

Pristine Sentence Translation: A New Approach to a Timeless Problem

Meenu Ahluwalia, Brian Coari, and Ben Brock

¹Master of Science in Data Science

Southern Methodist University

Dallas, Texas USA

{mahuwalia, bcoari, bbrock}@smu.edu

Abstract. Translating text from one language to another is a continuous technological challenge. Although many technologies, such as Google Translate, have used machine learning and neural networks to close the translation gap, there are still many translation problems to be solved. Issues such as multiple word meanings, proper sentence structure, slang, colloquialisms, and determining the literal meaning of words vs contextual intent of those words are areas where we sometimes still see Google Translate struggle. In this paper we explore an original strategy that provides a solution to these translation issues, demonstrate a proof-of-concept of the solution, and examine the feasibility of a large-scale solution. For our translation solution we populated a database with translations of entire sentences from one language to another, instead of the words in a sentence. Since a sentence represents an entire thought instead of an assembly of words, the translation did not suffer from the issues that plague Google Translate. We also used Natural Language Processing (NLP) and predictive modeling in order to find sentences close to the sentence requested, which provides the user examples of common grammatically-correct sentences. With these approaches we were able to translate sentences that seemed impossible using traditional translation methods.

1 Introduction

The ability to easily communicate with people in another language is one of the most powerful and satisfying experiences in life. Technology has come a long way from the discovery of the Rosetta Stone in 1799, which allowed us to translate Egyptian hieroglyphics to ancient Greek in a mere 23 years. In the modern day, tools such as Google translate can be used in real time to convert between languages and allow people to connect from different cultures [6]. The latest iterations of Google Translate even use machine learning and neural nets to parse more than just single words, delivering a more satisfying user experience [7].

As far as we have come, however, the areas where we struggle are still painfully obvious. While Google Translate will usually allow you to find a bathroom and order off a menu, the intricacies and complexities of a normal, native

conversation still can cause a non-fluent speaker issues. For example, if an American coworker mentions to a Brazilian coworker about their performance on a project with "You hit one out of the park.", the Brazilian coworker could translate the words, but without familiarity with the context of a baseball game, the Brazilian would be confused and would have to ask for clarification if it was possible. It would be even harder if the Brazilian was reading a book in English with a colloquialism, since there would be no human to ask for help.

If you consider these kinds of issues from a high level they might seem unsolvable. How can you train a translation tool to look at the meaning behind sentences using on the words provided? We think we have a possible answer: an original concept we are calling Pristine Sentence Translations (PSTs). The concept of PST is that instead of translating words or phrases using neural nets and machine learning, we simply store an entire sentence in a database, and we have entire sentences in other languages that represent the meaning of that sentence.

For example, using the example above we would have an entry for the English sentence "You hit one out of the park", and we would have an entry for a Portuguese sentence mapped to that English sentence that says "Você foi ótima" which translates in English to the meaning behind the phrase: "You did great". For another example, in Portuguese there's a sentence "Eu adoro Cafuné" Google Translate does not have a translation for "Cafuné", because it's a complicated word which loosely means "the act of running fingers through hair". Our program's goal is to return an English translation "I love the feeling of fingers running through my hair" when asked to translate "Eu adoro Cafuné" into English. Using this method there is no sentence or concept we will not be able to translate into another language given enough time and resources.

One main issue with the approach outlined above is that if we do not have an exact match for the sentence, our method return nothing. so if we tried to translate "You really hit one out of the park" from English into Portuguese we would not get any results. We decided to address this concern using Natural Language Processing (NLP) to filter out the noise in a sentence, and then use Predictive Modeling in order to find the sentence "most like" the input sentence. Using this method, "You really hit one out of the park" would ideally map most closely to "You hit one out of the park", and return the same translation: "Você foi ótima". The front-end will indicate that the translation is not for the original input sentence, instead it will indicate that it is "Showing Results for: You hit one out of the park."

Due to the strictly educational and academic nature of the project, we are not attempting to provide a full translation solution. We will limit our translations to English, Portuguese, and Hindi, and we will only provide translations for a few hundred phrases. This will be sufficient to demonstrate the appeal and power of this technique, and we will show how this solution could grow into a complete, living solution using crowdsourcing and time.

1.1 Existing Tools and Methods

The concept of translating one language to another is one of the most prevalent topics in data science, so it was all too easy to find existing tools and methods. This first example we would like to cover is one of the stronger forces in translation today: Google Translate. In 2016 Google started using Neural Nets to perform translations and saw vast reduction in errors. [18]. Even though this method was vastly superior to Google's previous system it is still dependent on the individual words in a sentence, and still subject to the same issues stated above with translations, since it would not handle colloquialisms or words that have no direct equivalent in the destination language.

The next method we found was much closer to our intent. The author developed a Neural Machine Translation System using Keras, and used a "closeness of match" system to translate between one language and another. [16] This approach is almost perfect for what we need, with the exception of only returning one sentence in the intended language. We would like an approach that allows us to find more than one sentence ID from the destination language. PST will use a similar approach, but aims to find a matching sentence ID for the input sentence, then lookup all sentences that are translations for that sentence in the requested language. More details will be provided in the "Solution Approach" section.

2 Solution Approach

Our approach to solving the issues around language translations takes five main steps. Each step is discrete and can be run and unit-tested independently. The "Full Demo" section will go over what the steps look like and how they relate to each other in a full iteration of the program.

First, we will form a pristine database of known, correct, contextual translations to sentences that focus on translating the intent of sentences rather than the word in that sentence. There is more detail on this step in the "Data Collection" section.

Second, we will take input from a user on the requested translation and the languages involved in the translation. This is covered in more detail in the "Full Demo" section.

Third, we plan to use a sequence to sequence model to find as close a match as possible to the input sentence and its associated translation. This section is still in progress in this iteration of the document, we will use it to verify our approach in step four which meets functional expectations for the project. More details of that can be found in the "Encoder-Decoder Long Short-Term Memory (LSTM) Networks" section.

Fourth, we will measure the closeness of our guessed sentence to our input sentence using Euclidean Distance and Cosine Similarity and return the result of the method that we determine is the best performing indicator, along with all sentences that relate to the matched sentence to the front end. This is covered

in more detail in the "Euclidean Distance and Cosine Similarity" and "Results and Analysis" sections.

Fifth and last, we will display our results to the user through the front end. This is covered in more detail in the "Full Demo" section.

2.1 Database Design

For our database design we wanted to leverage a cloud-based storage system so that we could both reference the same live database, as well as making the database updateable in realtime in case we want to utilize a realtime update strategy in the future. We chose Google Cloud Platform (GCP) as a host for our database due to its ease of setup and robust set of no-cost features.

Our GCP database consists of two tables:

1. Sentences: Stores a unique Sentence ID for a combination of text (max 500 characters) and a language key (EN=English, PT=Portuguese, HI=Hindi). There is also a place for a suggested replacement for possible future enhancements, but this is not used in our current version of the program.
2. Translations: A table that stores the Sentence IDs that are translations of each other, linking the sentences together. These Sentence IDs are foreign keys that link to Sentence IDs in the Sentece table

The use of the Translations lookup table allows us to associate any number of sentences to each other, so we can have more than one translation for a sentence, and we can maintain translations to sentences only that match up together. In the fuure we might have sentences in English that only have translations in Hindi, but not Portuguese, for example.

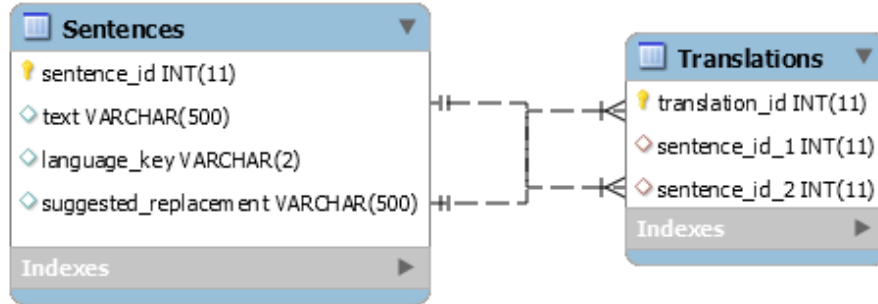


Fig. 1: Database Diagram

One of the problems we wanted to solve early was to be able to test our model offline, in case a network connect occurred we still wanted access to our model. To facilitate this we extracted the data in the database into dataframes through sqlalchemy [15] and exported the data to .pkl files. We then load the

.pkl files whenever we need the data for a model, which makes us immune to network interruptions. Examples of the SQL Queries used to extract the data can be found in Appendix A.

2.2 Data Collection

Covered in Ethical Considerations as well, I'll revisit this to make it unique.

For this project the initial set of 200 English Sentences data were gathered from three travel websites. [11,12,13]. Travel sites were chosen in order to get a good spread of what sentences people might find useful. Translations were then done using a combination of Google Translate [14] and a human expertise in Portuguese and Hindi in order to determine if the Google Translation made sense. If the sentence did not make sense as Google Translated it, we entered a more sensible translation. The original list of sentences and translations can be found in our public Github Repository here:

https://github.com/coarib/SMU_Masters_PST/blob/master/data/processed/PST_workbook.xlsx

3 Predictive Modeling

Automatic or machine translation is one of the most challenging AI tasks given the fluidity of human language. Classically, rule-based systems were used for this task, which were replaced in the 1990s with statistical methods. More recently, deep neural network models achieve state-of-the-art results in a field that is aptly named neural machine translation. [16]

Sequence to Sequence (often abbreviated to seq2seq) models are a special class of Recurrent Neural Network architectures typically used (but not restricted) to solve complex Language related problems like Machine Translation, Question Answering, creating Chat-bots, Text Summarization, etc. Our aim is to translate given sentence from one language to another. We will target sentence translations to and from English, Portuguese and Hindi languages only. Use of seq2seq (or Encoder-Decoder) architecture is appropriate in this case as the length of the input sequence does not has the same length as the output data

To summarize our model, the Encoder simply takes the input data, and train on it then it passes the last state of its recurrent layer as an initial state to the first recurrent layer of the decoder part. The Decoder takes the last state of encoder's last recurrent layer and uses it as an initial state to its first recurrent layer , the input of the decoder is the sequences that we want to get. We will use Keras API with Tensorflow backend to build our model.

3.1 Data Preparation

Before we start building the model, we need to clean up the text data (i.e. the sentences). We will remove all punctuation characters, normalize the case to

lowercase, normalize all Unicode characters to ASCII and remove any tokens that are not alphabetic. To build the model, we need to map words to Integers. We will use Keras Tokenize class for this. The Tokenizer must be constructed and then fit on either raw text documents or integer encoded text documents. Once fit, the Tokenizer provides four attributes that you can use to understand about your text., viz.,

1. word-counts: A dictionary of words and their counts
2. word-docs: A dictionary of words and how many documents each appeared in.
3. word-index: A dictionary of words and their uniquely assigned integers.
4. document-count: An integer count of the total number of documents that were used to fit the Tokenizer.

We will also compute the vocabulary sizes and the length of maximum sequence for both the languages. We need to encode each input and output sentences to integers and pad them to the maximum phrase length to make all sentences of the same length. This is because we will use word embedding for the input sentence and one hot encoding for the output. In one hot encoding, a document is represented as a sequence of integer values, where each word in the document is represented as a unique integer. One hot encoding is needed because the model will predict the probability of each word in the vocabulary as output.

3.2 Euclidean Distance and Cosine Similarity

For our comparison on closeness of strings we will leverage two methods: Euclidean Distance and Cosine Similarity. Euclidean Distance measures the distance between two points in a space, and therefore the minimum distance is the closest match. Cosine Similarity is the similarity between two vectors in a space, and therefore the maximum similarity is the closest match. In this section we will describe how these methods calculate their scores, and in the "Results" section we will show the result of the comparison and choose our preferred match.

Euclidean distance, also known as L2 distance or the Pythagorean metric, is a basic measure of distance in a space. given a number of dimensions in a space, the Euclidean Distance calculates the length of a line segment connecting those points. This distance is calculated by going through all of the dimensions, adding up the square of the difference between the points of each value for that dimension, then taking the square root of that sum. This method is called the Pythagorean Formula. [19]. This method is easy to understand, but it seems like it might struggle if the data are not of similar shapes, as is the case with our sentence sample.

Cosine Similarity tries to measure similarity in a fundamentally different way from Euclidean Distance. Instead of stepping through and comparing data point-by-point, Cosine Similarity instead steps through each sample and tries to build a directional vector for the shape of the data, then measures the

cosine of the angle between the two vectors. In the case of text analysis the cosine similarity uses the frequency of words in a sentence order to determine a directional vector. Cosine Similarity ignores the words that are not in common between two sentences, and focuses only on words that are in common, which avoids penalizing sentences that are very long and simply focuses on finding the closest possible match. Therefore cosine similarity is very resilient to data that is of different shapes, since the shape of the data would not necessarily affect the direction of the vector, instead the sentences with the closest ratios of like words would be considered a match. [21] This resilience to differently-shaped data is promising for our purposes, and seems to be a good fit for comparing sentences, and we will see in the "Results" section how it performs.

3.3 Encoder-Decoder Long Short-Term Memory (LSTM) Networks

A typical seq2seq model consists of 2 major components

1. Encoder
2. Decoder

Both these components are essentially two different Recurrent Neural Network models combined into one giant network.

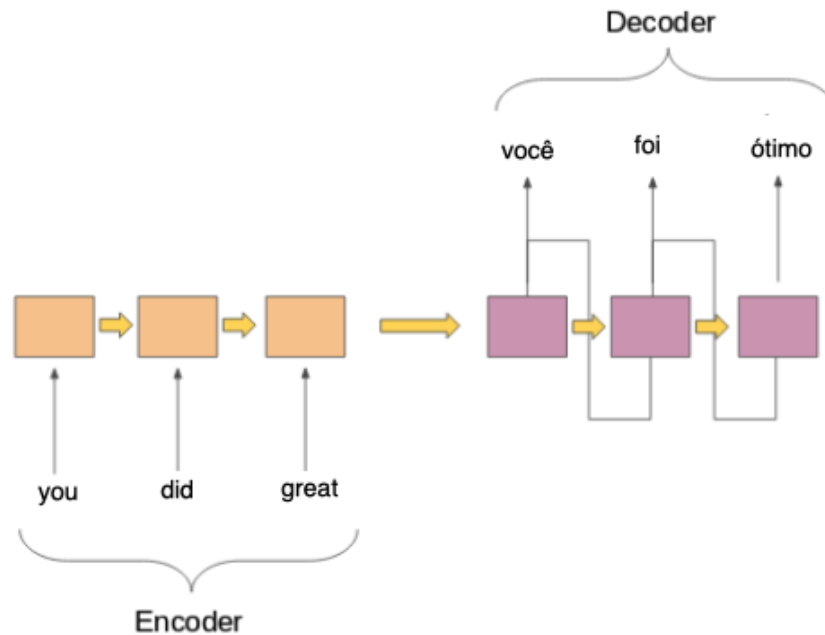


Fig. 2: Sequence to Sequence Modell
[4]

We will explain the Encoder and Decoder model in more detail.

Let's say we are trying to convert the following sentence from English to Portuguese.

Input sentence (English) - i have lost my passport

Output sentence (Portuguese) - eu perdi meu passaporte

A sentence can be seen as a sequence of words or characters. We will split the sentence by words. So, for the above example in English, there are 5 words which are fed to the encoder as shown in the figure below. The input is referred to as X and X_i is the input sequence at time step i . So we have the following input. $X_1 = i$, $X_2 = \text{have}$, $X_3 = \text{lost}$, $X_4 = \text{my}$, $X_5 = \text{passport}$. Each X_i is mapped to a fixed-length vector using the built-in embedding layer of Keras API.

The LSTM will read this sentence word by word in 5 time steps as shown in the figure.

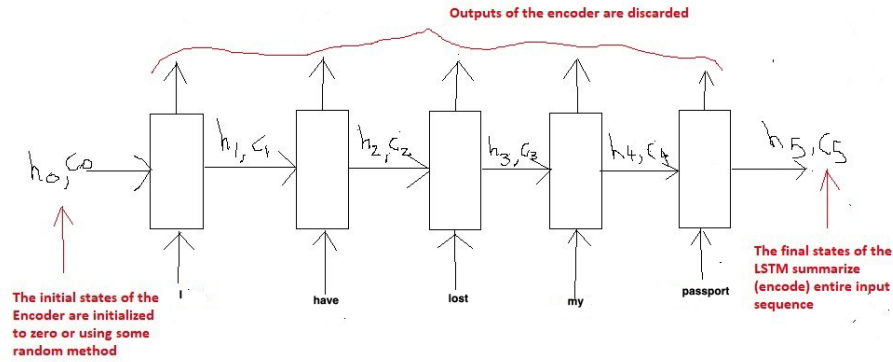


Fig. 3: Encoder LSTM
[5]

h_i and c_i in the figure above represent the internal state, viz., the hidden state and the cell state of the Encoder. In simple terms, they remember what LSTM has read till now. For example, h_3, c_3 vectors will remember that the network has read "I have lost" till now. Basically its the summary of information till time step 3 which is stored in the vectors h_3 and c_3 (thus called the states at time step 3). So, h_5, c_5 will contain the summary of the entire sentence. These states coming out of the last time step are also called as the "Thought vectors" as they summarize the entire sequence in a vector form. We initialize h_0, c_0 to zero as the model has not started to read the input.

Y_i is the output of the LSTM at each step. We discard the outputs of the encoder and only preserve the internal states as the model has nothing to output unless it has read the entire English sentence.

Next, we define the Decoder. Unlike the Encoder LSTM which has the same role to play in both the training phase as well as in the inference phase, the Decoder LSTM has a slightly different role to play in both of these phases. Recall that given the input sentence "i have lost my passport", the goal of the decoder is to output "eu perdi meu passaporte".

The initial states (h_0, c_0) of the Decoder are set to the final states of the Encoder. This intuitively means that the decoder is trained to start generating the output sequence depending on the information encoded by the encoder.

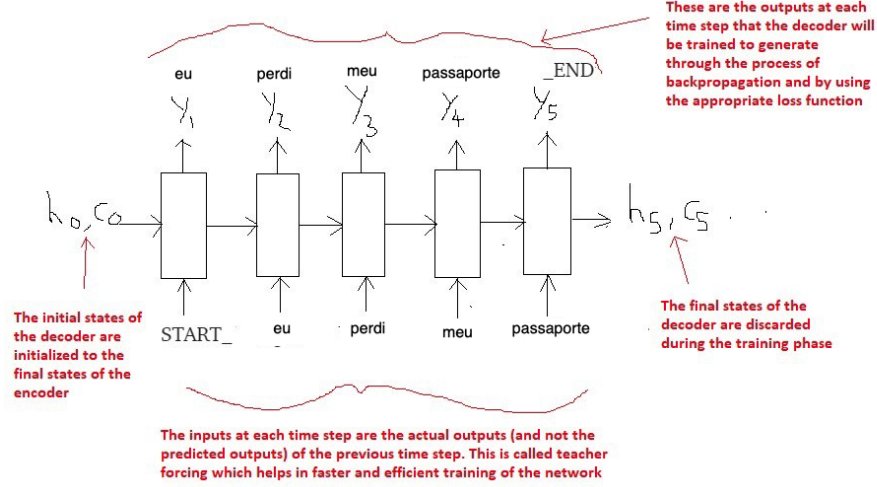


Fig. 4: Decoder LSTM
[6]

3.4 Building the Neural Translation Model

We will split our dataset into train and test set for model training and evaluation, respectively. Our seq2seq model is defined as the following

1. For the encoder, we will use an embedding layer and an LSTM layer
2. For the decoder, we will use another LSTM layer followed by a dense layer

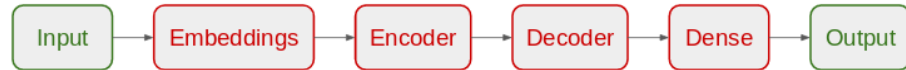


Fig. 5: Model Architecture
[7]

4 Full Demo

The front end for the Pristine Language Translation (PST) site is a simple proof-of-concept page with three boxes for input: The language and the text the user is translating from, and the language code the user is translating to. The user must fill in all three boxes with valid input before hitting the "Translate" button:

Pristine Sentence Translations

Enter the language you are translating from, the language you are translating to, and the requested translation. English = EN, Portuguese = PT, Hindi = HI

From Language:	<input type="text" value="EN"/>	Requested Translation: <div><div>You did good.</div></div>
To Language:	<input type="text" value="PT"/>	
<div>Translate</div>		

Fig. 6: Language Input

After the user has entered the data and hit the "Translate" button the information is passed to our python NLP processing. First our model finds the closest linguistic match to the sentence entered to a sentence for which the translation is known. Then we have a separate pos-processing step to calculate how close of a match the input sentence is to the matched sentence, and we retrieve that translations of that sentence for the requested language:

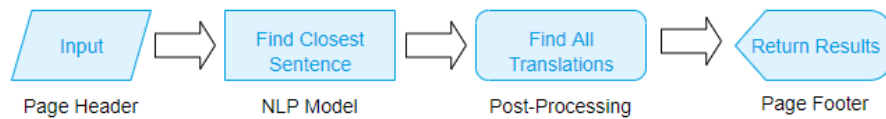


Fig. 7: Process Map

The results are then returned to the screen for display. Note that there can be more than one valid translation for an input sentence, provided our database has the mappings:

Pristine Sentence Translations

Enter the language you are translating from, the language you are translating to, and the requested translation. English = EN, Portuguese = PT, Hindi = HI

From Language:

To Language:

Requested Translation:

You did good.

Translate

Closest Available Match (63%):

you did great

Translations For Matched Sentence:

you did great

você fez muito bem

você foi ótimo

Fig. 8: Display Translations

If the user wants to try another sentence the user can change any of the three inputs above and hit the "Translate" button again to repeat the process.

4.1 Results and Analysis

[20]

4.2 Conclusions

[TODO: Add section here]

5 Ethical Considerations

The use of Natural Language Processing (NLP) in regards to language translations raises several ethical issues. Like many other fields in data science, practitioners of NLP must worry about the core questions around gathering data, as well as a second, more-specific ethical concern around their topic (in the case of this paper, translation of language). First we will address the core questions around ethical data science behaviors, then the more specific issue around our translation project.

In order to find a framework in which to address ethical considerations, it would be helpful to have a template, or some core questions to answer. Margot Mieske proposes in her paper "A Quantitative Study of Data in the NLP community" five key questions every NLP programmer must answer: [8]

- Has data been collected?
- How was this data collected and processed?
- Was previously available data used/extended – which one?
- Is a link or a contact given?
- Where does it point (private page, research institute, official repository)?

For this project the initial set of 200 English Sentences data were gathered from three travel websites. [11,12,13]. Travel sites were chosen in order to get a good spread of what sentences people might find useful. Translations were then done using a combination of Google Translate [14] and a human expertise in Portuguese and Hindi in order to determine if the Google Translation made sense. If the sentence did not make sense as Google Translated it, we entered a more sensible translation. The original list of sentences and translations can be found in our public Github Repository here:

https://github.com/coarib/SMU_Masters_PST/blob/master/data/processed/PST_workbook.xlsx

For the specifics around ethical considerations for NLP we turn to Jochen L. Leidner and Vassilis Plachouras and their paper "Ethical by Design: Ethics Best Practices for Natural Language Processing" [15]. Leidner and Plachouras propose that since NLP pertains to human language and touches every part of human life it has a specific ethics dimension, therefore automation and errors also become ethical topics. We need to make sure our translations and base sentences are unbiased and fair, without discrimination based on age, race, or gender. This is why we chose to target travel sites that hopefully won't favor certain demographics. At the very least we are transparent about from where we pulled the sentences.

This is also why we carefully scanned each translation for fairness to make sure our data was reviewed before putting it out to the public, and why we keep a tight control over the sentences that might be suggested for a user. As it is, this proof-of-concept system is not something I'd like to put out into the world for general use. With only 200 sentences the risks around mistranslating something are too high, there is the potential for widespread confusion if our suggestions

aren't close enough to the input request. Before going public we would need to follow Leidner and Plachouras's advice and establish an ethics review board that establish a process for reviewing and implementing new translation requests and helping monitor issues around translations on the site.

6 Conclusions and Other Work

[TODO: Add detail here]

Pristine Sentence Translations model is only built for about 200 sentences which could be translated from/to English, Portuguese and Hindi languages. This model could be expanded by adding more data and by incorporating more languages for translations.

One could try dropout and other forms of regularization techniques to mitigate over-fitting, or perform with hyperparameter tuning. Play with learning rate, batch-size, number of epochs etc.

It would be interesting to see how the model would perform when built using Attention.

References

1. Google's new translation software is powered by brainlike artificial intelligence (2016, September 27), Available at: https://www.sciencemag.org/news/2016/09/google-s-new-translation-software-powered-brainlike-artificial-intelligence?r3f_986=https://www.google.com/. Last accessed 4 Feb 2019
2. Found in translation: More accurate, fluent sentences in Google Translate (2016, November 15), Available at: <https://www.blog.google/products/translate/found-translation-more-accurate-fluent-sentences-google-translate>. Last accessed 4 Feb 2019
3. A ten-minute introduction to sequence-to-sequence learning in Keras, Available at: <https://blog.keras.io/a-ten-minute-introduction-to-sequence-to-sequence-learning-in-keras.html>. Last accessed 4 Feb 2019
4. Figure 1. Available at: https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/01/enc_dec_2.png. Last accessed 22 Mar 2019
5. Figure 2. Available at: https://cdn-images-1.medium.com/max/1600/1*37tR0olA8uW7Nz2YpFsWqA.jpeg. Last accessed 22 Mar 2019
6. Figure 3. Available at: https://cdn-images-1.medium.com/max/1600/1*SQAxk6gSUM_AK2d53neYmA.jpeg. Last accessed 22 Mar 2019
7. Figure 4, Available at: <https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/01/architecture.png>. Last accessed 22 Mar 2019
8. A Quantitative Study of Data in the NLP community, Margot Mieskes, Available at: <http://www.ethicsinnlp.org/workshop/EthNLP-2017.pdf>. Last accessed 25 Mar 2019
9. Google's new translation software is powered by brainlike artificial intelligence (2016, September 27), Available at: https://www.sciencemag.org/news/2016/09/google-s-new-translation-software-powered-brainlike-artificial-intelligence?r3f_986=https://www.google.com/. Last accessed 25 Mar 2019

10. Found in translation: More accurate, fluent sentences in Google Translate (2016, November 15), Available at: <https://www.blog.google/products/translate/found-translation-more-accurate-fluent-sentences-google-translate>. Last accessed 25 Mar 2019
11. 76 English Phrases for Traveling with Ease, Available at: <https://www.fluentu.com/blog/english/english-travel-phrases/>. Last accessed 25 Mar 2019
12. Travel English– 100 sentences you must know, Available at: <https://www.rong-chang.com/travel/>. Last accessed 25 Mar 2019
13. 100 Common English Phrases and Sentence Patterns (With Dialogue), Available at: <https://basicenglishspeaking.com/100-common-phrases-and-sentence-patterns/>. Last accessed 25 Mar 2019
14. Google Translate, Available at: <https://translate.google.com/>. Last accessed 25 Mar 2019
15. Ethical by Design: Ethics Best Practices for Natural Language Processing, Jochen L. Leidner and Vassilis Plachouras, Available at: <http://www.ethicsinnlp.org/workshop/pdf/EthNLP02.pdf/>. Last accessed 25 Mar 2019
16. How to Develop a Neural Machine Translation System from Scratch, by Jason Brownlee on January 10, 2018, Available at: <https://machinelearningmastery.com/develop-neural-machine-translation-system-keras/>. Last accessed 25 Mar 2019
17. SQL Alchemy, Available at: <https://www.sqlalchemy.org/>. Last accessed 25 Mar 2019
18. Found in translation: More accurate, fluent sentences in Google Translat, by Catherine Maticic on Sep. 27, 2016, Available at: <https://www.sciencemag.org/news/2016/09/google-s-new-translation-software-powered-brainlike-artificial-intelligence?r3f986=https://www.google.com/>. Last accessed 25 Mar 2019
19. What is Euclidean Distance?, On April 1, 2017, Available at: <https://machinelearningtutorials.com/euclidean-distance/>. Last accessed 25 Mar 2019
20. Text Similarities, Meenu Ahluwalia and Brian Coari, Available at: https://github.com/coarib/SMU_Masters_PST/blob/master/src/TextSimilarity_edited.ipynb. Last accessed 25 Mar 2019
21. Jiawei Han, Micheline Kamber and Jian Pei: Data Mining: Concepts and Techniques. 3rd ed. Elsevier Inc., Chapter 2.4.7 - Cosine Similarity, Page 77 (2012)

A SQL Queries

Select the text and sentence IDs for a language key. Language keys can be EN, PT, or HI for English, Portuguese, or Hindi, respectively.

```
select text, sentence_id from Sentences where language_key='[Language_Key]'
```

Select every combination of matching translations between all sentences. Storing this query in a dataframe allows us to get all the translations for a sentence just by matching a sentence ID to the input_sentence_id column.

```
SELECT a.sentence_id as input_sentence_id,
a.language_key as input_language_key, a.text as input_text,
```

```
b.text as output_text, b.sentence_id as output_sentence_id,  
b.language_key as output_language_key FROM masters.Translations  
left join masters.Sentences a on a.sentence_id=Translations.sentence_id_1  
left join masters.Sentences b on b.sentence_id=Translations.sentence_id_2
```