

Introducción al Procesamiento de Lenguaje Natural

Noviembre 2022

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.
- viii) El total de puntos es 70.

Ejercicio 1 [16 puntos]

Para cada afirmación diga si es Verdadera o Falsa. Justifique.

- i) En la frase *La noche de las luces* la palabra **de** tiene categoría adverbio.
- ii) Algo bueno que tiene representar documentos mediante el centroide de los word embeddings de sus palabras es que se captura la importancia del orden de las palabras en la oración.
- iii) Skip-gram es uno de los algoritmos usados para construir word embeddings.
- iv) Un modelo de lenguaje de n-gramas permite calcular la probabilidad de una secuencia de palabras.
- v) El *accuracy* (o exactitud) es una buena medida para la performance de un clasificador binario en un corpus muy desbalanceado.
- vi) La medida tf-idf sirve para medir qué tan importante es un término en un documento en el Modelo Booleano de Recuperación de Información.
- vii) Para la construcción de corpus paralelos se pueden definir distintos tipos de alineación, como por ejemplo por documento o por oración.
- viii) La medida BLEU tiene como principal objetivo dar la altura máxima de un árbol de derivación para una gramática libre de contexto.

Ejercicio 2 [20 puntos]

a) Dadas las siguientes oraciones:

Los perros comen huesos.

Los perros comen.

Un perro come lentamente.

- i. Utilizando la notación bracket (corchetes rectos) segmente cada oración en sintagmas e indique: núcleo de cada sintagma, categoría de cada núcleo, categoría de cada sintagma.
- ii. Escriba una Gramática Libre de Contexto $G:(V,T,P,O)$ que las genere.
- iii. Dé un ejemplo – por medio de una derivación – de una oración agramatical desde el punto de vista del español, pero que sea permitida por la gramática construida en la parte ii. Justifique.

b) Considere la siguiente gramática libre de contexto G:

- O → GN GV | GV
- GN → Nom | Det Nom | GN GP
- GV → V GN | V GN GP | V
- GP → Prep GN
- Det → una | un | el | la
- Prep → con | de
- Nom → María | amigos | comienzo | final
- V → mirando | mira | miro

Aplice el algoritmo CKY justificando su razonamiento a partir de la gramática G para cada una de las siguientes entradas, indicando qué salida devuelve en cada caso:

- 1) "Mirando la final con amigos"
- 2) "Mirando con amigos"

Ejercicio 3 [20 puntos]

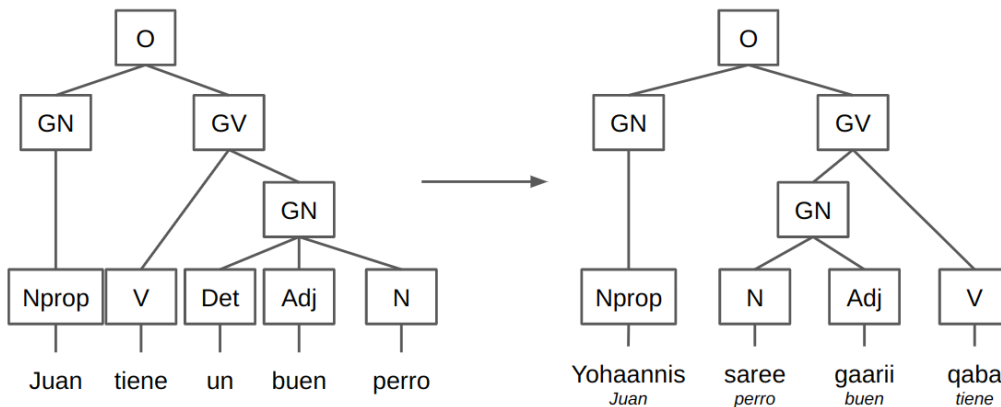
a) Sea la siguiente gramática con anotaciones semánticas:

- | | |
|-----------------|-----------------------------------|
| o → gn gv | o.sem = gn.sem(gv.sem) |
| gn → det nom | gn.sem = det.sem(nom.sem) |
| gn → npropio | gn.sem = npropio.sem |
| gv → v | gv.sem = v.sem |
| gv → neg v | gv.sem = neg.sem(v.sem) |
| nom → n | nom.sem = n.sem |
| det → un | det.sem = λP. λQ. ∃x P(x) ∧ Q(x) |
| neg → no | neg.sem = λP . λx . ¬P(x) |
| n → pájaro | n.sem = λx.pájaro(x) |
| npropio → Pablo | npropio.sem = λP.P(pablo) |
| v → habla | v.sem = λx. habla(x) |
| v → duerme | v.sem = λx. duerme(x) |
| adj → simple | adj.sem = λP.λx. simple(x) ∧ P(x) |

Utilizando las reglas anteriores dibuje el árbol sintáctico y derive la expresión lógica asociada a la oración:

Pablo no duerme

b) Suponga que se quiere construir un sistema de traducción automática del español al oromo basado en transferencia sintáctica. Escriba las reglas de transferencia para dicho sistema, basándose en la siguiente traducción de ejemplo:



Ejercicio 4 [14 puntos]

a) Se tiene un conjunto de noticias clasificadas en una de tres categorías (deportes, policial, economía) y los resultados de predicción de un sistema automático para esas tres categorías.

Noticia	Clase original	Predicción
n1	Deportes	Policial
n2	Deportes	Deportes
n3	Deportes	Deportes
n4	Policial	Policial
n5	Policial	Economía
n6	Deportes	Deportes
n7	Economía	Deportes
n8	Deportes	Economía
n9	Deportes	Policial
n10	Economía	Deportes

Construya la matriz de confusión. Calcule el accuracy total, y los valores de precisión, recall y medida-F por cada clase.

b) Sea el siguiente texto de una noticia:

Una turba enojada casi lincha a una turista que subió sin permiso los escalones del Castillo de Kukulcán, una de las nuevas siete maravillas del mundo moderno que se encuentra en la zona arqueológica de Chichén Itzá, al sureste de México.

El Instituto Nacional de Antropología e Historia (INAH) prohibió subirse al edificio sagrado de los mayas desde 2008, instaló un cordón de seguridad alrededor y anunció multas que van desde 50.000 (unos 2.558 dólares) a 100.000 pesos mexicanos (cerca de 5.115 dólares), dependiendo del daño que se cause a la estructura.

El director general del INAH, Diego Prieto, aún no brindan informes del incidente, por lo que la multitud sigue exigiendo cárcel y la expulsión de Yucatán, “y, si es del extranjero, que se vaya de México”, gritaban los presentes.

En incidentes anteriores, el INAH informó que sí se castigaría en acuerdo con el Ministerio de Turismo, quien establece las penas y sanciones contra aquellas personas que dañen o exploten monumentos arqueológicos inmuebles sin autorización del Instituto.

En un proceso de extracción de información sobre el texto anterior dé:

- i. al menos 3 entidades con nombre que pertenezcan a distintas clases.
- ii. al menos 2 relaciones binarias indicando sus argumentos.