

## Introducción al Procesamiento de Lenguaje Natural SOLUCIÓN - Diciembre de 2006

### Consideraciones generales

- i) La prueba es sin material escrito.
- ii) Escriba nombre y C.I. en todas las hojas.
- iii) Numere todas las hojas.
- iv) En la primera hoja, indique el total de hojas.
- v) Comience cada ejercicio en una hoja nueva.
- vi) Utilice las hojas de un solo lado.
- vii) Entregue los ejercicios en orden.

### Ejercicio 1 [6 puntos]

a) Describa brevemente las dos hipótesis principales sobre la forma en que se realiza el procesamiento morfológico humano.

*Existen dos hipótesis principales para el procesamiento morfológico (como se analiza la estructura de las palabras y como se generan las mismas) en humanos. La primera supone que se cuenta con una lista de todas las palabras posibles (sin considerar su estructura interna). La segunda, que las palabras se construyen a partir de una lista de morfemas (raíces y afijos) que se combinan de acuerdo a ciertas reglas de ordenamiento (morfológicas).*

b) Mencione dos ventajas de los transductores como herramienta para el desarrollo de analizadores morfológicos.

*Por su propia definición, los transductores permite realizar el análisis léxico y la generación como procesos simétricos, por lo que resuelto el problema del análisis se resuelve el de la generación, y viceversa. Por sus propiedades, permiten realizar un procesamiento eficiente de la entrada (en aquellos casos en se pueden secuencializar, es  $O(n)$ , siendo  $n$  el tamaño de la entrada). Los transductores pueden generarse fácilmente a partir de un álgebra de expresiones regulares que denota las relaciones regulares que ellos computan.*

### Ejercicio 2 [16 puntos]

Considere los **nombres** en el idioma ficticio Itio (ligeramente basado en el Esperanto), que se construyen según las siguientes reglas:

1. Comienzan con una raíz. Para este ejemplo consideraremos las siguientes: *hund, kat, bird, elephant* (en español: *perro, gato, pájaro, elefante*, respectivamente).
2. Una raíz por sí misma no es una palabra válida.
3. Todos los nombres tienen un sufijo *o*. Por lo tanto *elefanto* y *hundo* son palabras válidas.
4. Una palabra como *hundo* no tiene género marcado. Para marcarla como femenina, se agrega el sufijo *in* entre la raíz y el sufijo *o*. Por ejemplo *hundino* (perra) es una palabra válida.
5. Para marcar una palabra como diminutivo, se agrega un sufijo *et* entre la raíz y el sufijo *o*, por ejemplo *elefanteto* es una palabra válida. Para marcar un aumentativo, se agrega el sufijo *eg*.
6. Los sufijos femenino, diminutivo y aumentativo pueden co-ocurrir una cantidad arbitraria de veces.
7. Si una raíz termina en *t*, al agregarle el sufijo *in*, ésta se convierte en *v*, por cuestiones de pronunciación. Por lo tanto *gata* se escribe *kavino*.

a) Describa brevemente los tres componentes principales de un analizador morfológico. Describa cada componente para un hipotético analizador morfológico para los nombres en Itio.

Los tres componentes principales son: **Lexicón:** conjunto de morfemas del lenguaje (raíces y afijos). **Reglas morfológicas:** reglas para el ordenamiento de los morfemas. **Reglas ortográficas:** modificaciones en la forma de las palabras al combinarse los morfemas. En Itio el lexicón está dado por el conjunto formado por las raíces y los sufijos, las morfológicas son las reglas que dicen que los sufijos van después de las raíces, y la única regla ortográfica es la del punto 7.

b) Construya una expresión regular utilizando el álgebra de Xerox, que relacione un nombre en Itio con su análisis morfológico. Utilice las siguientes marcas léxicas para la salida: +Noun (Raíz), +Aug(Aumentativo), +Dim(Diminutivo),+Fem(Femenino). Por ejemplo, a la forma léxica *elefant+Noun+Fem* le corresponderá la palabra *elefanvino*.

La expresión regular itio denota, en el álgebra de expresiones regulares de Xerox, la relación regular que mapea palabras y estructuras de los nombres en Itio (primero se definen las raíces y los afijos), luego se los combina de acuerdo a las reglas morfológicas, para finalmente componer la relación con el reemplazo que denota los cambios ortográficos.:

```
define root [{hund}|{kat}|{bird}|{elefant}] %+Noun:0;
define sufijoo 0:o;
define diminutivo %+Dim:{et};
define aumentativo %+Aug:{eg};
define femenino %+Fem:{in};
define ortografia t -> v || _ {in};
define itio [root (femenino|diminutivo|aumentativo)* sufijoo] .o. ortografia;
```

### Ejercicio 3 [36 puntos]

a) Dé una definición de constituyente sintáctico.

Un constituyente sintáctico se puede definir como una palabra o secuencia de palabras que funciona en conjunto como una unidad dentro de la estructura jerárquica de una oración.

b) Sea el conjunto de oraciones siguiente:

- Festejemos
- Festejemos en el tren
- El ángel exterminador es una película de Buñuel
- Los suinos gruñen
- El tren de Roma a Viena partió tardíamente
- El sobre está sobre la mesa de luz

i) Modele este micro universo textual utilizando GLC.

Gramática		Léxico	
O	→ GV   GN GV	V	→ festejemos   es   gruñen   partió   está
GV	→ V   V GN   V GP   V Adv	Adj	→ exterminador
GN	→ N   Det N   Det N GP   Det N Adj   Det N GP GP   NomProp	Prep	→ de   a   sobre   en
GP	→ Prep GN	Det	→ el   una   los
		N	→ ángel   película   suinos   tren   sobre   mesa   luz
		NomProp	→ Buñuel   Roma   Viena
		Adv	→ tardíamente

Otras gramáticas son viables. Por ejemplo, "El ángel exterminador" puede manejarse como nombre propio. Se podría también manejar "mesa de luz" como un sustantivo único, aunque es más discutible.

ii) Indique al menos dos problemas que surgen al utilizar GLC para modelar lenguaje. Ejemplifique con la GLC construida en la parte anterior.

*Es costoso controlar concordancia. Con la gramática anterior puede generarse, por ejemplo, la oración "el película gruñen". También hay problemas con la subcategorización. Se podría generar, por ejemplo "los suinos gruñen Buñuel".*

c) Sea la GLC siguiente:

Gramática		Léxico	
O	→ GV   GN GV	N	→ vino   dueño   casa
GV	→ V GN   V GN GP	Det	→ el   la   los
GN	→ Det N   Det N GP	V	→ vino   dice
GP	→ Prep GN	Prep	→ de   a

i) Aplique el algoritmo de Earley a la siguiente entrada: *vino el dueño de la casa*

chart[0]			
1	$\gamma \rightarrow \bullet O$	[0,0]	Predecir
2	$O \rightarrow \bullet GV$	[0,0]	Predecir
3	$O \rightarrow \bullet GN GV$	[0,0]	Predecir
4	$GV \rightarrow \bullet V GN$	[0,0]	Predecir
5	$GV \rightarrow \bullet V GN GP$	[0,0]	Predecir
6	$GN \rightarrow \bullet Det N$	[0,0]	Predecir
7	$GN \rightarrow \bullet Det N GP$	[0,0]	Predecir
chart[1]			
8	$V \rightarrow vino \bullet$	[0,1]	Buscar
9	$GV \rightarrow V \bullet GN$	[0,1]	Completar (4,8)
10	$GV \rightarrow V \bullet GN GP$	[0,1]	Completar (5,8)
11	$GN \rightarrow \bullet Det N$	[1,1]	Predecir

12	GN → • Det N GP	[1,1]	Predecir
<b>chart[2]</b>			
13	Det → el •	[1,2]	Buscar
14	GN → Det • N	[1,2]	Completar (11,13)
15	GN → Det • N GP	[1,2]	Completar (12,13)
<b>chart[3]</b>			
16	N → dueño •	[2,3]	Buscar
17	GN → Det N •	[1,3]	Completar (14,16)
18	GN → Det N • GP	[1,3]	Completar (15,16)
19	GP → • Prep GN	[3,3]	Predecir
20	GV → V GN •	[0,3]	Completar (9,17)
21	GV → V GN • GP	[0,3]	Completar (10,17)
22	O → GV •	[0,3]	Completar (2,20)
<b>chart[4]</b>			
23	Prep → de •	[3,4]	Buscar
24	GP → Prep • GN	[3,4]	Completar (19,23)
25	GN → • Det N	[4,4]	Predecir
26	GN → • Det N GP	[4,4]	Predecir
<b>chart[5]</b>			
27	Det → la •	[4,5]	Predecir
28	GN → Det • N	[4,5]	Completar (25,27)
29	GN → Det • N GP	[4,5]	Completar (26,27)
<b>chart[6]</b>			
30	N → casa •	[5,6]	Predecir
31	GN → Det N •	[4,6]	Completar (28,30)
32	GN → Det N • GP	[4,6]	Completar (29,30)
33	GP → Prep GN •	[3,6]	Completar (24,31)
34	GN → Det N GP •	[1,6]	Completar (18,33)
35	GV → V GN GP •	[0,6]	Completar (21,33)
36	GV → V GN •	[0,6]	Completar (9,34)
37	GV → V GN • GP	[0,6]	Completar (10,34)
38	O → GV •	[0,6]	Completar (2,35)
39	O → GV •	[0,6]	Completar (2,36)

ii) Escriba el árbol o los árboles de análisis sintáctico correspondiente(s) al punto anterior.

Hay dos árboles posibles para la entrada, que corresponden a las filas 38 y 39 del chart anterior. Sus análisis son los siguientes:

[o [GV [v vino] [GN [Det el] [N dueño]]] [GP [Prep de] [GN [Det la] [N casa] ] ] ] ]

[o [GV [v vino] [GN [Det el] [N dueño] [GP [Prep de] [GN [Det la] [N casa] ] ] ] ] ]

#### Ejercicio 4 [6 puntos]

- Indique 2 problemas que aparecen al representar en Lógica de Primer Orden enunciados en lenguaje natural.
- Describa brevemente la aplicación de los métodos de aprendizaje automático supervisado y no supervisado para WSD. Mencione ventajas y desventajas de cada método.

*Comentar brevemente los métodos mencionados durante el curso. Las ventajas y desventajas se encuentran resumidas en la sección 19.1 del libro del curso.*

#### Ejercicio 5 [30 puntos]

Sea la siguiente gramática aumentada con anotaciones semánticas:

o → gn gv	[o.sem = gv.sem (gn.sem)]
gn → det nominal	[gn.sem = < det.sem x nominal.sem(x) > ]
gn → npropio	[gn.sem = npropio.sem]
nominal → n	[nominal.sem = $\lambda x$ isa(x,n.sem)]
nominal → nominal adj	
nominal → nominal pp	
gv → v	[gv.sem = v.sem]
gv → v gn	[gv.sem = v.sem (gn.sem)]
pp → prep gn	
v → lee	[v.sem = $\lambda x \lambda y \exists e$ leer (e) $\wedge$ lee(e,y) $\wedge$ leído(e,x)]
prep → de	[prep.sem = $\lambda x \lambda y$ de(y,x)]
prep → en	[prep.sem = $\lambda x \lambda y$ en(y,x)]

- Complete la gramática anterior con las 3 anotaciones semánticas que faltan. Estas anotaciones deben cubrir los casos en los que los adjetivos tienen una semántica intersectiva y las preposiciones una interpretación relacional (tal como aparece en las reglas para “de” y “en”).

Se agregan las partes subrayadas en la gramática:

o → gn gv	[o.sem = gv.sem (gn.sem)]
gn → det nominal	[gn.sem = < det.sem x nominal.sem(x) > ]
gn → npropio	[gn.sem = npropio.sem]
nominal → n	[nominal.sem = $\lambda x$ isa(x,n.sem)]
nominal → nominal adj	<u>[nominal.sem = <math>\lambda x</math> nominal.sem (x) <math>\wedge</math> isa (x, adj.sem)]</u>
nominal → nominal pp	<u>[nominal.sem = <math>\lambda x</math> nominal.sem (x) <math>\wedge</math> pp.sem(x)]</u>
gv → v	[gv.sem = v.sem]
gv → v gn	[gv.sem = v.sem (gn.sem)]
pp → prep gn	<u>[pp.sem = prep.sem(gn.sem)]</u>
v → lee	[v.sem = $\lambda x \lambda y \exists e$ leer (e) $\wedge$ lee(e,y) $\wedge$ leído(e,x)]
prep → de	[ prep.sem = $\lambda x \lambda y$ de(y,x)]
prep → en	[ prep.sem = $\lambda x \lambda y$ en(y,x)]
<u>npropio → Juan</u>	<u>[Juan]</u>
<u>npropio → Sartre</u>	<u>[Sartre]</u>

$n \rightarrow \text{autor} \quad [\text{autor}]$   
 $n \rightarrow \text{libro} \quad [\text{libro}]$   
 $\text{adj} \rightarrow \text{interesante} \quad [\text{interesante}]$   
 $\text{adj} \rightarrow \text{joven} \quad [\text{joven}]$

b) Utilizando las reglas anteriores (completadas con las entradas léxicas necesarias) realice las derivaciones para las representaciones semánticas de :

i. *Juan lee un libro interesante de Sartre.*

$n \rightarrow \text{libro} \quad [\text{libro}]$   
*libro*  
 $\text{nominal} \rightarrow n \quad [\lambda x \text{ isa}(x, \text{libro})]$   
*libro*  
 $\text{nominal} \rightarrow \text{nominal adj} \quad [\lambda x \text{ isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante})]$   
*libro interesante*  
 $\text{pp} \rightarrow \text{prep gn} \quad [\lambda y \text{ de}(y, \text{Sartre})]$   
*de Sartre*  
 $\text{nominal} \rightarrow \text{nominal pp} \quad [\lambda x \text{ isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante}) \wedge \text{de}(x, \text{Sartre})]$   
*libro interesante de Sartre*  
 $\text{gn} \rightarrow \text{det nominal} \quad [< \exists x \text{ isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante}) \wedge \text{de}(x, \text{Sartre}) >]$   
*un libro interesante de Sartre*  
 $\text{gv} \rightarrow v \text{ gn} \quad [\lambda y \exists e \text{ leer}(e) \wedge \text{lee}(e, y) \wedge \text{leído}(e, < \exists x \text{ isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante}) \wedge \text{de}(x, \text{Sartre}) >)]$   
*lee un libro interesante de Sartre*  
 $o \rightarrow \text{gn gv} \quad [\exists e \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \text{leído}(e, < \exists x \text{ isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante}) \wedge \text{de}(x, \text{Sartre}) >)]$   
*Juan lee un libro interesante de Sartre*

ii. *Juan lee un libro de un autor joven.*

$n \rightarrow \text{libro} \quad [\text{autor}]$   
*autor*  
 $\text{nominal} \rightarrow n \quad [\lambda x \text{ isa}(x, \text{autor})]$   
*autor*  
 $\text{nominal} \rightarrow \text{nominal adj} \quad [\lambda x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven})]$   
*autor joven*  
 $\text{gn} \rightarrow \text{det nominal} \quad [< \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) >]$   
*un autor joven*  
 $\text{pp} \rightarrow \text{prep gn} \quad [\lambda y \text{ de}(y, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) >)]$   
*de un autor joven*  
 $\text{nominal} \rightarrow \text{nominal pp} \quad [\lambda x \text{ isa}(x, \text{libro}) \wedge \text{de}(x, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) >)]$   
*libro de un autor joven*  
 $\text{gn} \rightarrow \text{det nominal} \quad [< \exists y \text{ isa}(y, \text{libro}) \wedge \text{de}(y, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) > > >]$   
*un libro de un autor joven*  
 $\text{gv} \rightarrow v \text{ gn} \quad [\lambda z \exists e \text{ leer}(e) \wedge \text{lee}(e, z) \wedge \text{leído}(e, < \exists y \text{ isa}(y, \text{libro}) \wedge \text{de}(y, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) > > >)]$   
*lee un libro de un autor joven*  
 $o \rightarrow \text{gn gv} \quad [\exists e \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \text{leído}(e, < \exists y \text{ isa}(y, \text{libro}) \wedge \text{de}(y, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) > > >)]$

$$\text{leído}(e, < \exists y \text{ isa}(y, \text{libro}) \wedge \\ \text{de}(y, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) > > >)]$$

Juan lee un libro de un autor joven

c) Expresé en lógica de primer orden las expresiones en *quasi logical form* obtenidas en el paso anterior.

i.

$$\begin{aligned} & \exists e \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \\ & \quad \text{leído}(e, < \exists x \text{ isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante}) \wedge \text{de}(x, \text{Sartre}) >) \\ \Rightarrow & \exists e \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \\ & \quad \exists x \text{ leído}(e, x) \wedge \text{isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante}) \wedge \text{de}(x, \text{Sartre}) \\ \Rightarrow & \exists e \exists x \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \\ & \quad \text{leído}(e, x) \wedge \text{isa}(x, \text{libro}) \wedge \text{isa}(x, \text{interesante}) \wedge \text{de}(x, \text{Sartre}) \end{aligned}$$

ii.

$$\begin{aligned} & \exists e \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \text{leído}(e, < \exists y \text{ isa}(y, \text{libro}) \wedge \\ & \quad \text{de}(y, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) > > >) \\ \Rightarrow & \exists e \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \exists y \text{ leído}(e, y) \wedge \text{isa}(y, \text{libro}) \wedge \\ & \quad \text{de}(y, < \exists x \text{ isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) >) \\ \Rightarrow & \exists e \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \exists y \text{ leído}(e, y) \wedge \text{isa}(y, \text{libro}) \wedge \\ & \quad \exists x \text{ de}(y, x) \wedge \text{isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) \\ \Rightarrow & \exists e \exists y \text{ leer}(e) \wedge \text{lee}(e, \text{Juan}) \wedge \text{leído}(e, y) \wedge \text{isa}(y, \text{libro}) \wedge \\ & \quad \text{de}(y, x) \wedge \text{isa}(x, \text{autor}) \wedge \text{isa}(x, \text{joven}) \end{aligned}$$

d) Discuta una posible extensión de la gramática anterior para que cubra oraciones como “Juan lee un libro en el parque”.

Esta extensión se puede hacer de modo similar a la de incorporar grupos preposicionales a los grupos nominales.

La regla en este caso sería :

$$gv \rightarrow gv \text{ pp} \quad [gv.sem = ?]$$

El problema que se presenta es que la relación asociada a la preposición (“en” en este caso) debe tomar como uno de sus argumentos la variable de evento asociada al verbo (“lee” en este caso), pero esta variable no está accesible (está ligada por el cuantificador existencial).

Una opción es definir un operador especial que nos permita acceder a la variable de evento. Otra opción puede ser no ligar existencialmente las variables de evento, sino con operadores lambda.

### Ejercicio 6 [6 puntos]

Defina brevemente las medidas *Precision*, *Recall*, *E* y *F*. Explique para que se usan. La medida *E* lleva un parámetro que tiene un cometido especial. ¿Cuál es y para qué se utiliza?

*Las medidas de Precision y Recall, son usadas en muchas tareas de procesamiento de lenguaje natural. En particular, en el área de la Recuperación de Información sirven para medir cuan bueno es la efectividad de un sistema.*

*Precisión: capacidad de mostrar fundamentalmente documentos relevantes*

*Recuperación: capacidad de encontrar **todos** los documentos relevantes.*

Son valores entre 0 y 1, y cuanto más próximas a 1, mejor es la evaluación de la técnica empleada. Típicamente, a mayor recall menor precisión.

$$Precision = \frac{|\text{documentos recuperados relevantes}|}{|\text{documentos recuperados}|}$$

$$Recall = \frac{|\text{documentos recuperados relevantes}|}{|\text{documentos relevantes}|}$$

En cuanto a las medidas E y F, las mismas se definen en base a Precision y Recall, son variantes de éstas. En particular, la medida E es variante de medida F. Permite poner mayor énfasis en precisión o en recuperación según el valor del parámetro

$\beta$ , el cual controla el balance entre P y R

$$F = \frac{2PR}{P+R} \qquad E = \frac{(1+\beta^2)PR}{\beta^2P+R}$$

$\beta = 1$ : Pesos iguales (En este caso  $E = F$ ).

$\beta > 1$ : Mayor peso a precisión

$\beta < 1$ : Mayor peso a recuperación