

Relazione esercitazione – False friends

L'obiettivo di questa esercitazione è stato quello di creare un algoritmo che fosse in grado di **individuare coppie di termini della lingua inglese che potessero essere considerati dei “false friends”**. Poiché in questo caso, consideriamo una sola lingua, possiamo considerare come false friends due termini che dal punto di vista della forma sono quasi identici ma che differiscono notevolmente nel significato che assumono.

Le risorse principali utilizzate per questa esercitazione sono state:

- **Wordnet**, per ottenere i possibili sensi associati ad un certo termine.
- **Semcor**, che invece è un corpus già utilizzato nella seconda parte del corso. Esso contiene al proprio interno circa 37000 frasi annotate a diversi livelli (per ogni parola in ogni frase contiene anche il synset di WN di riferimento), che è stato utilizzato per ottenere le frasi nelle quali andare a cercare le possibili parole che potrebbero essere false friends.

L'algoritmo sviluppato è composto dai seguenti passi:

- 1) L'utilizzo di una procedura chiamata **“individuazione_coppie_parole_synset”** che si preoccupa di andare ad individuare tutte le coppie **parola-synset** partendo da un insieme di frasi casuali prese da Semcor. Tutte le coppie trovate (senza duplicati) vengono serializzate. Il numero di frasi considerate per l'estrazione è stato 3000.
- 2) Una volta terminata la procedura descritta al passo 1), viene deserializzata la lista di coppie parola-synset e subito dopo viene chiamata la funzione **“get_coppie_parole_con_dist_giusta”** che prendendo in input la lista appena citata, si occuperà di trovare tutte le possibili coppie di parole che dal punto di vista della **distanza edit risultano essere molto simili tra loro**. La **distanza edit** è semplicemente il numero di operazioni di: **rimozione, inserimento e modifica** che sono necessarie per trasformare una parola in un'altra. Il numero massimo di operazioni richieste è un parametro che viene scelto a priori e in questa esercitazione il valore di soglia deciso è stato 2, quindi sono state selezionate tutte le coppie di parole che avevano una distanza edit minore di 2. La funzione utilizzata per calcolare l'edit è presente nella libreria **nlTK** ed è chiamata **“edit_distance”**. Ogni parola prima di essere confrontata con un'altra è stata lemmatizzata.
- 3) Dopodiché, una volta ottenute tutte le coppie di termini con edit distance minore di 2, per ciascuna coppia di parole sono stati eseguiti i seguenti passi:
 - 3.1) Sono stati presi da Wordnet tutti i possibili synsets associati a ciascuna parola della coppia.
 - 3.2) Per il task di Word Similarity, è stata utilizzata la misura di similarità chiamata **“Wu&Palmer”**. Il motivo di tale scelta è dovuto al fatto che durante le sperimentazioni fatte in un'esercitazione con il professor Radicioni, essa è risultata essere la metrica che restituiva risultati di similarità più vicini a quelli forniti dagli esseri umani.

Similarity Wu&Palmer:

$$sim(c1, c2) = \frac{2 * depth(LCS)}{depth(c1) + depth(c2)}$$

Ove:

LCS = Lowest Common Subsumer (è il primo antenato comune, presente come synset in Wordnet, tra i sensi c1 e c2).

depth(x) = è una funzione che misura la distanza tra la radice di Wordnet e il synset x.

3.3) Il valore di similarità tra ogni coppia di termini (w1, w2) è stato ottenuto eseguendo questo calcolo:

$$sim(w1, w2) = \underset{c1 \in S(w1), c2 \in S(w2)}{MAX} [sim(c1, c2)]$$

dove:

S(w1) = insieme di tutti i synsets associati al termine w1.

S(w2) = insieme di tutti i synsets associati al termine w2.

c1 = possibile synset associato al termine w1.

c2 = possibile synset associato al termine w2.

4. Se, il valore di similarità calcolato, **è al di sotto di una certa soglia che è stata fissata a 0.20**, allora, poiché questo vuol dire che le due parole dal punto di vista semantico sono abbastanza diverse tra loro, la coppia di termini viene memorizzata all'interno di una lista. Alla fine, dopo aver valutato tutte le varie coppie di parole da considerare, l'algoritmo stampa la lista di coppie di parole della lingua inglese che possono essere considerate false friends.

Risultati ottenuti

Il numero totale di false friends trovati dall'algoritmo è **120**, essi sono presenti nel file chiamato "**False friends – 120**".

Di seguito vengono riportati **tutti** i false friends trovati dall'algoritmo e il loro valore di similarità:

```
2      ['published', 'publisher', 0.181818181818182],
3      ['four', 'for', 0],
4      ['used', 'sed', 0],
5      ['pulled', 'pulley', 0.133333333333333],
6      ['many', 'may', 0.181818181818182],
7      ['many', 'mary', 0.181818181818182],
8      ['led', 'sed', 0],
9      ['brother', 'bother', 0.181818181818182],
10     ['aunt', 'taunt', 0.133333333333333],
11     ['set', 'sed', 0],
12     ['the', 'tie', 0],
13     ['victory', 'victor', 0.15384615384615385],
14     ['long', 'lung', 0.166666666666666],
15     ['sat', 'oat', 0.181818181818182],
16     ['bloat', 'boat', 0.181818181818182],
17     ['cattle', 'battle', 0.09523809523809523],
18     ['lung', 'hung', 0.166666666666666],
19     ['lung', 'flung', 0.15384615384615385],
20     ['lung', 'clung', 0.166666666666666],
21     ['find', 'finn', 0.181818181818182],
22     ['out', 'oat', 0.181818181818182],
23     ['navy', 'wavy', 0.181818181818182],
24     ['see', 'sed', 0],
25     ['lady', 'lay', 0.181818181818182],
26     ['lady', 'lacy', 0.181818181818182],
27     ['came', 'cafe', 0.166666666666666],
28     ['had', 'dad', 0.133333333333333],
29     ['wide', 'wife', 0.181818181818182],
30     ['thanking', 'thinking', 0.181818181818182],
31     ['form', 'for', 0],
32     ['join', 'goin', 0],
33     ['know', 'knox', 0.166666666666666],
34     ['sped', 'sed', 0],
35     ['law', 'lawn', 0.14285714285714285],
36     ['tricked', 'ticked', 0.166666666666666],
37     ['got', 'god', 0.181818181818182],
38     ['lawn', 'jawn', 0],
39     ['went', 'cent', 0.181818181818182],
40     ['feat', 'felt', 0.181818181818182],
41     ['feat', 'eat', 0.166666666666666],
42     ['feat', 'neat', 0.181818181818182],
```

43	['with', 'wish', 0],	85	['composer', 'composed', 0.16666666666666666],
44	['won', 'ion', 0.18181818181818182],	86	['peer', 'per', 0],
45	['royale', 'royal', 0],	87	['vast', 'vas', 0.18181818181818182],
46	['mary', 'marry', 0.15384615384615385],	88	['car', 'war', 0.15384615384615385],
47	['dad', 'bad', 0.13333333333333333],	89	['car', 'far', 0.16666666666666666],
48	['dad', 'did', 0.13333333333333333],	90	['quiet', 'quirt', 0.19047619047619047],
49	['dad', 'dead', 0.13333333333333333],	91	['sit', 'skit', 0.18181818181818182],
50	['word', 'sword', 0.125],	92	['guest', 'guess', 0.15384615384615385],
51	['sea', 'sed', 0],	93	['economics', 'economic', 0.15384615384615385],
52	['own', 'owl', 0.13333333333333333],	94	['poet', 'poem', 0.14285714285714285],
53	['later', 'liter', 0.18181818181818182],	95	['robe', 'role', 0.18181818181818182],
54	['gain', 'goin', 0],	96	['neutralism', 'neutralist', 0.11111111111111111],
55	['try', 'tray', 0.15384615384615385],	97	['earn', 'ear', 0.16666666666666666],
56	['creed', 'creek', 0.14285714285714285],	98	['per', 'pea', 0],
57	['creed', 'cried', 0.16666666666666666],	99	['per', 'peru', 0],
58	['care', 'car', 0.15384615384615385],	100	['per', 'pen', 0],
59	['care', 'cafe', 0.15384615384615385],	101	['ton', 'torn', 0.18181818181818182],
60	['knee', 'kneel', 0.16666666666666666],	102	['lurked', 'lured', 0.13333333333333333],
61	['backed', 'backer', 0.16666666666666666],	103	['near', 'ear', 0.18181818181818182],
62	['haste', 'hasty', 0.18181818181818182],	104	['coin', 'goin', 0],
63	['fist', 'fast', 0.16666666666666666],	105	['tutor', 'tumor', 0.125],
64	['fist', 'list', 0.15384615384615385],	106	['wipe', 'wife', 0.15384615384615385],
65	['fist', 'fit', 0.16666666666666666],	107	['send', 'sed', 0],
66	['bought', 'bough', 0.13333333333333333],	108	['swore', 'sword', 0.14285714285714285],
67	['eat', 'oat', 0.16666666666666666],	109	['familiar', 'familiar', 0],
68	['eat', 'ear', 0.16666666666666666],	110	['red', 'sed', 0],
69	['ever', 'fever', 0.18181818181818182],	111	['quirt', 'quit', 0.15384615384615385],
70	['finn', 'fine', 0.18181818181818182],	112	['gust', 'just', 0.15384615384615385],
71	['done', 'doe', 0.13333333333333333],	113	['thanked', 'shanked', 0.14285714285714285],
72	['add', 'adc', 0],	114	['whisked', 'whiskey', 0.14285714285714285],
73	['hear', 'heir', 0.18181818181818182],	115	['manufactured', 'manufacturer', 0.18181818181818182],
74	['hear', 'ear', 0.18181818181818182],	116	['lamp', 'damp', 0.16666666666666666],
75	['ion', 'ton', 0.125],	117	['sed', 'bed', 0],
76	['sworn', 'sword', 0.15384615384615385],	118	['sed', 'fed', 0],
77	['going', 'goin', 0],	119	['sed', 'sex', 0],
78	['for', 'far', 0],	120	['stir', 'stair', 0.15384615384615385],
79	['for', 'fox', 0],	121	['inner', 'sinner', 0.18181818181818182]
80	['for', 'fort', 0],	122	
81	['for', 'fog', 0],		
82	['boy', 'buy', 0.18181818181818182],		
83	['boy', 'boo', 0.15384615384615385],		
84	['like', 'pike', 0.18181818181818182],		