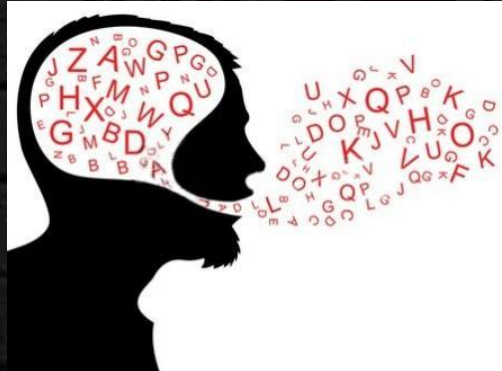




Gramáticas



Compiladores
Prof. Luis Enrique Hernández Olvera



Contenido

- Definición de gramática.
- Elementos de una gramática.
- Clasificación de gramáticas.
- Gramáticas independientes del contexto.
- Derivaciones.

Compiladores
Prof. Luis Enrique Hernández Olvera

2



Recordando

- Las **expresiones regulares** nos proporcionan una plantilla o patrón para las cadenas del lenguaje.
- Un **autómata finito** es un generador de cadenas del lenguaje. Especifica un lenguaje como el conjunto de todas las cadenas que lo hacen pasar del estado inicial a uno de aceptación.



Recordando

- Las expresiones regulares y los autómatas finitos nos proporcionan dos medios para especificar o definir lenguajes.
- Por lo tanto en ambos casos, todas las cadenas se corresponderán con un patrón en particular y dichas cadenas serán las únicas que formaran dicho lenguaje.



Gramática

- La **gramática** es el estudio de las reglas y principios que regulan el uso de las lenguas y la organización de las palabras dentro de una oración. También se denomina así al “conjunto de reglas y principios que gobiernan el uso de un lenguaje” así, cada lenguaje tiene su propia gramática.

Compiladores
Prof. Luis Enrique Hernández Olvera

5



Gramática

Una Gramática G se representa con una cuádrupla:

$$G=(N,\Sigma,S,P)$$

Donde:

- N es una colección finita de no terminales.
- Σ es un alfabeto (Conjunto de terminales).
- S es un no terminal llamado **Símbolo inicial**.
- P es una colección finita de **reglas de sustitución** llamadas **producciones**.

Compiladores
Prof. Luis Enrique Hernández Olvera

6



Elementos de una gramática

- **Símbolos terminales:** son elementos del alfabeto que no se pueden transformar, por eso se llaman terminales. Normalmente se denotan por letras minúsculas.
- **Variables o símbolos no terminales:** son elementos auxiliares que permiten poner restricciones sintácticas a un lenguaje. Las variables sí se pueden transformar, utilizando las reglas, en una cadena de variables y/o terminales. Por lo general se denotan por letras mayúsculas o por la notación *<variable>*.



Elementos de una gramática

- **Símbolo inicial:** es el símbolo a partir del cual se generan todas las palabras válidas.
- **Producciones:** permiten reemplazar variables para generar oraciones válidas de un lenguaje. Puede haber varias reglas de producción para una misma variable, en algunos casos, las distintas opciones se colocan en una sola regla con los distintos reemplazos separados por "|".
 - $\alpha \rightarrow \beta \mid \gamma \mid \delta$
 - abrevia las tres reglas
 - $\alpha \rightarrow \beta$
 - $\alpha \rightarrow \gamma$
 - $\alpha \rightarrow \delta$



Producciones

Cada producción consta de:

- Una variable a la que define (parcialmente) la producción. Esta variable a menudo se denomina cabeza de la producción.
- El símbolo de producción \rightarrow ("puede ser", "se compone de" o "es sustituido por").



Producciones

- Una cadena formada por cero o mas símbolos terminales y variables. Esta cadena, denominada *cuerpo* de la producción, representa una manera de formar cadenas pertenecientes al lenguaje de la variable de la cabeza. De este modo, dejamos solo los símbolos terminales y sustituimos cada una de las variables del cuerpo por una cadena que sabemos que pertenece al lenguaje de dicha variable



Ejemplo 1: Gramática

- $G = (N, \Sigma, S, P)$
 - $N = \{S, A, B\}$
 - $\Sigma = \{a, b, c\}$
 - $P: S \rightarrow AccA \quad A \rightarrow BA \mid \lambda \quad B \rightarrow a \mid b \mid c$

$w1 = abcc \in L(G)$

$w2 = acb \notin L(G)$



Ejemplo 1: Gramática

- $G = (N, \Sigma, S, P)$
 - $N = \{S, A, B\}$
 - $\Sigma = \{a, b, c\}$
 - $P: S \rightarrow AccA \quad A \rightarrow BA \mid \lambda \quad B \rightarrow a \mid b \mid c$

$w1 = abcc \in L(G)$

$w2 = acb \notin L(G)$

- $S \Rightarrow AccA \Rightarrow BAccA \Rightarrow aAccA \Rightarrow aBAccA \Rightarrow$
 $abAccA \Rightarrow ab\lambda ccA \Rightarrow abcc\lambda \Rightarrow abcc$



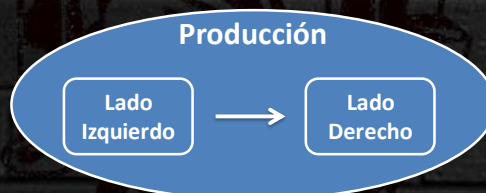
Lenguaje generado por una gramática

- El lenguaje $L(G)$ generado por una gramática G es el conjunto de todas las **cadenas** que puede generar G .
- Una **cadena** pertenece a $L(G)$ si :
 - Está compuesta de símbolos terminales
 - La cadena puede derivarse del símbolo inicial S aplicando las reglas de producción de la gramática.



Gramáticas Formales

- Una gramática formal consta de un conjunto finito de *símbolos terminales* (las palabras en un lenguaje formal), un conjunto finito de *símbolos no terminales*, un conjunto de *reglas de producción* con un lado izquierdo y otro derecho, y un *símbolo inicial* (No terminal " S ").





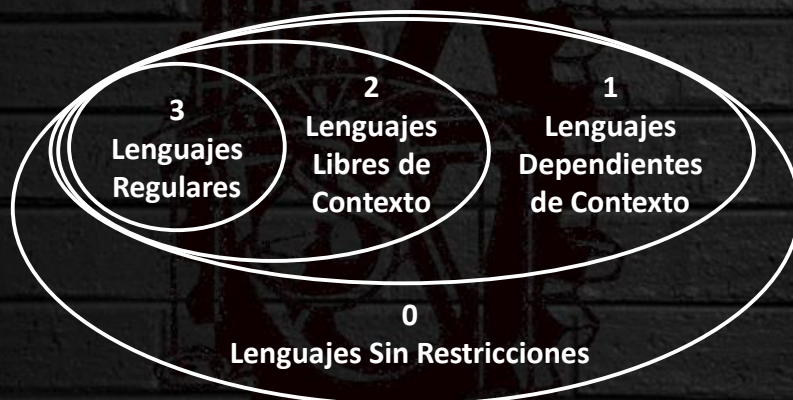
Gramáticas Formales

- Las reglas se aplican sustituyendo la parte de la izquierda por la parte de la derecha. Una derivación es una secuencia de aplicaciones de reglas.
- Cada gramática define el lenguaje formal de todas las sentencias que están formadas exclusivamente por los símbolos terminales a los que se puede llegar mediante derivación a partir del símbolo inicial.



Clasificación de Gramáticas

- Los cuatro tipos son:





Clasificación de Gramáticas

- En función de la forma de sus producciones, se puede caracterizar qué tan compleja es una gramática formal.
 - Gramáticas Tipo 0 (*sin restricciones*)
 - Gramáticas Tipo 1 (*dependientes de contexto*)
 - Gramáticas Tipo 2 (*independientes o libres de contexto*)
 - Gramáticas Tipo 3 (*gramáticas regulares*)



Gramáticas Tipo 3 (*gramáticas regulares*)

- Generan los **lenguajes regulares**. Las reglas (producciones) se restringen a un único no terminal en la parte izquierda y una parte derecha compuesta por un único terminal que puede estar seguido o no de un único no terminal. Es decir, normas del tipo:

$$A \rightarrow aB$$

$$A \rightarrow a$$

- Estos lenguajes son los que pueden ser decididos por un **autómata finito (*regular*)**. Los lenguajes regulares se utilizan para definir **estructura léxica** de los lenguajes de programación. **Definen la sintaxis** de los **identificadores**, **números**, **cadenas** y otros **elementos básicos** del lenguaje.



Gramáticas Tipo 2 (*independientes o libres de contexto*)

- Generan los lenguajes libres de contexto. Están definidas por reglas de la forma:

$$A \rightarrow \gamma$$

- A es un no terminal
- γ es una cadena de terminales y no terminales.
- Se denominan independientes de contexto porque A puede sustituirse por γ independientemente de las cadenas por las que esté acompañada.
- Estos lenguajes son todos los lenguajes que pueden ser reconocidos por los autómatas de pila.
- Los lenguajes independientes de contexto constituyen la base teórica para la sintaxis de la mayoría de los lenguajes de programación. Definen la sintaxis de las declaraciones, las proposiciones, las expresiones, etc. (i.e. la estructura)



Gramáticas Tipo 1 (*dependientes de contexto*)

- Generan los lenguajes dependientes de contexto. Contienen reglas de producción de la forma:

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

- A es un no terminal
- α , β y γ son cadenas de terminales y no terminales.
- α y β pueden ser vacíos, pero γ ha de ser distinto del vacío.
- Se denominan gramáticas dependientes del contexto, porque, como se observa, A puede ser sustituido por γ si está acompañada de α por la izquierda y de β por la derecha.
- Estos lenguajes son todos los lenguajes que pueden ser reconocidos por autómatas lineales acotados (Maquina de Turing Determinista).



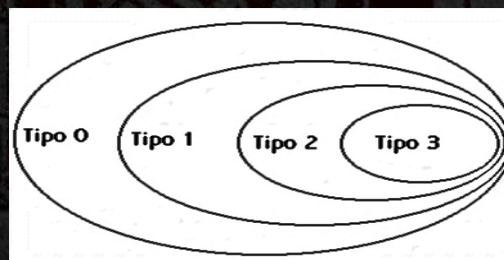
Gramáticas Tipo 0 (*sin restricciones*)

- Incluyen todas las gramáticas formales.
- El más general, al que pertenece la semántica de los lenguajes naturales y artificiales.
- A estos lenguajes no se les impone restricción alguna.
- Estos lenguajes son todos los lenguajes que pueden ser reconocidos por una *máquina de Turing*.



Clasificación de Gramáticas

- Todo lenguaje de tipo 3 es de tipo 2, todo lenguaje de tipo 2 es de tipo 1, y todo lenguaje de tipo 1 es de tipo 0.
- Se dice que un lenguaje es de tipo k [$k = 0, k = 1, k = 2, k = 3$] cuando existe una gramática de tipo k que genera ese lenguaje.





Clasificación de Gramáticas

- Para clasificar una gramática hemos de analizar una a una todas sus reglas de producción obteniendo el tipo de cada una de ellas. La clasificación de la gramática será la correspondiente al tipo de la producción de menor clasificación



Gramática	Lenguaje	Reglas de Producción	Si $\mu \rightarrow \varphi$, relación entre $ \mu $ y $ \varphi $	Solución
Tipo-0	Rekursivas	Sin restricciones		Máquinas de Turing
Tipo-1	Dependiente de contexto	$\alpha A \beta \rightarrow \alpha \gamma \beta$	$ \mu \leq \varphi $	Autómatas lineales acotados
Tipo-2	Independiente de contexto	$A \rightarrow \gamma$	$ \mu = 1$	Autómatas de pila
Tipo-3	Regular	$A \rightarrow aB$ $A \rightarrow a$	$ \mu = 1$	Autómatas finitos, regulares



Gramática independiente del contexto (GIC)

Una "Gramática independiente del contexto" (GIC) es una cuádrupla:

$$G=(N,\Sigma,S,P)$$

¡RECUERDA! El orden de los argumentos cambia según la referencia.

Donde:

- N es una colección finita de no terminales.
- Σ es un alfabeto (Conjunto de terminales).
- S es un no terminal llamado **Símbolo inicial**.
- P es una colección finita de **reglas de sustitución** llamadas **producciones**.

El lenguaje generado por la GIC G se denota por $L(G)$ y se llama lenguaje independiente del contexto (LIC).



GIC de un lenguaje de programación

- Para un ejemplo simplificado nos limitaremos a los operadores de $+$ y $*$, los identificadores serán letras seguidas por cero o más letras y dígitos $(a|b)(a|b|0|1)^*$.
- Se necesitarán 2 variables ("E" e "ID").
 - E representa las expresiones (el lenguaje de las expresiones que se van a definir).
 - ID representa los identificadores.



GIC de un lenguaje de programación

$$G = (N, T, P, E)$$

$$N = \{E, ID\}$$

$$T = \{+, *, (,), a, b, 0, 1\}$$

$$P =$$

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

$$ID \rightarrow a \mid b \mid IDa \mid IDb \mid ID0 \mid ID1$$


Interpretación de las Producciones

(1) $E \rightarrow ID$

(2) $E \rightarrow E+E$

(3) $E \rightarrow E * E$

(4) $E \rightarrow (E)$

(5) $ID \rightarrow a$

(6) $ID \rightarrow b$

(7) $ID \rightarrow IDa$

(8) $ID \rightarrow IDb$

(9) $ID \rightarrow ID0$

(10) $ID \rightarrow ID1$



Interpretación de las Producciones

- (1) Establece que una expresión puede ser un único identificador .
- (2) Establece que una expresión puede ser igual a dos expresiones concatenadas mediante un signo +
- (3) Establece una expresión similar a la regla anterior pero para el signo *
- (4) Establece que si tomamos cualquier expresión y la encerramos entre paréntesis, el resultado también es una expresión.

(1) $E \rightarrow ID$
 (2) $E \rightarrow E + E$
 (3) $E \rightarrow E * E$
 (4) $E \rightarrow (E)$



Interpretación de las Producciones

De la regla (5) hasta la (10) se describen los identificadores ID. El caso básico lo describen las reglas (5) y (6) que establecen que a y b son identificadores. Las cuatro reglas restantes establecen que si tenemos cualquier símbolo a, b, 0 o 1, podemos escribir un identificador antes de ellos y el resultado será otro identificador.

(5) $ID \rightarrow a$
 (6) $ID \rightarrow b$
 (7) $ID \rightarrow IDa$
 (8) $ID \rightarrow IDb$
 (9) $ID \rightarrow ID0$
 (10) $ID \rightarrow ID1$



Derivaciones

- Una cadena $w \in T$ es **derivable** a partir de la gramática G si y solo si existe una secuencia de derivación iniciando en el símbolo inicial (no terminal "S") y terminando en la cadena w .
- La doble flecha \Rightarrow se interpreta como "deriva", "produce" o "genera".
- La notación $E \Rightarrow^* w$ indica que la cadena w se deriva a partir de S en 0 o mas etapas.



Derivaciones

- Cuando derivamos una cadena, los no terminales representan la parte de la cadena que todavía no se ha generado. Cuando la derivación se completa, todos los trozos no generados habrán sido sustituidos por cadenas de símbolos terminales (hay la posibilidad de que sean cadenas vacías).



Ejemplo de derivación

Aplicando las producciones de la GIC, inferir que la cadena **$a^*(a+b00)$** pertenece al lenguaje.

$G = (N, T, P, E)$

$N = \{E, ID\}$

$T = \{+, *, (,), a, b, 0, 1\}$

$P =$

$E \rightarrow ID \mid E+E \mid E^*E \mid (E)$

$ID \rightarrow a \mid b \mid IDa \mid IDb \mid ID0 \mid ID1$



Ejemplo de derivación

Respuesta:

$E \Rightarrow E^*E \Rightarrow ID^*E \Rightarrow a^*E \Rightarrow a^*(E) \Rightarrow a^*(E+E) \Rightarrow$
 $A^*(ID+E) \Rightarrow a^*(a+E) \Rightarrow a^*(a+ID) \Rightarrow a^*(a+ID0) \Rightarrow$
 $a^*(a+ID00) \Rightarrow a^*(a+b00)$



Derivaciones por la izquierda y por la derecha

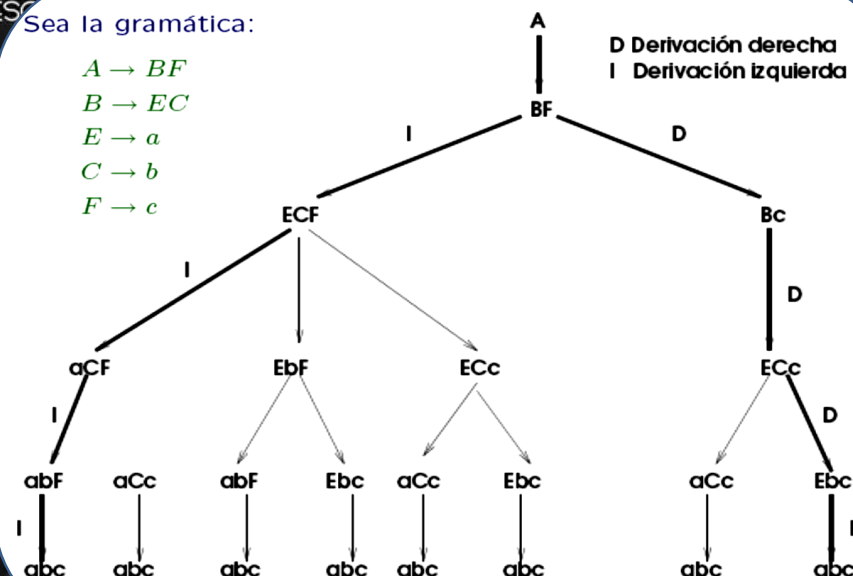
Con el fin de restringir el número de opciones disponibles en la derivación de una cadena, a menudo resulta útil seguir un patrón:

- **Derivación por la izquierda:** las reglas de reemplazo son aplicadas a la primera variable de izquierda a derecha.
- **Derivación por la derecha:** las reglas de reemplazo son aplicadas a la última variable de derecha a izquierda.



Sea la gramática:

$A \rightarrow BF$
 $B \rightarrow EC$
 $E \rightarrow a$
 $C \rightarrow b$
 $F \rightarrow c$





Ejemplo de derivación

Aplicando las producciones de la GIC, inferir que la cadena **a*(a+b00)** pertenece al lenguaje.

$G = (N, T, P, E)$

$N = \{E, ID\}$

$T = \{+, *, (,), a, b, 0, 1\}$

$P =$

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

$ID \rightarrow a \mid b \mid IDa \mid IDb \mid ID0 \mid ID1$



Ejemplo de derivación por la izquierda y por la derecha

- Derivación por la izquierda:

$E \Rightarrow E * E \Rightarrow ID * E \Rightarrow a * E \Rightarrow a * (E) \Rightarrow a * (E + E) \Rightarrow$
 $a * (ID + E) \Rightarrow a * (a + E) \Rightarrow a * (a + ID) \Rightarrow a * (a + ID0) \Rightarrow$
 $a * (a + ID00) \Rightarrow a * (a + b00)$

- Derivación por la derecha:

$E \Rightarrow E * E \Rightarrow E * (E) \Rightarrow E * (E + E) \Rightarrow E * (E + ID) \Rightarrow E * (E + ID0) \Rightarrow$
 $E * (E + ID00) \Rightarrow E * (E + b00) \Rightarrow E * (ID + b00) \Rightarrow$
 $E * (a + b00) \Rightarrow ID * (a + b00) \Rightarrow a * (a + b00)$



Árboles de Derivación

- A veces, es útil realizar un gráfico de la derivación, que indique de que manera ha contribuido cada no terminal a formar la cadena final de símbolos terminales. Tal gráfico tiene forma de árbol y se le conoce como árbol de derivación.



Árboles de Derivación

- Los árboles de derivación están estrechamente ligados a las derivaciones y a las inferencias recursivas.
- Una de sus mas importantes aplicaciones es en la ambigüedad en las gramáticas.



Árboles de Derivación

Sea $G = (V, T, P, S)$ una gramática. Los *árboles de derivación* para G son aquellos árboles que cumplen las condiciones siguientes:

1. Cada nodo interior está etiquetado con una variable de V .
2. Cada hoja está etiquetada bien con una variable, un símbolo terminal o ϵ . Sin embargo, si la hoja está etiquetada con ϵ , entonces tiene que ser el único hijo de su padre.
3. Si un nodo interior está etiquetado como A y sus hijos están etiquetados como X_1, X_2, \dots, X_k respectivamente, comenzando por la izquierda, entonces $A \rightarrow X_1 X_2 \dots X_k$ es una producción de P .

Observe que el único caso en que una de las X puede reemplazarse por ϵ es cuando es la etiqueta del único hijo y $A \rightarrow \epsilon$ es una producción de G .

Compiladores
Prof. Luis Enrique Hernández Olvera

41



Árboles de Derivación

Ejemplos:

Derivando $ID+E$ de E :



Derivación de

$P \Rightarrow^* 0110$



GIC para
palíndromos

- (1) $P \rightarrow \epsilon$
- (2) $P \rightarrow 0$
- (3) $P \rightarrow 1$
- (4) $P \rightarrow 0P0$
- (5) $P \rightarrow 1P1$

Compiladores
Prof. Luis Enrique Hernández Olvera

42



Árboles de Derivación

- Si nos fijamos en las hojas de cualquier árbol de derivación y las concatenamos empezando por la izquierda, de arriba hacia abajo, obtenemos una cadena denominada *resultado* del árbol, que siempre es una cadena que se deriva de la variable raíz.

Compiladores
Prof. Luis Enrique Hernández Olvera

43



Árboles de Derivación Ejemplos:

Derivando ID+E de E:



Derivación de

$P \Rightarrow^* 0110$



Compiladores
Prof. Luis Enrique Hernández Olvera

44



Árboles de Derivación

Los árboles de derivación cuyos resultados son cadenas pertenecientes al lenguaje de la gramática deben de cumplir con:

1. El resultado es una cadena terminal. Es decir, todas las hojas están etiquetadas con un símbolo terminal o con ϵ .
2. La raíz está etiquetada con el símbolo inicial.



Ejercicio

- De la cadena $a^*(a+b00)$ obtener:
- Su árbol de derivación por la izquierda que muestre que pertenece al lenguaje de nuestra gramática de expresiones.
- Su árbol de derivación por la derecha que muestre que pertenece al lenguaje de nuestra gramática de expresiones.



Gramáticas ambiguas

- Ciertas gramáticas permiten que una cadena terminal tenga más de un árbol de análisis. Esta situación hace que esa gramática sea inadecuada para un lenguaje de programación, ya que el compilador no puede decidir la estructura sintáctica de determinados programas fuentes y, por tanto, no podría deducir con seguridad cuál será el código ejecutable apropiado correspondiente al programa.



Gramáticas ambiguas

- Usando el ejemplo de gramática de expresiones observemos que las producciones
 - $E \rightarrow E * E \mid E + E$
- nos permite generar estas expresiones en cualquier orden que elijamos.



Gramáticas ambiguas Ejemplo

- consideremos la forma sentencial $E+E * E$.
Existen dos derivaciones de E:

$$(2) E \rightarrow E + E$$

$$(3) E \rightarrow E * E$$



Gramáticas ambiguas Ejemplo

- consideremos la cadena $a+a*a$. Existen dos derivaciones de E:

$$(2) E \rightarrow E + E$$

$$(3) E \rightarrow E * E$$



Eliminación de la ambigüedad de las gramáticas

- no existe un algoritmo que nos diga si una GIC es ambigua.
- Para los tipos de construcciones que aparecen en los lenguajes de programación comunes, existen técnicas bien conocidas que permiten eliminar la ambigüedad.



Eliminación de la ambigüedad de las gramáticas

- Analizando el ejemplo anterior vemos claramente que la precedencia de operadores no se respeta.
- La solución al problema de forzar la precedencia se resuelve introduciendo varias variables distintas, cada una de las cuales representa aquellas expresiones que comparten el mismo nivel de "fuerza de acoplamiento".



Eliminación de la ambigüedad de las gramáticas

- La gramática resultante para la ambigüedad en nuestra GIC queda de la siguiente forma:

$$E \rightarrow ID \mid E+E \mid E * E$$


$$\begin{aligned} E &\rightarrow E+M \mid M \\ M &\rightarrow M * M \mid ID \end{aligned}$$


GIC sin ambigüedad

$$G = (N, T, P, E)$$

$$N = \{E, ID\}$$

$$T = \{+, *, (,), a, b, 0, 1\}$$

$$P =$$

$$E \rightarrow E+M \mid M$$

$$M \rightarrow M * P \mid P$$

$$P \rightarrow (E) \mid ID$$

$$ID \rightarrow a \mid b \mid IDa \mid IDb \mid ID0 \mid ID1$$



Ejercicios

Encontrar una GIC que genere los lenguajes:

- a) $L = \{a^i b^i : i \geq 0\}$ sobre $\Sigma = \{a, b\}$
- b) Lenguaje de los palíndromos sobre $\Sigma = \{a, b\}$
- c) Todas las cadenas sobre $\Sigma = \{a, b\}$ que tienen un número par de símbolos.
- d) $(ab \cup ba)^*$ sobre $\Sigma = \{a, b\}$



Ejercicios (Soluciones posibles)

$$G = (N, \Sigma, S, P)$$

a)
 $\Sigma = \{a, b\}$
 $N = \{S\}$
 $P =$
 $S \rightarrow aSb$
 $S \rightarrow \lambda$

b)
 $\Sigma = \{a, b\}$
 $N = \{S\}$
 $P =$
 $S \rightarrow aSa \mid$
 $bSb \mid a \mid$
 $b \mid \lambda$

c)
 $\Sigma = \{a, b\}$
 $N = \{S\}$
 $P =$
 $S \rightarrow aSa \mid$
 $bSb \mid aSb \mid$
 $bSa \mid \lambda$

d)
 $\Sigma = \{a, b\}$
 $N = \{S\}$
 $P =$
 $S \rightarrow abS$
 $S \rightarrow baS$
 $S \rightarrow \lambda$