

# **DESIGN DOCUMENT**

## **GROUP NO: 3**

DISHITA MALAV: 2017A7PS0164P

IESHAAN SAXENA: 2017A7PS0180P

KUSHAGRA RAINA: 2017A7PS0161P

NITIN VINAYAK AGRAWAL: 2017A4PS0415P

SHEFALI TRIPATHI: 2017A7PS0139P

**TABLE OF CONTENTS**

1.Introduction	3
2.Assumptions	3
3.Limitations	4
4.Innovation	5
5.Algorithms and Formulae Utilized	5
6.Description of Various Functions Implemented	6

## INTRODUCTION

We have implemented a Cross Lingual Document Translator, using Statistical Machine Translation model. The statistical model, IBM Model 1, has been trained for alignment and translation. Performance metrics such as cosine similarity and Pearson's correlation coefficient, have also been implemented in the translator.

## ASSUMPTIONS

1. Taking into consideration the large size of the data set, as a part of preprocessing the documents, we have:

- Removed all punctuation marks and replaced them with "".
- Converted all Uppercase characters to Lowercase for uniformity.
- Removed all numerals and replaced them with "".
- Removed a list of stop-words.

**List of Stop Words for English:** ['of', 'the', 'I', 'on', 'and', 'would', 'to', 'you', 'a', 'in', 'that', 'as', 'have', 'for', 'be', 'from', 'it', 'at', 'can', 'an', 'has', 'The', 'It', 'is', 'not', 'with', 'We', 'by', 'This', 'we', 'are', 'more', 'our', 'or', 'also', 'these', 'but', 'must']

**List of Stop Words for Dutch:** ['van', 'de', 'ik', 'Aan', 'en', 'Zou', 'naar', 'u', 'een', 'in', 'dat', 'als', 'hebben', 'voor', 'worden',

'van', 'het', 'Bij', 'kan', 'een', 'heeft',  
 'De', 'Het', 'is', 'niet', 'met', 'Wij',  
 'door', 'Deze', 'wij', 'zijn', 'meer', 'onze',  
 'of', 'ook', 'deze', 'maar', 'moet']

2. An SQL Database is created to store translational probabilities for mapping English to Dutch (Foreign Language). The Database Columns are:

ENG_WORD	DUT_WORD	PROBABILITY
----------	----------	-------------

Primary Key: The composite key of ENG\_WORD and DUT\_WORD.

The Clustered Index is based on the Primary Key.

### **LIMITATIONS**

Space and Time were the two major limitations in the project.

1. Due to the large size of the dataset, we found that storing the translational probabilities, counts and other metadata related to the datasets, in data structures implemented in python (like a dictionary) was highly space consuming. Hence, we stored them in files and SQL Databases.
2. The EM Algorithm used in training the IBM Model also takes a significant amount of time. In order to reduce the time taken, we used SQL Database in which reduced the time taken in seeks. Stop-words were also removed in order to reduce the amount of time taken in processing these commonly occurring words.

## **INNOVATION:**

The large size of the dataset increases the numbers of seeks, in order to make the model most time and space efficient we used an SQL Database to store the translational probabilities  $t(e|f)$ . We also used a clustered index in the database to further reduce the retrieval time.

## **ALGORITHMS AND FORMULAE UTILIZED:**

### 1.Expectation-Maximization Algorithm:

```

initialize  $t(e|f)$  uniformly
do until convergence
    set  $\text{count}(e|f)$  to 0 for all  $e, f$ 
    set  $\text{total}(f)$  to 0 for all  $f$ 
    for all sentence pairs  $(e\_s, f\_s)$ 
        set  $\text{total\_s}(e) = 0$  for all  $e$ 
        for all words  $e$  in  $e\_s$ 
            for all words  $f$  in  $f\_s$ 
                 $\text{total\_s}(e) += t(e|f)$ 
        for all words  $e$  in  $e\_s$ 
            for all words  $f$  in  $f\_s$ 
                 $\text{count}(e|f) += t(e|f) / \text{total\_s}(e)$ 
                 $\text{total}(f) += t(e|f) / \text{total\_s}(e)$ 
    for all  $f$ 
        for all  $e$ 
             $t(e|f) = \text{count}(e|f) / \text{total}(f)$ 

```

### 2.Pearson Coefficient Calculation:

$$sim(a, b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2}}$$

3. Cosine Similarity Calculation:

$$\cos(\vec{q}, \vec{d}) = \vec{q} \bullet \vec{d} = \sum_{i=1}^{|V|} q_i d_i$$

### DESCRIPTION OF VARIOUS MODULES IMPLEMENTED:

#### 1. Function: readfiles

Description: Function opens the source files for English and Dutch languages, loads all their sentences to return variables 'eng' and 'dutch' respectively.

Input Variables:

*Eng\_File*: Name of the file containing English Corpus.

*Dutch\_File*: Name of the file containing Dutch Corpus.

*no\_of\_sentences*: Total number of sentences in the Corpus'.

Return Variables:

*eng*: List of all English sentences.

*dutch*: List of all Dutch sentences.

#### 2. Function: remove\_punc

Description: Function calls `remove_stopwords` and `remove_stuff` to remove stopwords and punctuation, sentence by sentence, respectively.

Input Variables:

*l*: List of English or Dutch sentences

*lang*: Name of the language passed.

Return Variables:

*l*: Containing all the sentences without the stop-words and punctuation.

### 3. Function: `remove_stopwords`

Description: Removes stop-words in a sentence.

Input Variables:

*eng\_sw*: List of English stop-words.

*dut\_sw*: List of Dutch stop-words.

*lang*: Name of the language the input sentence belongs to.

*l*: A single English or Dutch sentence.

Return Variables:

*l*: Sentence without stop-words.

### 4. Function: `remove_stuff`

Description: Removes all punctuation in a sentence.

Input Variables:

*l*: A single English or Dutch sentence.

Return Variables:

*l*: Sentence with punctuations removed.

#### 5.Function: assign\_line\_no

Description: Creates an inverted index where each word is matched to all the lines it is present in.

Input Variables:

*doc*: Contains list of the cleaned sentences of either Dutch or English (i.e. without the punctuations or stop-words).

Return Variables:

*dict\_lo*: Contains the inverted index (in the form of a dictionary) of all the words in the document.

#### 6.Function: initialize

Description: Creates Dutch-English word pairs then initializes their probabilities. It then writes then in an SQL Database.

Input Variables:

*foreign\_no\_of\_words* : Total number of Dutch words.

*foreign\_l*: All Dutch words.

*english\_l*: All English words.

Return Variables:



Returns True.

#### 7.Function: finding\_probabilities

Description: Updates count variable for each Dutch-English word pair.

Input Variables:

*dutch\_sentences*: List of clean Dutch sentences retrieved

*eng\_sentences*: List of english sentences retrieved

*no\_of\_sentences*: Total number of sentences

*total*: Normalized fractional counts

Return Variables:

*total*: Normalized fractional counts of every dutch-english word pair.

#### 8.Function: running\_function

Description: Function repeatedly calls finding\_probabilities to update count for every Dutch-English word pair.

Input Variables:

*foreign\_l*: List of all Dutch words.

*english\_l* : List of all English words.

*dutch\_sentences* : List of all clean Dutch sentences retrieved

*eng\_sentences* : List of all clean English sentences retrieved

*no\_of\_iterations* : Number of times the count is to be updated.

Return Variables:

None.

#### 9. Function: retrieve\_max

Description: Function to retrieve the translation with maximum probabilities for each word to be translated, i.e. to obtain the most probable translation for each word

Input Variables:

*num\_dict\_dutch* : Dictionary of Inverted index of all Dutch words

*num\_dict\_eng* : Variable which stores line number of English words.

Return Variables:

*translation\_etof*: Translation from English to Dutch.

*translation\_ftoe*: Translation from English to Dutch.

#### 10. Function: pearson\_coefficient

Description: Calculates the Pearson coefficient. It calculates normalized weight of individual sentences by calling `maintain_count()`.

Input Variables:

*dutchWord\_Line\_no*: Dictionary of Dutch words mapping to individual line numbers.

*dut\_cleaned*: Dictionary of cleaned Dutch lines.

*result\_dut*: List of lines produced by our model.

Return Variable:

Returns Pearson Coefficient between two Documents

#### 11. Function: `cosine_similarity`

Description: Function calculates cosine similarity. It calculates normalized weight of individual sentences by calling `maintain_count()`.

Input Variables:

*dutchWord\_Line\_no*: Dictionary of Dutch words mapping to individual line numbers.

*dut\_cleaned*: List of cleaned Dutch lines.

*result\_dut*: List of lines produced by our model.

Return Variables:

*Total\_sim*: Cosine similarity between two dutch lines.

#### 12. Function `maintain_normalized_tf`:

**Description:** Function to find normalized  $T_f$  weights of a Dutch sentence with a sentence produced by our model.

**Input Variables:**

*doc\_dict\_line\_no*: Single Dutch sentence.

*data\_clean*: List that contains Dutch sentences produces by the IBM model.

**Return Variables:**

*og\_dutch\_tf*: Normalized  $T_f$  weights 2 sentences at a time.

13. **Function:** produce\_sentence

**Description:** Translates an English sentence to Dutch and adds the punctuation to form a complete sentence.

**Input Variables:**

*eng\_sentence*: A single English sentence.

*translated\_dict*: Dictionary of Dutch translation of all English words.

*doc\_lang*: Variable gives user option to either translate English to Dutch or vice-versa. Default is English to Dutch.

**Return Variables:**

*s*: Dutch translation of an English sentence.

#### 14. Function: translate\_doc

Description: Function calls produce\_sentence repeatedly to translate the entire English document to Dutch.

##### Input Variables:

*eng\_doc*: Variable contains either Dutch or English list of sentences.

*translated\_dict*: Variable contains Dutch translation of all English words.

*rewrite\_file*: The file in which all the translated words are stored.

*doc\_lang*: Variable gives user option to either translate English to Dutch or vice-versa. Default is English to Dutch.

##### Return Variables:

*result\_doc*: List of all translated sentences.

#### 15. Function: delete\_files

Description: Deleting previously present files responsible for maintaining translational probabilities and count.