

## Project: Non-lexical sounds in dialogue utterances

Student name: *Dominik Künkele*

---

Course: *Machine Learning 2 (LT2326)*

Due date: *November 30th, 2022*

### Background

In one of their papers, Edlund and his colleagues describe human interaction with machines using metaphors on a spectrum. On the one side of the spectrum is the so called *interface metaphor*. The interface of the machine is built in a way that users are aware, they are talking to a machine. Consequently, they also adapt their language, and use more command like utterances (e.g. "Call John", "Set the timer") as they would fill slots in an imaginary web form. On the other end of the spectrum resides the *human metaphor*. Here, the users don't really know, they are talking to a machine and therefore use a "normal" language in the sense of utterances, they would use in a human-human dialogue. While the interface metaphor is still used very commonly, many dialogue systems like *Alexa* or *Siri* lie somewhere in between these extremes. The human metaphor is implemented rarely and is seen more often in science fiction.

Nevertheless, there are a few strong reasons for using the human metaphor and let the users talk in natural language. Natural language is

**easy to use.** Since we use natural language all the time in all human interactions, it is very natural for us, to also use it in machine interaction.

**flexible.** Natural language allows us to express everything we want to express. We can express for instance thoughts, feelings or facts with different certainties or also for example talk about things that were, things that are, and some things that have not yet come to pass. There are only very few things in a human mind that cannot be represented in natural language (partly in combination with mimic and gestures).

**resilient to error handling.** Furthermore, it allows us, to correct things we said or also specify certain parts more if we realize the listener is not understanding very easily. In the role of the listener, we can also verify, what we understood. Both can be done in various ways, by for instance rephrasing the utterance, repeating all or parts of it or even ask questions about it.

**enjoyable** Using Natural Language is finally also much more enjoyable than using a command like language. This of course depends on the person, but in general, humans evolved to use and interact in this language.

The goal is now, not to create a machine that is very close to humans, but to create "[a] machine that acts human enough that we respond to it as we respond to another human" as Cassel puts it in his paper in 2007.

There are endless factors that could define, *what* human-like interaction is. They include for example a change in pitch (upwards, to signal turn-keeping; downwards, to signal turn-yielding) or choices of words ("This is the automated booking system" will obviously signal the user that it speaks to a machine, while "Hello" does not). One big factor are non-lexical sounds (NLS) like hesitations ("uh", "hmm", ...) or repetitions of words/n-grams ("if I'm home then I I definitely watch her"). This could be famously observed in the *Google Duplex* demo in 2018, when a machine was calling a hair saloon using NLS.

This project aims to make machine utterances more human-like, by adding NLS to an utterance at natural positions. Additionally, the utterances could also be enhanced, by adding repetitions of n-grams.

### Data resource

As a resource for this project, I use the *Switchboard corpus*. There are various dialogue corpora published, but the Switchboard corpus contains NLS in a mostly structured way, while keeping the rest of the dialogue mostly undistilled, but annotated and as close as possible to the recordings. This means that the corpus includes

**non-grammatical sentences:** "yeah i think it is too it's gonna get better"

**repetitions of n-grams:** "oh she didn't she didn't do something"

**annotated partial words:** "bec- because we're so"

**anomalies of words with correct word** "bettle/better"

Furthermore, the corpus classifies the NLS based on their sound into thirteen classes: ah, eh, hm, huh, huh-uh, hum-um, ooh, uh, uh-huh, uh-hum, uh-oh, um, um-hum. These classes will then be used for predicting the NLS in the generated sentences.

I cleansed the corpus, by removing annotations, *silences* and *laughters*, replacing partial with complete words and NLS with a specific token per class. In the end, I could extract 247 123 utterances. The length of these utterances ranges from 1 to 81 words. Since I only used each utterance by itself to generate NLS and was not using any relations in the dialogue, I also excluded utterances shorter than two words. That yielded me around 150 000 usable utterances. Around 70 000 of these utterances contain at least one NLS.

Almost half of the utterances contain a repetition. Most of the repetitions are only on word, but I could identify also repetitions of up to 9-grams.

### Methods

The idea of this project was to predict for each word in a sentence if and which NLS should follow. For the sentence *uh-huh i think so and uh*, the representation for the model would look like this:

input		<SOS>	i	think	so	and	<EOS>
output		uh-huh	<NO-NLS>	<NO-NLS>	<NO-NLS>	uh	<NO-NLS>

For this, I tokenized the sentences, using *nltk's word\_tokenize* function, added start-of-sequence and end-of-sequence tokens and aligned it with the NLS. The input was then encoded with

### Results

My model has very high scores for all the metrics *Accuracy*, *Precision* and *Recall*, even though the predicted sentences are not really good. The reason for that is only a very small proportion of the slots are NLS, while the biggest proportion is the <NO-NLS> token. Therefore, my model that often only predicts <NO-NLS> tokens for all of the slots performs very well. To make this problem better visible, I introduced a new metric, called the *NLS score*. It is a weighted accuracy that weights <NO-NLS> tokens antiproportionally to its occurrence in the test corpus. More specifically the weights are calculated as following for the whole test corpus:

$$W_{NLS} = \frac{\text{number of NO\_NLS tokens}}{\text{number of slots}} \quad (1)$$

$$W_{NO\_NLS} = 1 - W_{NLS} \quad (2)$$

The weights are then applied to each token of a sentence, depending the gold label. If the gold label was a <NO\_NLS> token,  $W_{NO\_NLS}$  is applied, otherwise  $W_{NLS}$ . The sum of the weighted scores is then normalized over the length of the sentence and averaged over all sentences.

$$A_{WS} = \frac{W_{NLS} * \sum_0^t \text{correct NLS} + W_{NO\_NLS} * \sum_0^t \text{correct NO\_NLS}}{t} \quad (3)$$

$$A_W = \frac{\sum_0^n A_{WS}}{n} \quad (4)$$

Finally, the NLS score is the mean of the resulting weighted accuracy with the average ‘normal’ accuracy.

$$NLS \text{ score} = \frac{\sum_0^n A_W + \sum_0^n A}{2n} \quad (5)$$

This new metric takes the inbalance of <NO\_NLS> and NLS tokens better into account. It is also flexible in respect to a changing ratio, since it could also handle the opposite case if the biggest portion of slots were NLS tokens.

As expected the NLS score is much lower for all different models.

## Discussion