

# Documentation:

## Module: 1.py

### **load** ( *path* )

Loads parallel corpus which is in json format

**Parameters:**    **path : str**

The file name (string) from where the matrix will be loaded.

**Returns:**        **data : 2D list**

2D list with Dictionary of both languages as keys

### **get\_tokens** ( *data* )

Splits sentences into words

**Parameters:**    **Data: 2D list**

2D list with Dictionary of both languages as keys

**Returns:**        **Tokens: List**

List contains sentence tokenized into words.

**iteration** ( *data, tokens, total, prev\_prob*)

Loads parallel corpus which is in json format

**Parameters:**    **Data: 2D list**

2D list with Dictionary of both languages as keys

**Tokens: List**

Of words by language

**total:**

counts of the destination words, weighted according to their translation probabilities  $t(e|s)$

**Prev\_prob:**

Translation probabilities from previous iteration

**Returns:**    **curr\_prob: List**

**distance** ( *table\_1,table\_2*)

Finds RMS between both the tables.

**Parameters:**    **table\_1**

**table\_2**

**Returns:**    **sqrt(res): Integer**



### **check\_convergence (prev\_prob, curr\_probs, epsilon)**

Checks when to stop iterating

|                    |                                                                |
|--------------------|----------------------------------------------------------------|
| <b>Parameters:</b> | <b>Prev_prob, curr_probs: 2d Dict</b>                          |
|                    | Translation probability matrices before and after an iteration |
|                    | <b>Epsilon: int</b>                                            |
|                    | Boundary for convergence.                                      |
| <b>Returns:</b>    | <b>Bool</b>                                                    |
|                    | True if $\text{delta} < \text{epsilon}$                        |

### **train ( data, epsilon)**

Trains the dataset

|                    |                                                                |
|--------------------|----------------------------------------------------------------|
| <b>Parameters:</b> | <b>Data: 2D list</b>                                           |
|                    | 2D list with Dictionary of both languages as keys              |
|                    | <b>epsilon:int</b>                                             |
| <b>Returns:</b>    | <b>Curr_probs</b>                                              |
|                    | Translation probability matrices before and after an iteration |

## **Module: 2.py**

### **load ( path )**

Loads parallel corpus which is in json format

|                    |                   |
|--------------------|-------------------|
| <b>Parameters:</b> | <b>path : str</b> |
|--------------------|-------------------|

The file name (string) from where the matrix will be loaded.

**Returns:**      **data : 2D list**

2D list with Dictionary of both languages as keys

**tokData**      ( *data*, *target='fr'* )

Tokenizes the data and loads it appropriately as 2D list containing Aligned Sent

**Parameters:**      **data : 2D list**

2D list with Dictionary of both languages as keys

**target: str**

Language other than english used in corpus loaded.

**Returns:**      **data : 2D list**

2D list of tokenized words

**main**      ( )

The main function that carries out all the major steps to produce the results.

**No parameters or return values.**

## **Module: 3.py**

**load**      ( *path* )

Loads parallel corpus which is in json format

**Parameters:**      **path : str**

The file name (string) from where the matrix will be loaded.

**Returns:**      **data : 2D list**

2D list with Dictionary of both languages as keys

**phrase\_bases\_extraction** ( *filename*, *foreign* )

Tokenizes the data and loads it appropriately as 2D list containing Aligned Sent

**Parameters:**      **filename : str**

Name of the dataset file

**foreign: str**

Language other than english used in corpus loaded.

**main** ( )

The main function that carries out all the major steps to produce the results.

**No parameters or return values.**