal{A} = (Q, \Sigma, \mathcal{X}, \Delta, I, F)$  un VA y  $d = a_0 \dots a_{n-1} \in \Sigma^*$  un documento.

## Definiciones

Para una ejecución  $\rho$  de  $\mathcal{A}$  sobre  $d$ :

$$\rho : (p_0, i_0) \xrightarrow{o_1} (p_1, i_1) \xrightarrow{o_2} \dots \xrightarrow{o_m} (p_m, i_m)$$

decimos que  $\rho$  es **válida** si, y solo si, para todo  $\mathbf{x} \in \mathcal{X}$ :

- existe un único  $k_1 \leq m$  tal que  $o_{k_1} = \langle \mathbf{x}$
- existe un único  $k_2 \leq m$  tal que  $o_{k_2} = \mathbf{x} \rangle$  y
- $k_1 < k_2$ .

Si  $\rho$  es **válido** se define el **mapping de  $\rho$**   $\text{map}(\rho) : \mathcal{X} \rightarrow \text{Spans}(d)$  tal que:

$$[\text{map}(\rho)](\mathbf{x}) = [i_{k_1}, i_{k_2}]$$

para todo  $\mathbf{x} \in \mathcal{X}$  y  $k_1, k_2 \leq m$  con  $o_{k_1} = \langle \mathbf{x}$  y  $o_{k_2} = \mathbf{x} \rangle$ .

# El mapping de una ejecución válida

Si  $\rho$  es válido se define el **mapping de  $\rho$**   $\text{map}(\rho) : \mathcal{X} \rightarrow \text{Spans}(d)$  tal que:

$$[\text{map}(\rho)](\mathbf{x}) = [i_{k_1}, i_{k_2}]$$

para todo  $\mathbf{x} \in \mathcal{X}$  y  $k_1, k_2 \leq m$  con  $o_{k_1} = \langle \mathbf{x}$  y  $o_{k_2} = \mathbf{x} \rangle$ .

¿cuál es el mapping de las siguientes ejecuciones válidas?

$$\rho_2 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{e} (p_1, 1) \xrightarrow{\langle \mathbf{x} \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{\langle \mathbf{x} \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

$$\rho_3 : (p_0, 0) \xrightarrow{a} (p_0, 1) \xrightarrow{\langle \mathbf{x} \rangle} (p_2, 1) \xrightarrow{b} (p_2, 2) \xrightarrow{\langle \mathbf{x} \rangle} (p_3, 2) \xrightarrow{a} (p_3, 3)$$

$$\text{map}(\rho_2) = \text{map}(\rho_3) = [\mathbf{x} \mapsto [1, 2]]$$

# Función de extracción de un vset autómatas

Sea  $\mathcal{A} = (Q, \Sigma, \mathcal{X}, \Delta, I, F)$  un vset autómatas.

## Definición

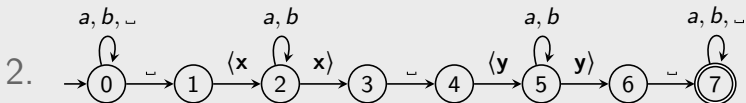
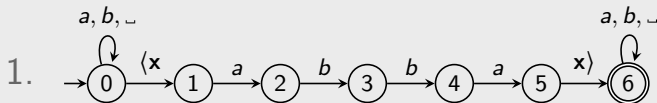
Se define la función  $\llbracket \mathcal{A} \rrbracket$  tal que para todo documento  $d \in \Sigma^*$ :

$$\llbracket \mathcal{A} \rrbracket(d) = \left\{ \text{map}(\rho) \mid \begin{array}{l} \rho \text{ es una ejecución} \\ \textbf{válida} \text{ y de } \textbf{aceptación} \text{ de } \mathcal{A} \text{ sobre } d \end{array} \right\}$$

VA nos entrega otra forma de extraer información de un documento.

# Función de extracción de un vset autómatata

¿qué función de extracción define cada VA?



# Outline

Spans

Regex

Vset automata

Desde regex a VA

# Vset automata funcionales

Sea  $\mathcal{A} = (Q, \Sigma, \mathcal{X}, \Delta, I, F)$  un vset autómat.

$$\llbracket \mathcal{A} \rrbracket(d) = \left\{ \text{map}(\rho) \mid \begin{array}{l} \rho \text{ es una ejecución} \\ \textbf{válida} \text{ y de } \textbf{aceptación} \text{ de } \mathcal{A} \text{ sobre } d \end{array} \right\}$$

## Definición

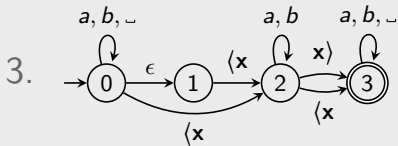
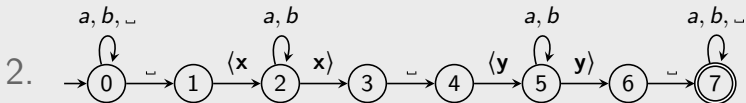
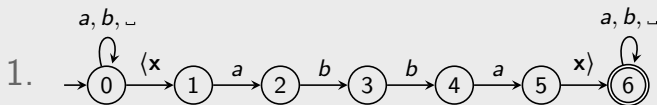
Decimos que un vset autómat  $\mathcal{A}$  es **funcional** si, y solo si, para todo documento  $d$  y para toda ejecución  $\rho$  de  $\mathcal{A}$  sobre  $d$ :

si  $\rho$  es de **aceptación**, entonces  $\rho$  es **válida**.

Para funcional solo necesitamos verificar que la ejecución es de aceptación.

# Vset automata funcionales

¿cuáles VA son funcionales?





# Desde regex a vset autómatas

## Teorema

Para toda regex  $R$  válida, existe un vset autómatas **funcional**  $\mathcal{A}_R$  de **tamaño lineal** en  $|R|$  tal que para todo documento  $d$ :

$$\llbracket R \rrbracket(d) = \llbracket \mathcal{A}_R \rrbracket(d)$$

Demostración: ejercicio (similar al Teorema de Kleene)

¿es verdad la otra dirección?