

DEVELOPMENT OF WORD BASED MACHINE TRANSLATION SYSTEM

A Mini Project Report

Submitted

To

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

**BY
SUVAM BASAK**



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

**JIS COLLEGE OF ENGINEERING
BLOCK 'A' PHASE 'III' KALYNAI NADIA-741235
December, 2014**



CERTIFICATE

We hereby notify that the work which is being presented (Report) entitled “**project name**” and submitted to the Department **Computer Science and Engineering** of JIS college of Engineering is an authentic record of our own work carried out during a period from “**duration**” under the supervision of **Name & Designation of supervisor(s)**, CSE Department.

Signature of Candidates

R.No.

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Signature of Supervisor(s)

Date: 20th November

**Name & Designation of
Supervisor(s)**

Head
CSE Department

JIS College of Engineering
(An Autonomous Institution)

Department of Computer Science and Engineering

UG Project Proposal'2015-2016

Title: DEVELOPMENT OF WORD BASED MACHINE TRANSLATION SYSTEM
USING “**NATURAL LANGUAGE PROCESSING IN PYTHON**”.

Abstract:

“Natural language processing” here refers to the use and ability of system to process sentences in a natural language such as English. NLP is highly essential in artificial intelligence. We used Python for it, because Python is a simple yet powerful programming language which excellent functionality for processing data.

Software Requirement : Some several free software packages is required for our processing purpose. Python version 2.7, Natural language processing toolkit (NLTK), NLTK-Data, NumPy etc.

Language translation : At first we used a English-Spanish Parallel Corpus [The Europarl parallel corpus is extracted from the proceedings of the European Parliament]. Then we extracted words from the corpus to another file using a python program which was written by us. After that we needed to map English words with Spanish words. For the mapping purpose we used Anymalign.py and Treetagger was used for finding parts of speech. After Anymalign processing we extracted first and second column (English and Spanish) from the output file. Then we made a Program to translate English text to Spanish text. The program takes input only in English text and split it word by word and then searches for corresponding Spanish word then displays.

MEMBER DETAIL :

- 1.Suvam Basak
- 2.Shaif Aslam
- 3.Sayantan Pal
- 4.Sandip Pore

Mentor's Name : Sainik Kr. Mahata

CONTENT

- 1 Introduction
- 2 Types of Machine Translation
- 3 Things required
- 4 Assignments
- 5 Advantage of machine translation
- 6 Disadvantage of machine translation
- 7 Resource
- 8 Reference
- 9 Conclusion

Introduction:

Natural language processing (NLP) is a field of computer science, artificial intelligence, and computational linguistics concerned with the interactions between computers and human (natural) language. As such, NLP is related to the area of human-computer interaction. Many challenges in NLP involve Natural language understanding, that is, enabling computer to derive meaning from human or natural language input, and others involve natural language generation. Now here our job is Machine translation, so Machine translation some thing referred to by the abbreviation MT is a sub-field of computational linguistics that investigates the use of software to translate text or speech from one language to another.

On a basic level, MT performs simple substitution of words in one language for words in another, but that alone usually cannot produce a good translation of a text because recognition of whole phrases and their closest counterparts in the target language is needed. Solving this problems in linguistic typology, translation of idioms, and the isolation of anomalies.

Current machine translation software often allows for customization by domain or profession , improving output by limiting the scope of allowable substitutions. This technique is particularly effective in domain where formal or formulaic language is used . It follows that machine translation of government and legal documents more readily produces output that conversation or less standardised text.

Improved output quality can also be achieved by human intervention: for example, some systems are able to translate more accurately if the user has unambiguously identified which words in the text are proper names. With the assistance of these techniques, MT has proven useful as a tool to assist human translators and, in a very limited number of cases, can even produce output that can be used as is.

The progress and potential of machine translation have been debated much through its history. Since the 1950s, a number of scholars have questioned the possibility of achieving fully automatic machine translation of high quality. Some critics claim that there are in-principle obstacles to automatizing the translation process.

Types of machine translation :

Total four types of machine translation are available.

Those are :

1. Word based machine translation. [WBMT]
2. Example based machine translation. [EBMT]
3. Phrase based machine translation. [PBMT]
4. Statistical machine translation. [SMT]

From four types of machine translation our topic is Word based machine translation [WBMT].

Things required :

1. Python programming language : Python is a widely used general-purpose, high-level programming language. Its design philosophy emphasizes code readability, and its syntax allows programmers to express concepts in fewer lines of code than would be possible in languages such as C++ or Java. The language provides constructs intended to enable clear programs on both a small and large scale.

Python supports multiple programming paradigms, including object-oriented, imperative and functional programming or procedural styles. It features a dynamic type system and automatic memory management and has a large and comprehensive standard library.

Python interpreters are available for installation on many operating systems, allowing Python code execution on a wide variety of systems. Using third-party tools, such as Py2exe or Pyinstaller, Python code can be packaged into stand-alone executable programs for some of the most popular operating systems, allowing the distribution of Python-

based software for use on those environments without requiring the installation of a Python interpreter.

CPython, the reference implementation of Python, is free and open-source software and has a community-based development model, as do nearly all of its alternative implementations. CPython is managed by the non-profit Python Software Foundation.

2.Natural Language Toolkit : NLTK is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to over 50 corpora and lexical resource such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, wrappers for industrial-strength NLP libraries, and an active discussion forum. Thanks to a hands-on guide introducing programming fundamentals alongside topics in computational linguistics, plus comprehensive API documentation, NLTK is suitable for linguists, engineers, students, educators, researchers, and industry users alike. NLTK is available for Windows, Mac OS X, and Linux. Best of all, NLTK is a free, open source, community-driven project.

NLTK has been called “a wonderful tool for teaching, and working in, computational linguistics using Python,” and “an amazing library to play with natural language.”

Natural Language Processing with Python provides a practical introduction to programming for language processing. Written by the creators of NLTK, it guides the reader through the fundamentals of writing Python programs, working with corpora, categorizing text, analyzing linguistic structure, and more. The book is being updated for Python 3 and NLTK 3. (The original Python 2 version is still available at http://nltk.org/book_1ed. And in case of our project we used Python 2.7 version)

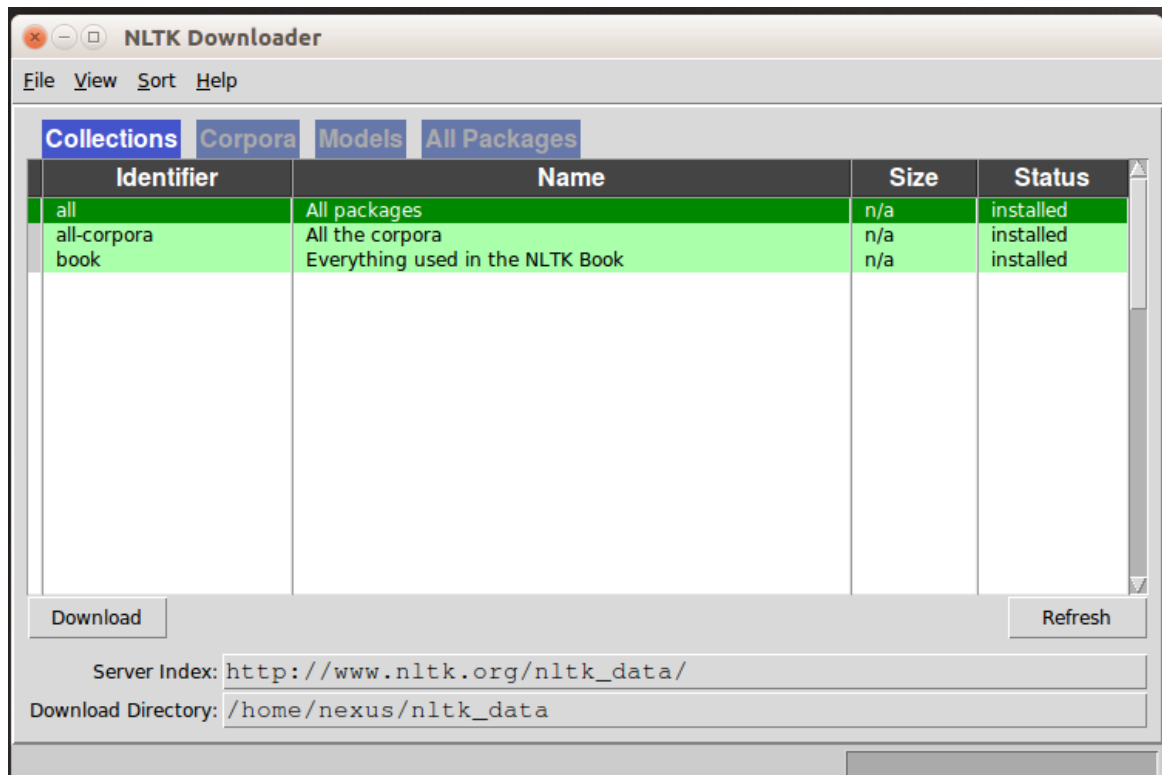
NLTK installation command :

1. `sudo apt-get install python-pip`
2. `sudo pip install -U numpy`
3. `sudo pip install -U nltk`
4. For testing installation : Run python then `import nltk`

3.NLTK_data : NLTK comes with many corpora, toy grammars, trained models, etc. A complete list is posted at:http://nltk.org/nltk_data/. After installing NLTK then used NLTK's data downloader to download nltk data.

```
nexus@nexus-Inspiron-3543: ~  
nexus@nexus-Inspiron-3543:~$ python  
Python 2.7.10 (default, Oct 14 2015, 16:09:02)  
[GCC 5.2.1 20151010] on linux2  
Type "help", "copyright", "credits" or "license" for more information.  
>>> import nltk  
>>> nltk.download()
```

Installation command



NLTK Downloader

4.Parallel corpus : A parallel corpus is a corpus that contains a collection of original texts in language L_1 and their translations into a set of languages $L_2 \dots L_n$. In most cases, parallel corpora contain data from only two languages. Closely related to parallel corpora are 'comparable corpora', which consists of texts from two or more languages which are similar in genre, topic, register etc. without, however, containing the same content.

<http://www.statmt.org/europarl/> this link used to download the English to Spanish parallel corpus.

5.Word alignment : There are three types of word alignment are available.

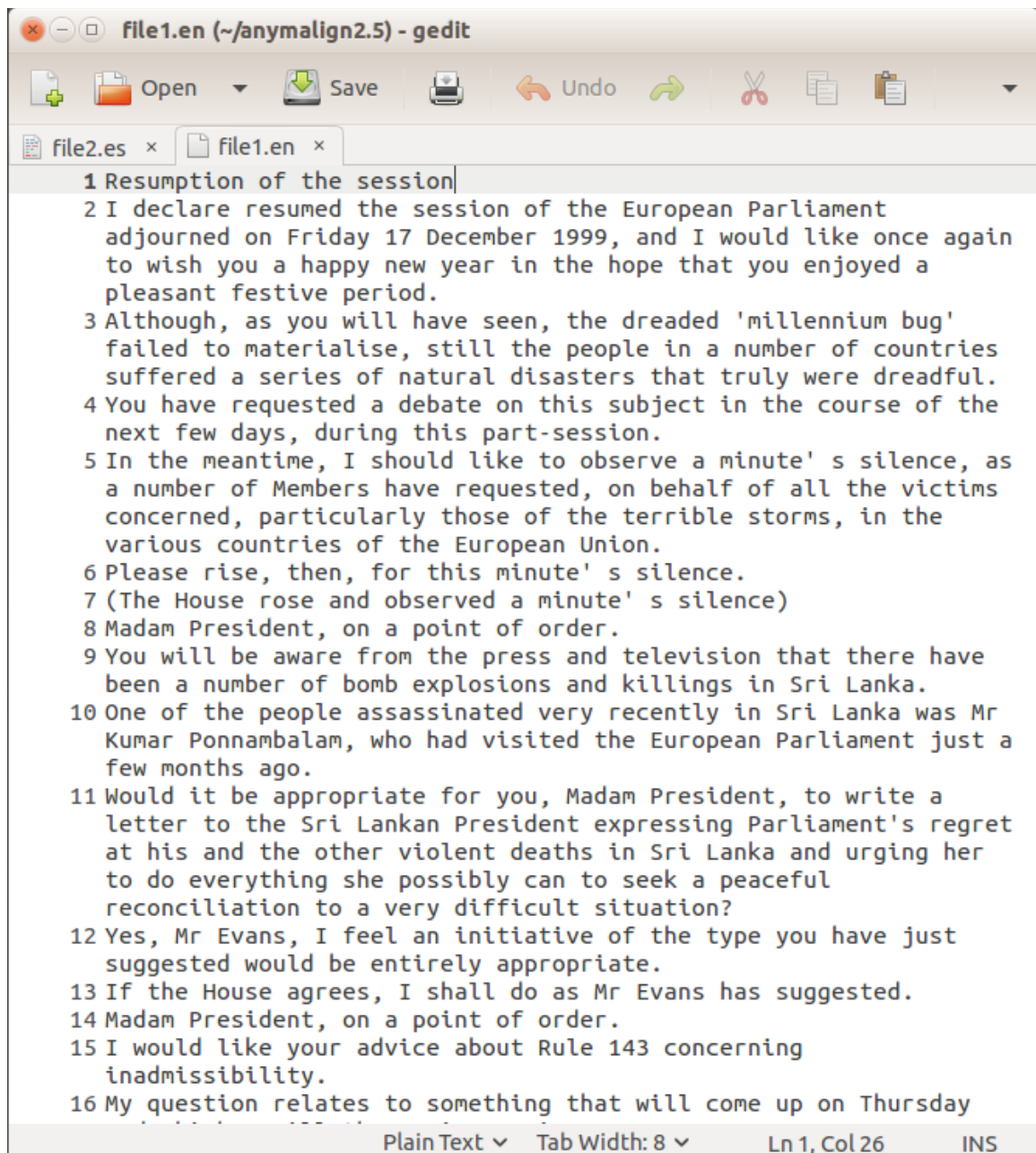
Those are :

1. GIZA++
2. MGIZA++
3. Anymalign

In case of our project, we used anymalign2.5 . Anymalign is a multilingual sub-sentential aligner. It can extract lexical equivalences from sentence-aligned parallel corpora. Its main advantage over other similar tools is that it can align any number of languages simultaneously. This package is downloaded from <https://anymalign.limsi.fr/>. Other info like how to use it, input file format, output file format also available in this link.

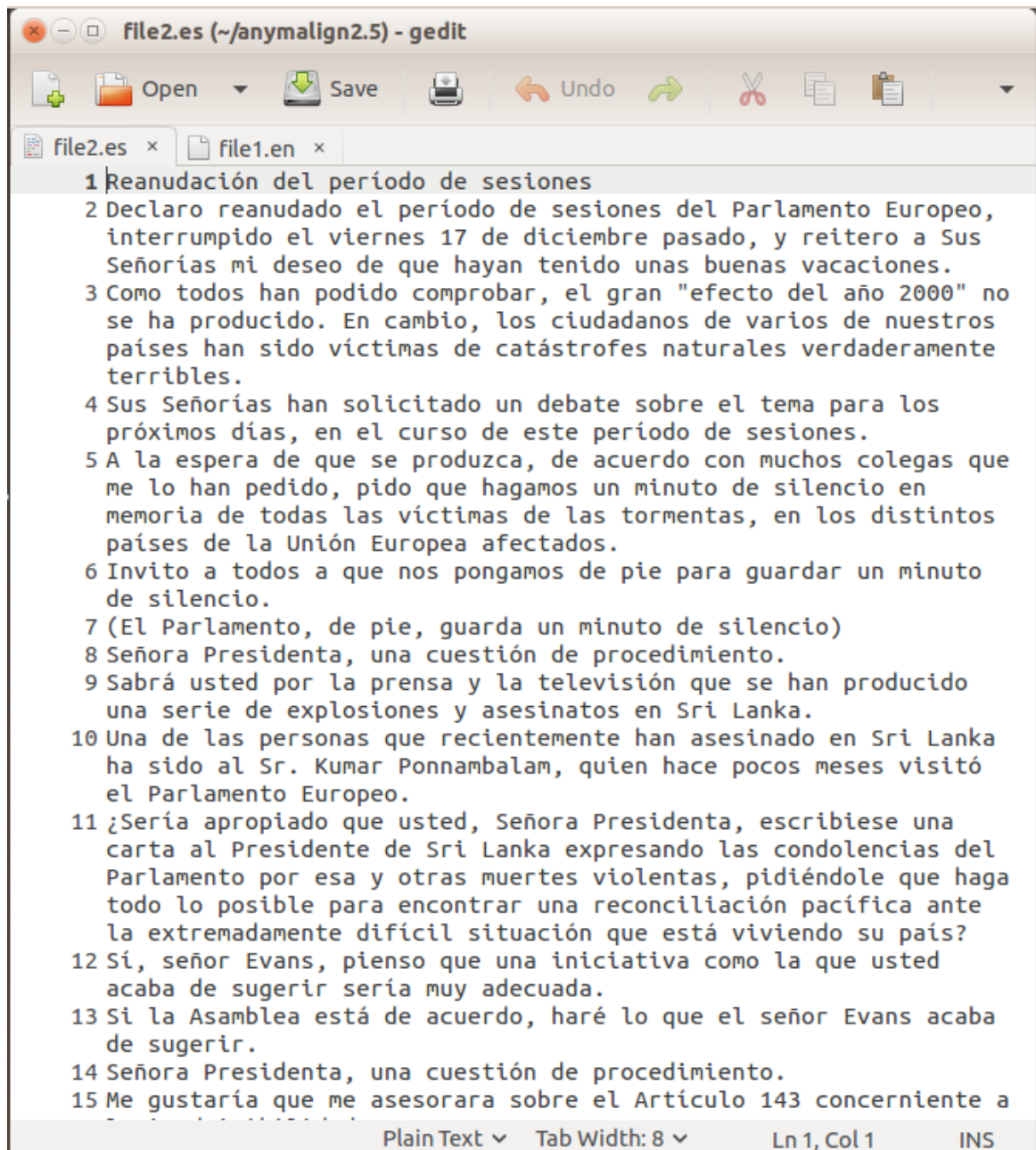
Input to the system : anymalign.py can read input data in separate files, where each file may contain one or more languages. Typically, multilingual texts are available in separate files, one file per language, one sentence per line, the corresponding lines being translations (all files have the same number of lines). For instance, you may have a tiny trilingual corpus in English, French, and German, where each file is made of 2 lines:

Here our input is English and Spanish text file.



```
file1.en (~/anymalign2.5) - gedit
Open Save Undo
file2.es x file1.en x
1 Resumption of the session|
2 I declare resumed the session of the European Parliament
  adjourned on Friday 17 December 1999, and I would like once again
  to wish you a happy new year in the hope that you enjoyed a
  pleasant festive period.
3 Although, as you will have seen, the dreaded 'millennium bug'
  failed to materialise, still the people in a number of countries
  suffered a series of natural disasters that truly were dreadful.
4 You have requested a debate on this subject in the course of the
  next few days, during this part-session.
5 In the meantime, I should like to observe a minute' s silence, as
  a number of Members have requested, on behalf of all the victims
  concerned, particularly those of the terrible storms, in the
  various countries of the European Union.
6 Please rise, then, for this minute' s silence.
7 (The House rose and observed a minute' s silence)
8 Madam President, on a point of order.
9 You will be aware from the press and television that there have
  been a number of bomb explosions and killings in Sri Lanka.
10 One of the people assassinated very recently in Sri Lanka was Mr
  Kumar Ponnambalam, who had visited the European Parliament just a
  few months ago.
11 Would it be appropriate for you, Madam President, to write a
  letter to the Sri Lankan President expressing Parliament's regret
  at his and the other violent deaths in Sri Lanka and urging her
  to do everything she possibly can to seek a peaceful
  reconciliation to a very difficult situation?
12 Yes, Mr Evans, I feel an initiative of the type you have just
  suggested would be entirely appropriate.
13 If the House agrees, I shall do as Mr Evans has suggested.
14 Madam President, on a point of order.
15 I would like your advice about Rule 143 concerning
  inadmissibility.
16 My question relates to something that will come up on Thursday
  ...
Plain Text Tab Width: 8 Ln 1, Col 26 INS
```

Frist input for anymalign.py



```
file2.es (~/anymalign2.5) - gedit
Open Save Print Undo
file2.es x file1.en x
1 Reanudación del período de sesiones
2 Declaro reanudado el período de sesiones del Parlamento Europeo,
  interrumpido el viernes 17 de diciembre pasado, y reitero a Sus
  Señorías mi deseo de que hayan tenido unas buenas vacaciones.
3 Como todos han podido comprobar, el gran "efecto del año 2000" no
  se ha producido. En cambio, los ciudadanos de varios de nuestros
  países han sido víctimas de catástrofes naturales verdaderamente
  terribles.
4 Sus Señorías han solicitado un debate sobre el tema para los
  próximos días, en el curso de este período de sesiones.
5 A la espera de que se produzca, de acuerdo con muchos colegas que
  me lo han pedido, pido que hagamos un minuto de silencio en
  memoria de todas las víctimas de las tormentas, en los distintos
  países de la Unión Europea afectados.
6 Invito a todos a que nos pongamos de pie para guardar un minuto
  de silencio.
7 (El Parlamento, de pie, guarda un minuto de silencio)
8 Señora Presidenta, una cuestión de procedimiento.
9 Sabrá usted por la prensa y la televisión que se han producido
  una serie de explosiones y asesinatos en Sri Lanka.
10 Una de las personas que recientemente han asesinado en Sri Lanka
  ha sido al Sr. Kumar Ponnambalam, quien hace pocos meses visitó
  el Parlamento Europeo.
11 ¿Sería apropiado que usted, Señora Presidenta, escribiese una
  carta al Presidente de Sri Lanka expresando las condolencias del
  Parlamento por esa y otras muertes violentas, pidiéndole que haga
  todo lo posible para encontrar una reconciliación pacífica ante
  la extremadamente difícil situación que está viviendo su país?
12 Sí, señor Evans, pienso que una iniciativa como la que usted
  acaba de sugerir sería muy adecuada.
13 Si la Asamblea está de acuerdo, haré lo que el señor Evans acaba
  de sugerir.
14 Señora Presidenta, una cuestión de procedimiento.
15 Me gustaría que me asesorara sobre el Artículo 143 concerniente a
  ...
Plain Text Tab Width: 8 Ln 1, Col 1 INS
```

second input for anymalign.py Output of the system :

Output files have basically the same format as the all-languages-in-one input file.
After running this process almost 1 hour we get this kind of output.

op-file (~/anymalign2.5) - gedit

Open

Save

Undo

op-file x

1	and	y	-	0.867202	0.898377	406794	
2	Commission	Comisión	-	0.904087	0.835330		
	115791						
3	not	no	-	0.783933	0.646228	112463	
4	Council	Consejo	-	0.953519	0.949122	102079	
5	President,	Señor Presidente,	-	0.524922			
	0.619656	82809					
6	report	informe	-	0.782963	0.909649	77161	
7	the	la	-	0.279681	0.444671	69651	
8	Parliament	Parlamento	-	0.852837	0.853843		
	64425						
9	or	o	-	0.876107	0.885031	62038	
10	the	de	-	0.245201	0.318387	61064	
11	two	dos	-	0.976055	0.957316	60735	
12	President,	Presidente,	-	0.380983	0.811992		
	60102						
13	in	en	-	0.386365	0.360306	55557	
14	Mrs	Sra.	-	0.934419	0.959019	47575	
15	of	de	-	0.313667	0.228273	43781	
16	three	tres	-	0.989738	0.980813	43400	
17	but	pero	-	0.666392	0.795720	41313	
18	is	es	-	0.272799	0.403877	39561	
19	behalf	nombre	-	0.979852	0.965779	39539	
20	Mr President,	Señor Presidente,	-	0.766008			
	0.291656	38976					
21	Member	Estados	-	0.635379	0.530203	37725	
22	proposal	propuesta	-	0.943853	0.882089		
	37218						
23	Union	Unión	-	0.608039	0.612916	36759	
24	between	entre	-	0.904811	0.809984	36263	
25	countries	países	-	0.903368	0.795440	34122	
26	debate	debate	-	0.935942	0.867540	33985	
27	my	mi	-	0.805529	0.808409	33013	
28	also	también	-	0.631924	0.610537	32831	

Plain Text Tab Width: 8 Ln 733052, Col 9 INS

Output of anymalign

Assignments :

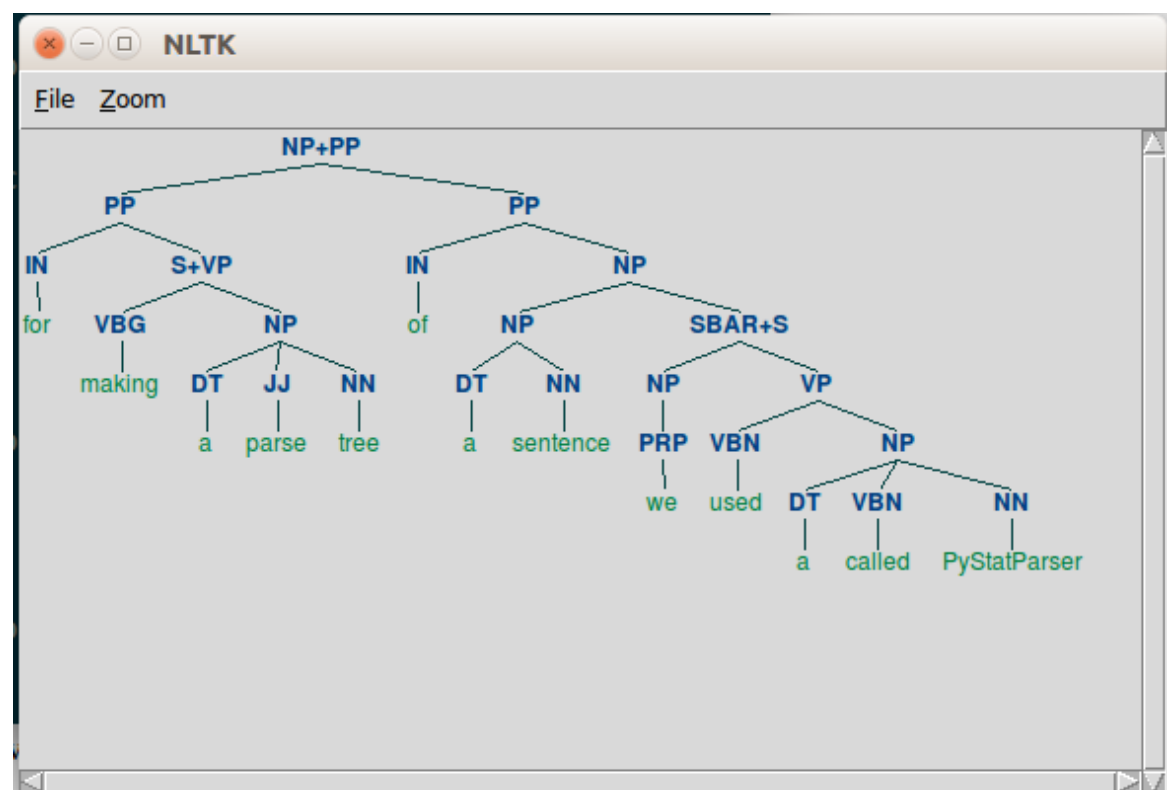
POS tagger : Our first job was POS tagger. POS tagger is build-in feature of the nltk module. For determine parts of speech of a sentence at first we need to split the

sentence word by word. So at first we take input the sentence as a string. Then split it word by word using word tokenize (built it function). Then using pos_tagger we got a list of words with parts of speech tag.

```
[('In', 'IN'), ('lexical', 'JJ'), ('analysis', 'NN'), (',', ','), ('tokenization', 'NN'), ('is', 'VBZ'), ('the', 'DT'), ('process', 'NN'), ('of', 'IN'), ('breaking', 'VBG'), ('a', 'DT'), ('stream', 'NN'), ('of', 'IN'), ('text', 'NN'), ('up', 'RB'), ('into', 'IN'), ('words', 'NNS'), (',', ','), ('phrases', 'NNS'), (',', ','), ('symbols', 'NNS'), (',', ','), ('or', 'CC'), ('other', 'JJ'), ('meaningful', 'JJ'), ('elements', 'NNS'), ('called', 'VBN'), ('tokens', 'NNS'), ('.', '.'), ('The', 'DT'), ('list', 'NN'), ('of', 'IN'), ('tokens', 'NNS'), ('becomes', 'NNS'), ('input', 'VBP'), ('for', 'IN'), ('further', 'JJ'), ('processing', 'NN'), ('such', 'JJ'), ('as', 'IN'), ('parsing', 'NN'), ('or', 'CC'), ('text', 'NN'), ('mining', 'NN'), ('.', '.')]
>>>
```

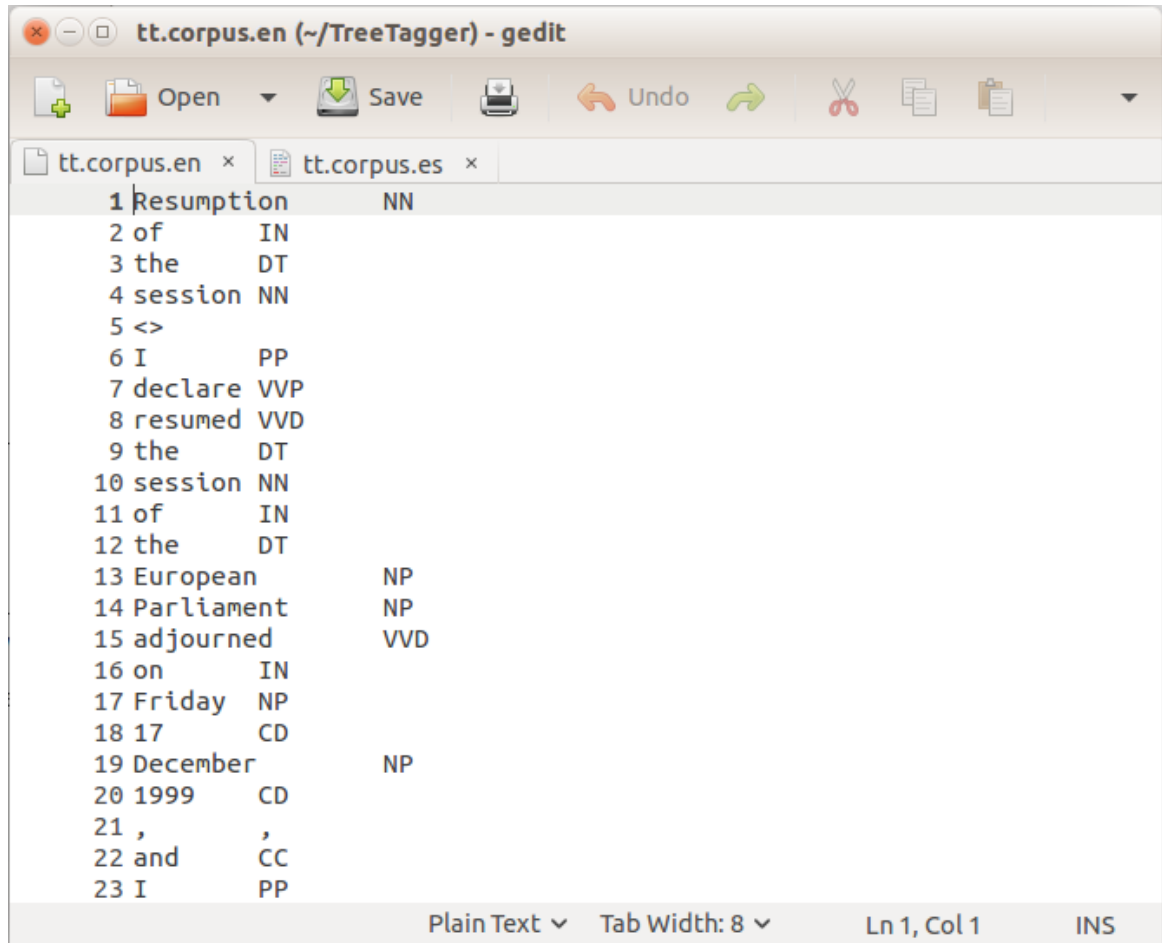
Output of pos tagger

Parse tree : For making a parse tree of a sentence we used a called PyStatParser. Which is available in www.github.com. Using this module we made a parse tree of a sentence.

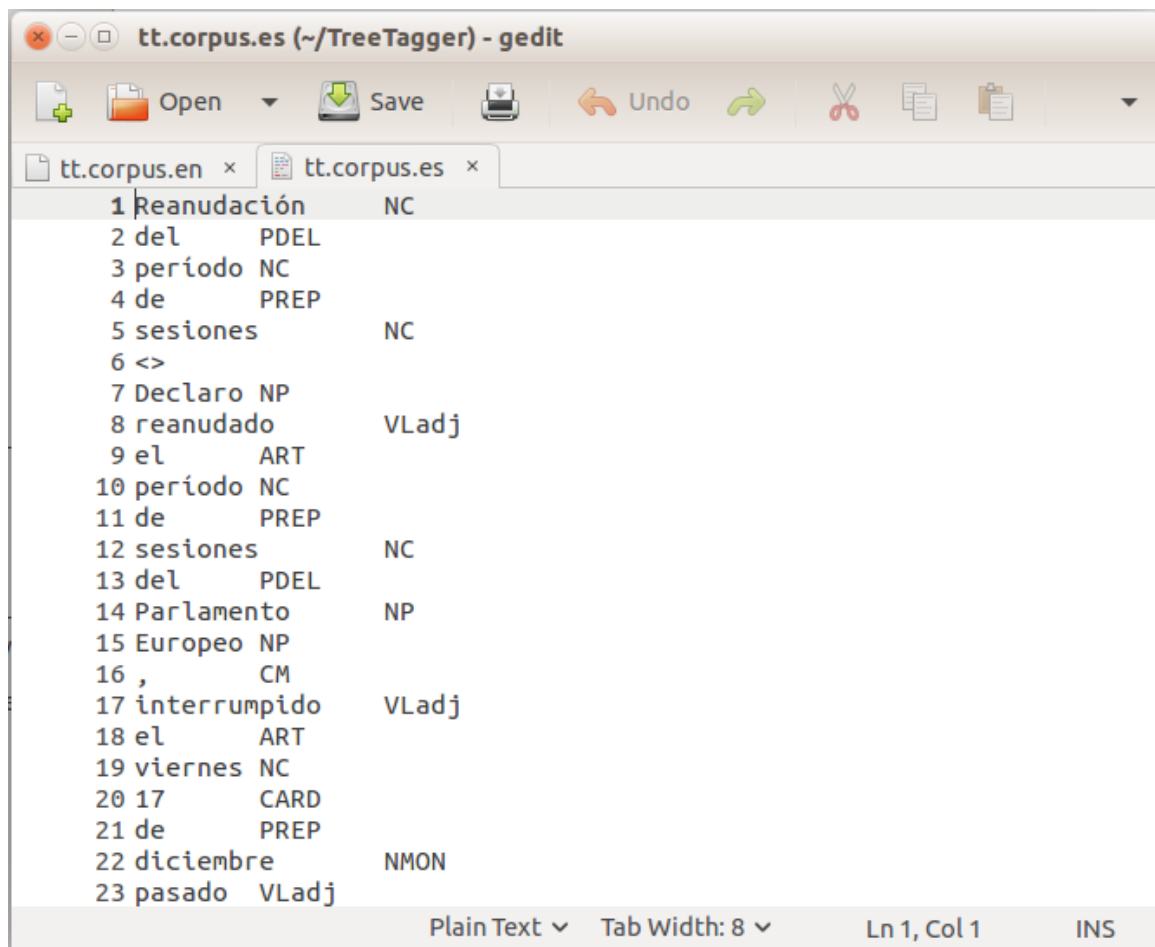


Parse tree

Tree tagger : Before starting word based translation we need to find out the parts of speech of English and Spanish words. Using TreeTagger we did this job.



Tree tagging of English



Tree tagging of Spanish

English to Spanish translator : In case of translator first job is making a file which contains two columns. The first column contains English words and second column contains Spanish words which is same meaning with left column's English word. Then our second job is a program to translate English word in Spanish word.

The program takes input as a string then split that string in word by word. The program search in the first column for each word. If the word is available in that list then picked up the Spanish word which is present right side of that English word. And then displays.

If the English word is not present in the list then it prints the English word (which is input) as output.

```
Python 2.7.10 Shell
File Edit Shell Debug Options Window Help
Python 2.7.10 (default, Oct 14 2015, 16:09:02)
[GCC 5.2.1 20151010] on linux2
Type "copyright", "credits" or "license()" for more information.
>>> ===== RESTART =====
>>>
Enter TEXT [english] : I have a book.
SPANISH :
que han un correctamente.
>>> |
```

Translation of English words into Spanish

```
Python 2.7.10 Shell
File Edit Shell Debug Options Window Help
Python 2.7.10 (default, Oct 14 2015, 16:09:02)
[GCC 5.2.1 20151010] on linux2
Type "copyright", "credits" or "license()" for more information.
>>> ===== RESTART =====
>>>
Enter TEXT [english] : I have a book.
SPANISH :
que han un correctamente.
>>> ===== RESTART =====
>>>
Enter TEXT [english] : hello world
hello
mundo
```

When English word not found

Advantage of Machine Translation :

The rate of machine translation is exponentially faster than that of human translation. The average human translator can translate around 2,000 words a day. Multiple translators can be assigned to a given project to increase that output, but it pales in comparison to machine translation. Machine translation can generate thousands of words each minute. One should note that the output of machine translation is not in its final useable form right away, but in certain scenarios it can be quite useful. Even when adding a post-editing step, machine translation takes a fraction of the time that human translation takes.

Cost : The cost of machine translation is also a mere fraction of the cost of human translation. While the major cost for standard translation projects is the cost of the

translators, the biggest cost for machine translation projects is in the post-editing process, which ends up saving the client a great deal.

Adaptation : Machine translation can memorize key terms and phrases that are used within a given industry. This leads to translations that are very consistent across the entire file, something that is more difficult to achieve when using multiple human translators.

Disadvantages of Machine Translation :

1. Accuracy is not offered by the machine translation on a consistent basis. You can get the gist of the draft or documents but machine translation only does word to word translation without comprehending the information which might have to be corrected manually later on.
2. Systematic and formal rules are followed by machine translation so it cannot concentrate on a context and solve ambiguity and neither makes use of experience or mental outlook like a human translator can.

Resource :

Reference :

1. <http://www.tutorialspoint.com/>
2. <http://www.nltk.org/>
3. <http://www.nltk.org/book/>
4. <https://github.com/>
5. <http://www.statmt.org/europarl/>
6. <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

7. <http://stackoverflow.com/>

Conclusion :

So I would like to conclude that benefits of machine translation all over the world. It is too difficult for a human to learn lots of languages. But a machine translation program can translate any language to another language. This technology makes our life familiar with unknown languages. A popular translator Google Translate is a free multilingual Statistical Machine Translation service provided by Google to translate text, speech, images, or real-time videos from one language into another. It offers a web interface, mobile interface (for Android & iOS) it can be used as a browser extension. It supports 90 languages at various levels.