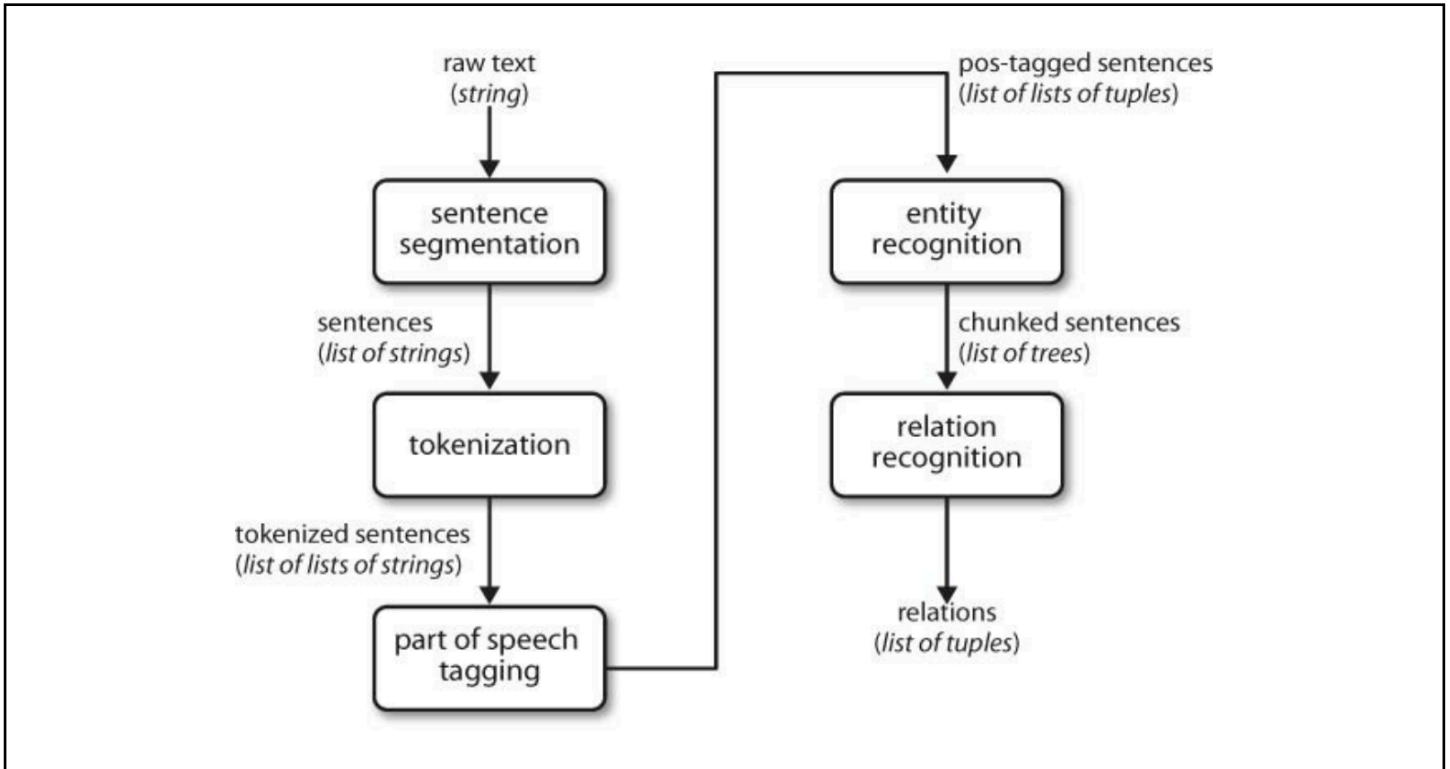




Student:	García González Aarón Antonio
Group:	3CV9
Learning unit:	Natural language processing
Teacher:	Kolesnikova Olga
Homework #12:	Information extraction
Date:	Monday, July 9th, 2020

We are going to use the next architecture for to do this activity:



To get sentences from text files, sentence segmentation, tokenization, I encoded the next function:

```
1. def getSentences(PATH, OBJECT):
2.     result = []
3.
4.     for comment in os.listdir(PATH+OBJECT+"/"):
5.         path_file = PATH+OBJECT+"/"+comment
6.
7.         f = open(path_file, 'r', encoding = 'ISO-8859-1', errors="ignore")
8.         text = f.read()
9.         f.close()
10.
11.         text = re.sub('\n', ' ', text)
12.         text = text.strip()
13.
14.         sentences = nltk.sent_tokenize(text)
15.         sentences = [sentence.strip() for sentence in sentences if len(sentence) > 12]
16.         sentences = [nltk.word_tokenize(sentence) for sentence in sentences]
17.
18.         result += sentences
19.
20.     return result
```

To do POS tag, I used the “Stanford Log-linear Part-Of-Speech Tagger”, In this time the reason why I used this POS tagger instead of my POS tagger developed in class is only to know how functions this and so for next times can I use the best for my opinion or needs.

```
1. tagger = "/Users/aarongarcia/desktop/12_Info_extraction/stanford-tagger-4.0.0/models/spanish-ud.tagger"
2. jar = "/Users/aarongarcia/desktop/12_Info_extraction/stanford-tagger-4.0.0/stanford-postagger.jar"
3. reference = "https://nlp.stanford.edu/software/spanish-faq.shtml#tagset"
4. etiquetador = StanfordPOSTagger(tagger,jar)
5. saveFilePKL([etiquetador.tag(oracion) for oracion in oraciones],"etiquetas.pkl")
6. oraciones = getFilePKL("etiquetas.pkl")
```

This tool has only one issue, is so slow, that is the reason why I use pickle files.

The first point of this homework is getting noun phrases, I obtained a text file, I show some of these phrases:

1. (NP hoteles/NOUN)	41. (NP calidad/NOUN del /ADP hotel/NOUN)
2. (NP estrellas/NOUN)	42. (NP las/DET estrellas/NOUN)
3. (NP mejores/ADJ)	43. (NP caro/ADJ)
4. (NP no/ADV creo/VERB)	44. (NP euros/noche/NOUN)
5. (NP este/DET hotel/NOUN)	45. (NP Realmente/ADV conozco/VERB)
6. (NP Las/DET habitaciones/NOUN)	46. (NP hoteles/NOUN)
7. (NP pequeñas/ADJ)	47. (NP estrellas/NOUN)
8. (NP no/ADV tienen/VERB)	48. (NP una/DET pasada/ADJ)
9. (NP camas/NOUN de/ADP matrimonio/NOUN)	49. (NP Queria/PROPN)
10. (NP terraza/NOUN)	50. (NP mi/DET novio/NOUN)
11. (NP decoradas/ADJ)	51. (NP un/DET finde/NOUN)
12. (NP manera/NOUN antigua/ADJ)	52. (NP un/DET hotel/NOUN)
13. (NP algunas/DET grietas/NOUN)	53. (NP estrellas/NOUN)
14. (NP Los/DET mandos/NOUN)	54. (NP una/DET gran/ADJ desilusión/NOUN)
15. (NP la/DET tele/NOUN)	55. (NP Pocos/DET establecimientos/NOUN hoteleros /ADJ)
16. (NP no/ADV funcionan/VERB)	56. (NP Madrid/PROPN)
17. (NP no/ADV tienen/VERB)	57. (NP mejor/ADJ situación/NOUN)
18. (NP canales/NOUN)	58. (NP Tryp/PROPN)
19. (NP no/ADV comentar/VERB)	59. (NP Ambassador/PROPN)
20. (NP además/ADV estan/VERB)	60. (NP establecimiento/NOUN)
21. (NP sucias/NOUN)	61. (NP estrellas/NOUN)
22. (NP los/DET baños/NOUN)	62. (NP la/DET cadena/NOUN)
23. (NP la/DET habitación/NOUN)	63. (NP Sol-Meliá/PROPN)
24. (NP Las/DET instalaciones/NOUN)	64. (NP ubicado/ADJ)
25. (NP viejas/ADJ y/CCONJ muy/ADV descuidadas/ADJ)	65. (NP Cuesta/PROPN)
26. (NP el/DET gimnasio/NOUN)	66. (NP Santo/PROPN)
27. (NP la/DET piscina/NOUN climatizada/ADJ solo/ADJ)	67. (NP Dominfgo/PROPN)
28. (NP la/DET tarde/NOUN)	68. (NP un/DET paso/NOUN del /ADP Palacio/PROPN)
29. (NP masajes/NOUN)	69. (NP Real/PROPN)
30. (NP no/ADV había/VERB)	70. (NP Senado/PROPN)
31. (NP masajista/NOUN)	71. (NP Ópera/PROPN)
32. (NP mala/ADJ)	72. (NP la/DET Gran/ADJ)
33. (NP el/DET entorno/NOUN)	73. (NP Vía/PROPN)
34. (NP bonito/ADJ)	74. (NP Plaza/PROPN)
35. (NP personal/NOUN del /ADP hotel/NOUN)	75. (NP España/PROPN)
36. (NP amable/ADJ)	76. (NP el/DET antiguo/ADJ)
37. (NP las/DET zonas/NOUN comunes/ADJ)	77. (NP Palacio/PROPN)
38. (NP desayuno/NOUN tipo/NOUN buffet/NOUN)	78. (NP Duques/PROPN)
39. (NP lo/DET mejor/ADJ)	79. (NP Granada/PROPN)
40. (NP una/DET revisión/NOUN)	80. (NP Ega/PROPN)

The grammar that I used to do this:

```
1. grammar_NP = r"""NP:
2.             {<DET><NOUN><ADP><PROPN>}
3.             {<NOUN><NOUN><NOUN>*}
4.             {<NOUN><ADP><NOUN>}
5.             {<DET>?<ADJ><NOUN>}
6.             {<DET>?<NOUN><ADJ>*}
7.             {<DET><ADJ>}
8.             {<ADV><VERB><NOUN>?}
9.             {<DET><ADV>?<ADJ>}
10.            {<ADJ><CCONJ><ADV>?<ADJ>}
11.            {<ADJ>?<NOUN>?}
12.            {<PROPN><ADJ>*}
13.            """
```

The second point is getting sentence extraction, the grammar that I used to do this:

```
1. grammar_CHUNK = r"""CHUNK:
2.             {<DET><NOUN><AUX><ADJ>}
3.             {<NOUN><ADP><NOUN>}
4.             {<ADJ><ADP><NOUN><ADJ>}
5.             {<DET>?<NOUN><VERB><ADJ>}
6.             {<DET><NOUN><ADJ>}
7.             {<ADV><VERB><NOUN>}
8.             {<DET><NOUN><ADV><VERB><NOUN|VERB>*}
9.             {<NOUN>+<ADP|DET>+<ADJ|NOUN>+}
10.            {<DET>?<NOUN><ADP><DET>?<NOUN|VERB|ADJ>+}
11.            {<DET><NOUN><ADP><DET><NOUN>}
12.            {<DET><NOUN><AUX><DET>?<NOUN>}
13.            {<PROPN><ADP><PROPN|NOUN>+}
14.            {<DET><NOUN><VERB><DET>+<NOUN>}
15.            """
```

I obtained a text file; I show some of these phrases:

```
1. (CHUNK Las/DET habitaciones/NOUN son/AUX pequeñas/ADJ)
2. (CHUNK camas/NOUN de/ADP matrimonio/NOUN)
3. (CHUNK decoradas/ADJ de/ADP manera/NOUN antigua/ADJ)
4. (CHUNK la/DET tele/NOUN no/ADV funcionan/VERB)
5. (CHUNK además/ADV estan/VERB sucias/NOUN)
6. (CHUNK Las/DET instalaciones/NOUN estan/VERB viejas/ADJ)
7. (CHUNK la/DET piscina/NOUN climatizada/ADJ)
8. (CHUNK no/ADV habia/VERB masajista/NOUN)
9. (CHUNK el/DET entorno/NOUN es/AUX bonito/ADJ)
10. (CHUNK personal/NOUN del/ADP hotel/NOUN)
11. (CHUNK las/DET zonas/NOUN comunes/ADJ)
12. (CHUNK desayuno/NOUN tipo/NOUN buffet/NOUN de/ADP lo/DET mejor/ADJ)
13. (CHUNK calidad/NOUN del/ADP hotel/NOUN)
14. (CHUNK Realmente/ADV conozco/VERB hoteles/NOUN)
15. (CHUNK novio/NOUN un/DET finde/NOUN)
16. (CHUNK Pocos/DET establecimientos/NOUN hoteleros/ADJ)
17. (CHUNK estrellas/NOUN de/ADP la/DET cadena/NOUN)
18. (CHUNK Cuesta/PROPN de/ADP Santo/PROPN Domingó/PROPN)
19. (CHUNK Plaza/PROPN de/ADP España/PROPN)
20. (CHUNK Duques/PROPN de/ADP Granada/PROPN)
21. (CHUNK su/DET estructura/NOUN original/ADJ)
22. (CHUNK un/DET establecimiento/NOUN muy/ADV frecuentado/VERB)
23. (CHUNK personajes/NOUN de/ADP la/DET pequeña/ADJ pantalla/NOUN)
24. (CHUNK habitaciones/NOUN para/ADP sus/DET estancias/NOUN)
25. (CHUNK zona/NOUN de/ADP teatros/NOUN)
26. (CHUNK Tv/NOUN de/ADP pago/NOUN)
```