

# Procesamiento del Lenguaje Natural

# Procesamiento del Lenguaje Natural

Los agentes envían señales a otros agentes para...

- Informar
- Pedir ayuda
- Compartir conocimiento
- Preguntar
- Ordenar
- Comprometerse
- Etc.

# Procesamiento del Lenguaje Natural

- La comunicación requiere ciertas convenciones que forman el lenguaje
- El lenguaje está formado por signos con significados acordados
- No es una habilidad propia de los humanos
- Los humanos parecen usar la gramática para producir infinitos mensajes estructurados

# Procesamiento del Lenguaje Natural

## Lenguajes formales

- Son un conjunto de cadenas
- Cada cadena es una concatenación de símbolos terminales (palabras)
- Ejemplos: lógica de predicados, Python
- **Tienen definiciones estrictas**
- Asocian un significado o semántica a cada cadena válida
- El significado de una cadena sólo está influido por su forma

# Procesamiento del Lenguaje Natural

## Lenguajes naturales

- Generado espontáneamente en un grupo de hablantes
- **No tiene definiciones estrictas**
- El significado específico y contextual de sus componentes intervienen en la validez o no de la frase
- El lenguaje natural no puede ser fácilmente caracterizado porque:
  - No existen reglas gramaticales claras (o no se respetan)
  - Hay ambigüedad
  - La relación entre símbolo y objetos no está formalmente definida
- Pragmática: el significado real de la cadena cuando es dicha en una situación determinada.
- La teoría de lenguajes formales es útil para el estudio de los lenguajes naturales

# Procesamiento del Lenguaje Natural

## Gramáticas

- Una gramática es un conjunto finito de reglas que especifican un lenguaje
- Define la sintaxis de las cadenas válidas
- Las cadenas están compuestas de subcadenas llamadas frases, las cuales pueden pertenecer a distintas categorías.
- Se usan tanto para interpretar el texto como para generarlo

# Procesamiento del Lenguaje Natural

## Categorías gramaticales

- Frases nominales (FN): Hacen referencia a objetos del mundo
- Frases verbales (FV): Indican Acciones
- A frases de cierta categoría, como las nominales y las verbales se puede combinar para crear frases del tipo "sentencia" (S).
- A las categorías (FV, FN, S, etc.) se las llama "símbolos no terminales"
- Las gramáticas definen los no terminales por medio de reglas de reescritura.
- Notación de la forma Backus-Naur (BNF):  $S \rightarrow FN FV$  (una sentencia puede consistir en una FN seguida por cualquier FV)

# Procesamiento del Lenguaje Natural

Ejemplo del libro



# Procesamiento del Lenguaje Natural

## Ejemplo del libro

<i>Sustantivo</i>	→	<b>hedor</b>   <b>brisa</b>   <b>resplandor</b>   <b>nada</b>   <b>agente</b>   <i>wumpus</i>   <b>foso</b>   <b>oro</b>   <b>este</b>   ...
<i>Verbo</i>	→	<b>es</b>   <b>ver</b>   <b>oler</b>   <b>disparar</b>   <b>sentir</b>   <b>heder</b>   <b>ir</b>   <b>tomar</b>   <b>llevar</b>   <b>matar</b>   <b>girar</b>   ...
<i>Adjetivo</i>	→	<b>derecho</b>   <b>izquierdo</b>   <b>oriental</b>   <b>muerto</b>   <b>posterior</b>   <b>hediondo</b>   ...
<i>Adverbio</i>	→	<b>aquí</b>   <b>allí</b>   <b>cerca</b>   <b>adelante</b>   <b>correctamente</b>   <b>a la izquierda</b>   <b>al este</b>   <b>al sur</b>   <b>atrás</b>   ...
<i>Pronombre</i>	→	<b>mí</b>   <b>tú</b>   <b>yo</b>   <b>ello</b>   ...
<i>Nombre</i>	→	<b>Juan</b>   <b>María</b>   <b>Alicante</b>   <b>Aristóteles</b>   ...
<i>Artículo</i>	→	<b>el</b>   <b>un</b>   ...
<i>Preposición</i>	→	<b>a</b>   <b>en</b>   <b>sobre</b>   ...
<i>Conjunción</i>	→	<b>y</b>   <b>o</b>   <b>pero</b>   ...
<i>Dígito</i>	→	<b>0</b>   <b>1</b>   <b>2</b>   <b>3</b>   <b>4</b>   <b>5</b>   <b>6</b>   <b>7</b>   <b>8</b>   <b>9</b>

**Figura 22.3** El léxico para  $\mathcal{E}_0$ .

# Procesamiento del Lenguaje Natural

$S$	$\rightarrow$	$FN\ FV$ $S\ Conjunción\ S$	Yo + siento una brisa Yo siento una brisa + y + yo huelo a <i>wumpus</i>
$FN$	$\rightarrow$	$Pronombre$ $Nombre$ $Sustantivo$ $Artículo\ del\ sustantivo$ $Dígito\ Dígito$ $FN\ FP$ $FN\ OraciónRel$	Yo Juan foso el + <i>wumpus</i> 3 4 el <i>wumpus</i> + al este el <i>wumpus</i> + que es hediondo
$FV$	$\rightarrow$	$Verbo$ $FV\ FN$ $FV\ Adjetivo$ $FV\ FP$ $FV\ Advverbio$	huele mal siento + una brisa es hediondo gira + al este ve + adelante
$FP$	$\rightarrow$	$Preposición\ FN$	a + el este
$OraciónRel$	$\rightarrow$	<b>que</b> $FV$	que + es hediondo

**Figura 22.4** La gramática para  $\mathcal{E}_0$ , con frases de ejemplo para cada regla.

# Procesamiento del Lenguaje Natural

<i>Sustantivo</i>	→	hedor   brisa   resplandor   nada   agente   wumpus   foso   oro   este   ...
<i>Verbo</i>	→	es   ver   oler   disparar   sentir   heder   ir   tomar   llevar   matar   girar   ...
<i>Adjetivo</i>	→	derecho   izquierdo   oriental   muerto   posterior   hediondo   ...
<i>Adverbio</i>	→	aquí   allí   cerca   adelante   correctamente   a la izquierda   al este   al sur   atrás   ...
<i>Pronombre</i>	→	mí   tú   yo   ello   ...
<i>Nombre</i>	→	Juan   María   Alicante   Aristóteles   ...
<i>Artículo</i>	→	el   un   ...
<i>Preposición</i>	→	a   en   sobre   ...
<i>Conjunción</i>	→	y   o   pero   ...
<i>Dígito</i>	→	0   1   2   3   4   5   6   7   8   9

Figura 22.3 El léxico para  $\mathcal{E}_0$ .

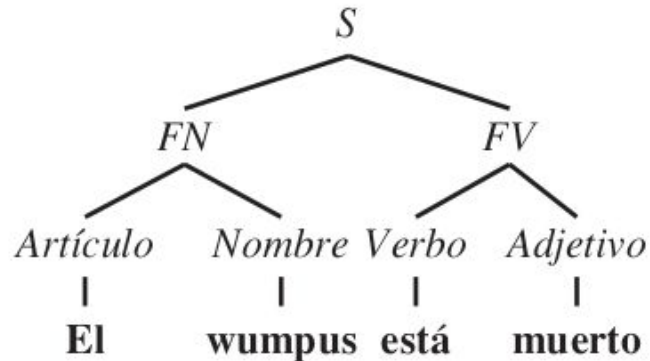
<i>S</i>	→	<i>FN FV</i>	Yo + siento una brisa
		<i>S Conjunción S</i>	Yo siento una brisa + y + yo huelo a wumpus
<i>FN</i>	→	<i>Pronombre</i>	Yo
		<i>Nombre</i>	Juan
		<i>Sustantivo</i>	foso
		<i>Artículo del sustantivo</i>	el + wumpus
		<i>Dígito Dígito</i>	3 4
		<i>FN FP</i>	el wumpus + al este
		<i>FN OracionRel</i>	el wumpus + que es hediondo
<i>FV</i>	→	<i>Verbo</i>	huele mal
		<i>FV FN</i>	siento + una brisa
		<i>FV Adjetivo</i>	es hediondo
		<i>FV FP</i>	gira + al este
		<i>FV Adverbio</i>	ve + adelante
<i>FP</i>	→	<i>Preposición FN</i>	a + el este
<i>OracionRel</i>	→	<i>que FV</i>	que + es hediondo

Figura 22.4 La gramática para  $\mathcal{E}_0$ , con frases de ejemplo para cada regla.

# Procesamiento del Lenguaje Natural

## Análisis sintáctico (parsing)

- Descubrimiento de la estructura de las frases de acuerdo con las reglas de una gramática
- Búsqueda de un árbol de análisis válido cuyas hojas son las palabras de la cadena



# Procesamiento del Lenguaje Natural

## Gramáticas aumentadas

- Hasta ahora hemos visto gramáticas libres de contexto
- Las desventajas se hacen evidentes en la generación
  - No todas las frases sustantivas (por ejemplo) pueden aparecer en cualquier lugar. "Me comí una manzana" vs "Me comí una consulta de SQL"
  - Tampoco da lo mismo usar cualquier artículo. "La casa" vs "Los casa"

# Procesamiento del Lenguaje Natural

## Gramáticas aumentadas

- Hay más reglas que podríamos agregar. Por ejemplo:
  - Una gramática simple puede generar “Yo huelo un hedor”, pero también genera “Mi huelo un hedor”
  - “Yo” se emplea en el caso subjetivo y “mi” en el caso objetivo

$$\begin{aligned} S &\rightarrow FN_s FV \mid \dots \\ FN_s &\rightarrow Pronombre_s \mid Nombre \mid Sustantivo \mid \dots \\ FN_o &\rightarrow Pronombre_o \mid Nombre \mid Sustantivo \mid \dots \\ FV &\rightarrow FV FN_o \mid \dots \\ FP &\rightarrow Preposición FN_o \\ Pronombre_s &\rightarrow \mathbf{Yo} \mid \mathbf{tu} \mid \mathbf{el} \mid \mathbf{ella} \mid \mathbf{ello} \mid \dots \\ Pronombre_o &\rightarrow \mathbf{mi} \mid \mathbf{tu} \mid \mathbf{su} \mid \dots \end{aligned}$$

# Procesamiento del Lenguaje Natural

## Gramáticas aumentadas

- Con los pronombres por ejemplo, además sabemos que "Yo" es la 1ra persona del singular, "Ustedes" es la segunda del plural, etc.
- Una categoría como Pronombre que ha sido aumentada con características como el género y número es llamada subcategoría.
- Estas gramáticas pueden representar ese tipo de conocimiento para hacer una distinción más específica de qué tan probable es una oración.

# Procesamiento del Lenguaje Natural

## Gramáticas aumentadas

- Aumentar la cantidad de reglas para incluir a cada subcategoría haría crecer exponencialmente el tamaño de la gramática.
- En vez de introducir nuevas reglas. Las reglas aumentadas permiten parámetros para las categorías no terminales.



# Procesamiento del Lenguaje Natural

## Gramáticas aumentadas

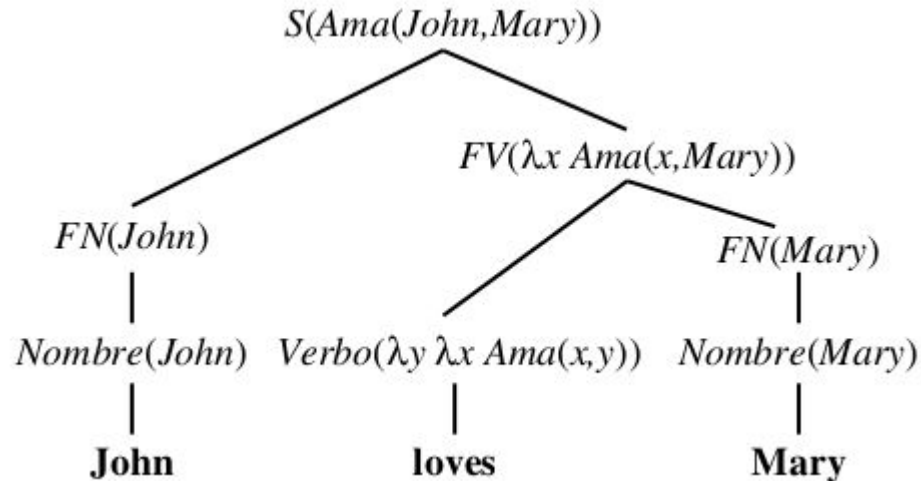
- Aumentar la expresividad de las gramáticas formales, permitiendo una mayor flexibilidad en la generación de oraciones.
- En vez de incrementar el número de parámetros de la gramática, se permiten

$$\begin{aligned} S &\rightarrow FN_s FV \mid \dots \\ FN_s &\rightarrow Pronombre_s \mid Nombre \mid Sustantivo \mid \dots \\ FN_o &\rightarrow Pronombre_o \mid Nombre \mid Sustantivo \mid \dots \\ FV &\rightarrow FV FN_o \mid \dots \\ FP &\rightarrow Preposición FN_o \end{aligned}$$
$$\begin{aligned} Pronombre_s &\rightarrow \text{Yo} \mid \text{tu} \mid \text{el} \mid \text{ella} \mid \text{ello} \mid \dots \\ Pronombre_o &\rightarrow \text{mi} \mid \text{tu} \mid \text{su} \mid \dots \end{aligned}$$
$$\begin{aligned} S &\rightarrow FN(\text{Subjetivo}) FV \mid \dots \\ FN(\text{caso}) &\rightarrow Pronombre(\text{caso}) \mid Nombre \mid Sustantivo \mid \dots \\ FV &\rightarrow FV FN(\text{Objetivo}) \mid \dots \\ FP &\rightarrow Preposición FN(\text{Objetivo}) \end{aligned}$$
$$\begin{aligned} Pronombre(\text{Subjetivo}) &\rightarrow \text{Yo} \mid \text{tu} \mid \text{el} \mid \text{ella} \mid \text{ello} \mid \dots \\ Pronombre(\text{Objetivo}) &\rightarrow \text{mi} \mid \text{tu} \mid \text{su} \mid \dots \end{aligned}$$

# Procesamiento del Lenguaje Natural

## Interpretación semántica

- Se incorporan en la gramática equivalencias semánticas



# Procesamiento del Lenguaje Natural

## Modelos probabilísticos de lenguaje

- En lenguaje natural no se puede definir (como verdadero/falso) si una cadena es gramatical, pero podemos calcular qué tan probable es que una lo sea
- Un modelo de lenguaje es una distribución que describe la probabilidad de cualquier cadena de texto.
- ¿Probabilidad de qué? De ser gramaticalmente correcta, de ser una respuesta adecuada, de representar una situación dada, una traducción, etc..

# Procesamiento del Lenguaje Natural

## Modelos probabilísticos de lenguaje

- Bag-of-words
- Modelos n-gram
- Gramáticas

# Procesamiento del Lenguaje Natural

## Bag-of-words

- Permite clasificar textos
- Es un modelo generativo

$$\mathbf{P}(Class | w_{1:N}) = \alpha \mathbf{P}(Class) \prod_j \mathbf{P}(w_j | Class)$$

# Procesamiento del Lenguaje Natural

## N-gram word models

- Una de las limitaciones de bag-of-words es que hay palabras que pueden pertenecer a más de un dominio, por ejemplo, “unicornio”
- Si se toman conjuntos de  $n$  palabras adyacentes es más fácil determinar la clase o predecir la próxima palabra

$$P(w_j | w_{1:j-1}) = P(w_j | w_{j-n+1:j-1})$$

$$P(w_{1:N}) = \prod_{j=1}^N P(w_j | w_{j-n+1:j-1})$$

# Procesamiento del Lenguaje Natural

## Gramáticas

$S$	$\rightarrow NP VP$	[0.90]	I + feel a breeze
	$S Conj S$	[0.10]	I feel a breeze + and + It stinks
$NP$	$\rightarrow Pronoun$	[0.25]	I
	$Name$	[0.10]	Ali
	$Noun$	[0.10]	pits
	$Article Noun$	[0.25]	the + wumpus
	$Article Adjs Noun$	[0.05]	the + smelly dead + wumpus
	$Digit Digit$	[0.05]	3 4
	$NP PP$	[0.10]	the wumpus + in 1 3
	$NP RelClause$	[0.05]	the wumpus + that is smelly
	$NP Conj NP$	[0.05]	the wumpus + and + I
$VP$	$\rightarrow Verb$	[0.40]	stinks
	$VP NP$	[0.35]	feel + a breeze

# Procesamiento del Lenguaje Natural

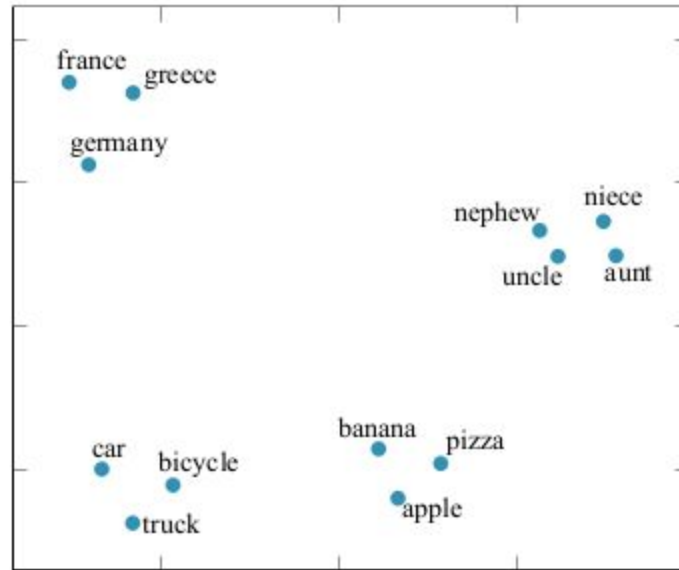
## Deep learning

- Dada la cantidad de texto disponible, tenía sentido considerar enfoques basados en aprendizaje automático orientado a datos.
- Algunas alternativas:
  - Word embeddings
  - Redes recurrentes
  - LSTM
  - Transformers



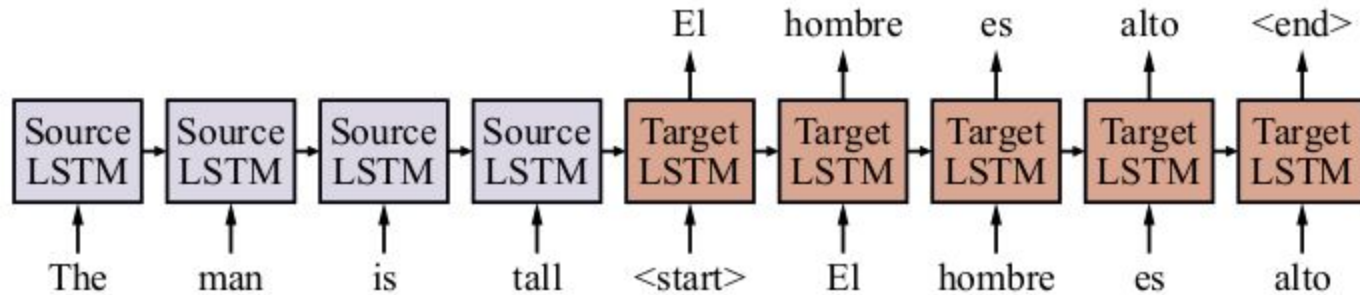
# Procesamiento del Lenguaje Natural

## Word embeddings



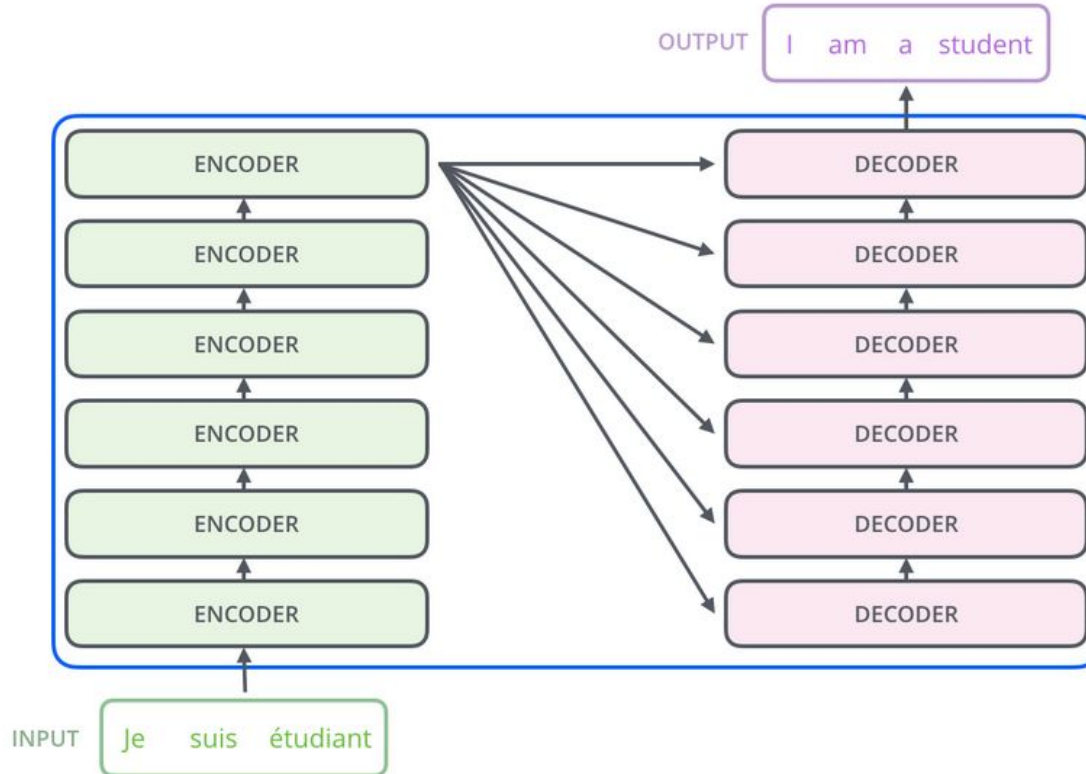
# Procesamiento del Lenguaje Natural

## LSTM - Sequence-to-sequence



# Procesamiento del Lenguaje Natural

## Transformers



Fin