

Algorithms for speech and natural language processing: TD #1

Due on February 25, 2019

ACHER Clément

Part 1

The goal of this part is to build a classifier that given a voice command with a fixed length can recognize the spoken command. To achieve the best results, quite a few tricks can be used at various levels.

Feature extraction and feature engineering

The first thing I have done is to remove the cap on the number of samples used. Note that for the validation and test sets, this is almost necessary : the proposed implementation will only extract samples from the first couple classes. They then can't be used to validate a model.

Tuning the speech features

It is pretty rare in ASR to feed a classifier with the raw signals. Features like MFCCs and Mel Filter Banks are more suited for such tasks. They however both rely on the choice of a few hyperparameters. To tune this features, I have mostly rely on the resources from this website¹ and the default parameters from another feature extracting library² that I had used before. Table 1 shows the huge impact the choice of these parameters has : it shows the accuracy obtained using the two models initially included (no model tuning done at this point).

Accuracy on val. set	MFCC - Initial parameters	Tuned MFCC - no normalization	Tuned MFCC - w/ normalization	Mel Filter Banks - w/ normalization
Log. Reg.	3.6	14.92	30.46	10.42
MLP	39.4	53.68	65.92	52.66

Table 1: Accuracy on vanilla models using different features

Normalization

Another factor of improvement that can be done at the feature level is about normalization. The idea is to get all the 13 features of the MFCC (or Mel filter banks) within the same range. As it can be seen on the Figure 1a, the MFCC features are heavily dominated by the first coefficient. The 3 datasets are then normalized using the mean and the standard deviations of each coefficients across the whole training set. Note that there are other possibilities to normalize the datasets. Normalized MFCCs can be seen on the Figure 1b. Normalization has a positive impact on the accuracy of the classifiers as shown in table 1. The MLP classifier also converges faster (21 iterations vs. 29).

As MFCC features were giving the best results, I have decided to only use these.

Data augmentation

The dataset also provides various sounds that can be used to perform data augmentation. However, the voice command audio signal and the noise should not simply added : many noise are way louder than the voice command. Simply summing these signals makes the command inaudible. To weight the “quantity” of noise to add, I use the ratio of the RMS of the two signals.

Part 2

Question 2.1

The numerator is a sum of positive integers, so the WER can't be negative. It can be greater than 100, for instance if the reference is 'I like apple' and hypothesis “She really loves dogs and cats”, the WER is 200.

¹<http://practicalcryptography.com/miscellaneous/machine-learning/guide-mel-frequency-cepstral-coefficients-mfccs/>

²<https://python-speech-features.readthedocs.io/en/latest/>

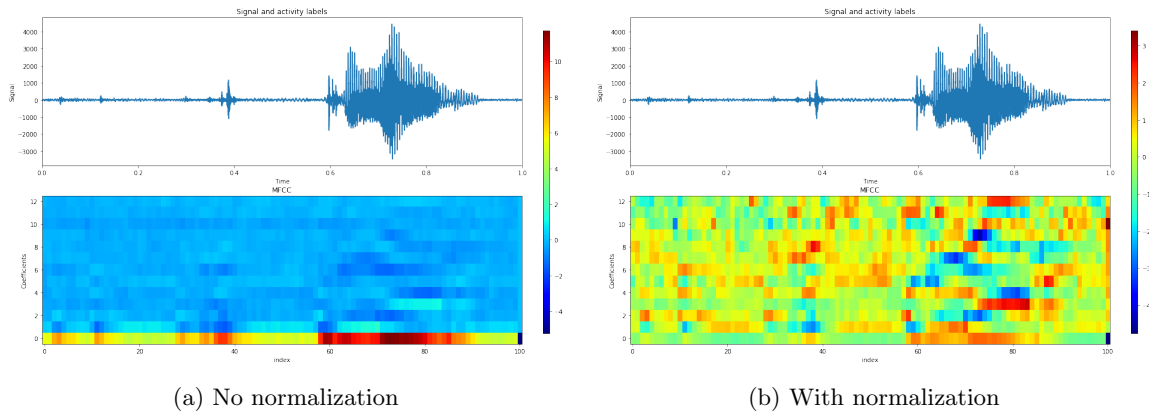


Figure 1: Raw signal (top) and MFCC feature (bottom) with and without normalization.

Question 2.2

`train_labels.count(label) < nb_ex_per_class` ensures that the dataset used for training is balanced and then that the prior probability of each word is equal.

Question 2.3

We have the following output:

True sentence: go marvin one right stop
 Predicted sentence with greedy search: go marvin on five stop

There are 2 substitutions out of the 5 words. The WER is then $\frac{2}{5} = 0.4$.

Question 2.4

In the bigram model, we have:

$$P(W) = \prod_{k=1}^T P(W_k | W_{k-1})$$

Question 2.5

The bigrams are stored in matrix $(P(w_k = j | w_{k-1} = i))_{i,j}$. To compute these probabilities, we go through all the sequences of the training set, and increment a the count of the encountered bigrams. The probabilities are then obtained by normalizing by the sum of the rows.

On top of that, there are two extra tricks : Laplace smoothing (i.e. add 1 to every bigram counts) is used to avoid 0-counts. We also add an extra “starting word” at the beginning of each sequence so that there is a bigram for the first word of the sequence. We can notice this way that go is often the first word of the sequence.

Question 2.6

Increasing N leads to a more accurate language model. However the number of entries in the transition matrix is M^N where M is the number of different possible words (31 here). The space complexity increases exponentially, and as the training dataset is rather small, it makes the matrix very sparse.

Question 2.7

We denote L the length of a sequence ($L = 5$), B the beam width and M the number of different words. The time complexity is then $\Theta(LBM(1 + \log(BM)))$ if a sorting algorithm with complexity $\Theta(n \log(n))$ is

used.

Question 2.8