

Speech Command Model

Pradeep Moturi (es16btech11016)

May 2020

Contents

1 Introduction

This is the report of my work done under Dr. G V V Sharma on building a speech command recognition model for a voice bot. The model used is based on concepts of Convolution, LSTM and Attention and is derived from [?].

2 Create Data

1. Set 16KHz as sampling rate
2. Record 80 utterances of each command.
3. Trim each utterance to one second.
4. Save samples of each command in different folders
Dataset/forward
Dataset/back
Dataset/left
Dataset/right
Dataset/stop

Used Audacity to do this.

3 Loading Data

I've used soundfile package to read the .wav. You may choose to use any other package like wavefile, librosa etc to do the same job.

As this is one of the slowest part, I've stored the loaded data as a numpy file for ease and speed of access. Now, I can load the data from npy file if repeating the experiment.

4 Split dataset

Do a stratified split of the dataset into train and test set with 20% as test samples.
Set a random seed for reproducing the split.

5 Augment data

Augment each audio sample by time shifting in 25000 length vectors filled with zeros.
Take steps of 500 to create 18 files per sample

6 Feature Extraction

MFCCs are most prominent features used in audio processing. Normalizing the MFCCs over the frequency axis is found to reduce effect of noise. Kapre is a python package that provides layers for

audio processing that are compatible with keras and utilize GPU for faster processing. Kapre provides us with a layer basically

```
Melspectrogram (padding='same', sr=16000, n_mels=39, n_dft = 1024, power_melgram=2.0, return_decibel_melgram=True, trainable_fb=False, trainable_kernel=False, name='mel_stft')
```

Arguments to the layer

padding: Padding when convoluting

sr: Sampling rate of audio provided

n_mels: number of coefficients to return

n_dft: width

power_melgram: exponent to raise log-mel-amplitudes before taking DCT. Using power 2 is shown to increase performance by reducing effect of noise

return_decibel_melgram: If to return log over values

trainable_fb: If filter bank trainable

trainable_kernel: If the kernel is trainable

7 Building Model

7.1 Concept

1. Using Convolutional layers ahead of LSTM is shown to improve performance in several research papers.
2. BatchNormalization layers are added to improve convergence rate.
3. Using Bidirectional LSTM is optimal when complete input is available. But this increases the runtime two-fold.
4. Final output sequence of LSTM layer is used to calculate importance of units in LSTM using a FC layer.
5. Then take the dot product of unit importance and output sequences of LSTM to get Attention scores of each time step.
6. Take the dot product of Attention scores and the output sequences of LSTM to get attention vector.
7. Add an additional FC Layer and then to output Layer with SoftMax Activation.

7.2 Hyper parameters

- **sparse_categorical_crossentropy** is used as **Loss** because only output which should be 1 is given instead of One Hot Encoding.

- **sparse_categorical_accuracy** is used as performance **Metric** for the above reason.
- **Adam** is used as **Optimizer**. Adam is adaptive learning rate optimization algorithm. This is shown to achieve a faster convergence because of having all the features of other optimization algorithms.

7.3 Notations

Operators:

- \times indicate matrix multiplication
- $*$ denote convolution (0 padding to same size)
- $.$ denote dot product
- $+, -$ can expand dimensions of their arguments

Format:

Layer Name (Layer Type) (Output Size).

Parameters

Equations

Output

$= \text{equation}$

Layer name indicates output of the corresponding layer.

Let us understand the maths behind the model.

7.4 Math

You can have a overall look at the architecture of the model in Fig [??]. Lets observe the math in each layer below.

0. Input (InputLayer) (49, 39, 1)

1. Conv1 (Conv2D) (49, 39, 10)

Parameters:

Kernel = (5, 1, 1, 10), Bias = (10)

Conv1[:, :, i]

$$= \text{Kernel}[:, :, :, i] * \text{Input} + \text{Bias}[i]$$

2. BN1 (BatchNormalization) (49, 39, 10)

Parameters:

Trainable: $\gamma = (10)$, $\beta = (10)$,

Non-Trainable: Mean = (10), Std = (10)

Equations:

$$\text{Mean}[i] = \text{mean}(\text{Conv1}[:, :, i])$$

$$\text{Std}[i] = \text{std}(\text{Conv1}[:, :, i])$$

BN1[i]

$$= (\text{Conv1}[:, :, i] - \text{Mean}[i]) \frac{\gamma[i]}{\text{Std}[i]} + \beta[i]$$

3. Conv2 (Conv2D) (49, 39, 1)

Parameters:

Kernel = (5, 1, 10, 1), Bias = (1)

Conv2[:, :, 1]

$$= \text{Kernel}[:, :, :, 1] * \text{BN1} + \text{Bias}$$

4. BN2 (BatchNormalization) (49, 39, 1)

Parameters:

Trainable: $\gamma = (1)$, $\beta = (1)$,

Non-Trainable: Mean = (1), Std = (1)

Equations:

$$\text{Mean}[i] = \text{mean}(\text{Conv2}[:, :, i])$$

$$\text{Std}[i] = \text{std}(\text{Conv2}[:, :, i])$$

BN2[i]

$$= (\text{Conv2}[:, :, i] - \text{Mean}[i]) \frac{\gamma[i]}{\text{Std}[i]} + \beta[i]$$

5. Squeeze (Reshape) (49, 39)

Squeeze

$$= \text{BN2.reshape}(49, 39)$$

6. LSTMSequences (LSTM) (49, 64)

Parameters:

$$U^i = U^f = U^o = U^g = (39, 64),$$

$$W^i = W^f = W^o = W^g = (64, 64),$$

$$B^i = B^f = B^o = B^g = (64)$$

Equations:

$$i_t = \sigma(\text{Squeeze}[:, t] \times U^i + h_{t-1} \times W^i + B^i)$$

$$f_t = \sigma(\text{Squeeze}[:, t] \times U^f + h_{t-1} \times W^f + B^f)$$

$$o_t = \sigma(\text{Squeeze}[:, t] \times U^o + h_{t-1} \times W^o + B^o)$$

$$\tilde{C}_t = \tanh(\text{Squeeze}[:, t] \times U^g + h_{t-1} \times W^g + B^g)$$

$$C_t = \sigma(f_t * C_{t-1} + i_t * \tilde{C}_t)$$

$$h_t = \tanh(C_t) * o_t$$

LSTMSequences[t]

$$= h_t$$

7. FinalSequence (Lambda) (64)

FinalSequence

$$= \text{LSTMSequences}[-1, :]$$

8. UnitImportance (Dense) (64)

Parameters:

Weights = (64,64), Bias = (64)

UnitImportance

$$= \text{Weights} \times \text{FinalSequence} + \text{Bias}$$

9. AttentionScores (Dot) (49)

AttentionScores[i]

$$= \text{UnitImportance.LSTMSequences}[i, :]$$

10. AttentionSoftmax (Softmax) (49)

AttentionSoftmax[i]

$$= \frac{\exp(\text{AttentionScores}[i])}{\sum_j \exp(\text{AttentionScores}[j])}$$

11. AttentionVector (Dot) (64)

AttentionVector[i]

$$= \text{AttentionSoftmax.LSTMSequences}[:, i]$$

12. FC (Dense) (32)

Parameters:

Weights = (64,64), Bias = (64)

FC

$$= \text{Weights} \times \text{AttentionVector} + \text{Bias}$$

13. Output (Dense) (5)

Parameters: