

Hidden Markov Model for speech processing

Nguyen Van Duong¹⁺ and Le Thanh Tung²⁺

¹K62-CACLC3

²K62-CACLC2

⁺these authors contributed equally to this work

ABSTRACT

Modern general-purpose speech recognition systems are based on hidden Markov models. These are statistical models that output a sequence of symbols or quantities. HMMs are used in speech recognition because a speech signal can be viewed as a piecewise stationary signal or a short-time stationary signal. In a short time-scale (e.g., 10 milliseconds), speech can be approximated as a stationary process. Speech can be thought of as a Markov model for many stochastic purposes.

1 Introduction

Hidden Markov models are popular is because they can be trained automatically and are simple and computationally feasible to use. In speech recognition, the hidden Markov model would output a sequence of n-dimensional real-valued vectors (with n being a small integer, such as 10), outputting one of these every 10 milliseconds. The vectors would consist of cepstral coefficients, which are obtained by taking a Fourier transform of a short time window of speech and decorrelating the spectrum using a cosine transform, then taking the first (most significant) coefficients. The hidden Markov model will tend to have in each state a statistical distribution that is a mixture of diagonal covariance Gaussians, which will give a likelihood for each observed vector. Each word, or (for more general speech recognition systems), each phoneme, will have a different output distribution; a hidden Markov model for a sequence of words or phonemes is made by concatenating the individual trained hidden Markov models for the separate words and phonemes. In this work, we only focus on buiding a HMM model to recognize one single word in Vietnamese language.

2 Dataset

Dataset is collected from [Vnexpress](#) with a lot of work from all members of speech processing class. First, we crop the interested pieces of speech: "benh_nhan", "co_the", "khong", "nguoi", "duoc" from the large dataset above. Then, we split the dataset in train/test with proportion of 70/30.

3 Results

Top1 accuracy is properly used here for this task.

	Accuracy
duoc	0.89
nguoi	0.8
co_the	0.96
benh_nhan	0.9
khong	0.8
MeanAcc	0.87

4 Method

4.1 Train

As mention above, we crop 100 .wav files and split them in 70 samples for train and 30 samples for validation. Training phase follows step:

4.1.1 Read raw wav file

The raw wav is read by librosa library

4.1.2 Get MFCC features

We define 25ms per window length then slide it over 10ms per hop. Then we subtract the mean to normalize the data.

To get more features, we calculate the 1st and 2nd order of original mfcc feature. Finally we concatenate all the features to a vector

4.1.3 Clustering

We use Kmean clustering to separate the data in 10 part before put it into the HMM model

4.2 Validation

In validation phase, we do exactly the same as in training phase, but we try within or without Kmean clustering. We found out that Kmean actually improve the accuracy of validation data

4.3 Test

Finally, we self-record our own data to test the model.

Although our model performs well on validation set, it is not as good as on test set. We realize that the 2 data distribution are so different that our model may overfit on training set.

5 Future work

To be more accurate on test set, we manage to collect more data and improve the model.