

## Audio Modelling

June 2, 2022

## 1 Modelling and Deployment using MLOps

Now that we have audio input data & corresponding labels in an array format, it is easier to consume and apply Natural language processing techniques. We can convert audio files labels into integers using label Encoding or One Hot Vector Encoding for machines to learn. The labeled dataset will help us in the neural network model output layer for predicting results. These help in training & validation datasets into nD array. At this stage, we apply other pre-processing techniques like dropping columns, normalization, etc. to conclude our final training data for building models. Moving to the next stage of splitting the dataset into train, test, and validation is what we have been doing for other models. We can leverage CNN, RNN, LSTM, CTC etc. deep neural algorithms to build and train the models for speech applications like speech recognition. The model trained with the standard size few seconds audio chunk transformed into an array of n dimensions with the respective labels will result in predicting output labels for test audio input. As output labels will vary beyond binary, we are talking about building a multi-class label classification method.

```
[5]: import pandas as pd
import numpy as np
import os, sys
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
sys.path.append(os.path.abspath(os.path.join('../scripts')))
import tensorflow as tf
from clean import Clean
from utils import vocab
from deep_learner import DeepLearn
from modeling import Modeler
from evaluator import CallbackEval
```

```
[6]: AM_ALPHABET='          auīāēəo'
      EN_ALPHABET='abcdefghijklmnopqrstuvwxyz'
```

```
[7]: cleaner = Clean()
char to num,num to char=vocab(AM ALPHABET)
```

```
2022-06-02 08:09:18,883:logger:Successfully initialized clean class
```

The vocabulary is: [' ', '!', '"', '#', '\$', '%', '&', '\'', '(', ')', '\*', '+', ',', '-', '.', ':', ';', '<', '=', '>', '?', '@', '[', '\\', '^', '\_', '`', '{', '|', '}', '~', '0', '1', '2', '3', '4', '5', '6', '7', '8', '9', 'A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'I', 'J', 'K', 'L', 'M', 'N', 'O', 'P', 'Q', 'R', 'S', 'T', 'U', 'V', 'W', 'X', 'Y', 'Z', 'a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z']

```

'', '', '', '', '', '', '', '', '', 'a', 'u', 'i', 'ā', 'e', 'ə', 'o']
(size =44)

```

## 2 Deep Learning Model

**objective:** Build a Deep learning model that converts speech to text.

```

[8]: swahili_df = pd.read_csv("../data/swahili.csv")
    amharic_df = pd.read_csv("../data/amharic.csv")

[9]: pre_model = Modeler()

[10]: swahili_preprocessed = pre_model.preprocessing_learn(swahili_df, 'key', 'file')

[11]: amharic_preprocessed = pre_model.preprocessing_learn(amharic_df, 'key', 'file')

[12]: train_df, val_df, test_df = amharic_preprocessed

[13]: batch_size = 32
    # Define the training dataset
    train_dataset = tf.data.Dataset.from_tensor_slices(
        (list(train_df["file"]), list(train_df["text"])))
    )
    train_dataset = (
        train_dataset.map(cleaner.encode_single_sample, num_parallel_calls=tf.data.
        ↪AUTOTUNE)
        .padded_batch(batch_size)
        .prefetch(buffer_size=tf.data.AUTOTUNE)
    )

    # Define the validation dataset
    validation_dataset = tf.data.Dataset.from_tensor_slices(
        (list(val_df["file"]), list(val_df["text"])))
    )
    validation_dataset = (
        validation_dataset.map(cleaner.encode_single_sample, num_parallel_calls=tf.
        ↪data.AUTOTUNE)
        .padded_batch(batch_size)
        .prefetch(buffer_size=tf.data.AUTOTUNE)
    )

```

### 2.1 LSTM Deep Learning

```

[14]: learn = DeepLearn(input_width=1, label_width=1, shift=1, epochs=5,
    train_df=train_df, val_df=val_df, test_df=test_df,
    label_columns=['mfcc-0'])
fft_length = 384

```

```

model = learn.build_asr_model(
    input_dim=fft_length // 2 + 1,
    output_dim=char_to_num.vocabulary_size(),
    rnn_units=512,
)
model.summary(line_length=110)

```

Model: "DeepSpeech\_2"

Layer (type) Param #	Output Shape
input (InputLayer) 0	[(None, None, 193)]
expand_dim (Reshape) 0	(None, None, 193, 1)
conv_1 (Conv2D) 14432	(None, None, 97, 32)
conv_1_bn (BatchNormalization) 128	(None, None, 97, 32)
conv_1_relu (ReLU) 0	(None, None, 97, 32)
conv_2 (Conv2D) 236544	(None, None, 49, 32)
conv_2_bn (BatchNormalization) 128	(None, None, 49, 32)
conv_2_relu (ReLU) 0	(None, None, 49, 32)
reshape (Reshape) 0	(None, None, 1568)
bidirectional_1 (Bidirectional) 6395904	(None, None, 1024)
dropout (Dropout) 0	(None, None, 1024)

bidirectional_2 (Bidirectional)	(None, None, 1024)
4724736	
dropout_1 (Dropout)	(None, None, 1024)
0	
bidirectional_3 (Bidirectional)	(None, None, 1024)
4724736	
dropout_2 (Dropout)	(None, None, 1024)
0	
bidirectional_4 (Bidirectional)	(None, None, 1024)
4724736	
dropout_3 (Dropout)	(None, None, 1024)
0	
bidirectional_5 (Bidirectional)	(None, None, 1024)
4724736	
dense_1 (Dense)	(None, None, 1024)
1049600	
dense_1_relu (ReLU)	(None, None, 1024)
0	
dropout_4 (Dropout)	(None, None, 1024)
0	
dense (Dense)	(None, None, 45)
46125	

```
=====
=====
Total params: 26,641,805
Trainable params: 26,641,677
Non-trainable params: 128
-----
-----
```

### 3 Evaluation

**objective:** Evaluate your model.

```
[ ]: epochs = 1
    # Callback function to check transcription on the val set.
```

```
validation_callback = CallbackEval(model,validation_dataset)
# Train the model
history = model.fit(
    train_dataset,
    validation_data=validation_dataset,
    epochs=epochs,
    callbacks=[validation_callback],
)
```

```
2022-06-02 08:10:28.076462: W
tensorflow/core/framework/cpu_allocator_impl.cc:82] Allocation of 39458048
exceeds 10% of free system memory.
2022-06-02 08:10:29.155589: W
tensorflow/core/framework/cpu_allocator_impl.cc:82] Allocation of 39458048
exceeds 10% of free system memory.
2022-06-02 08:10:29.223129: W
tensorflow/core/framework/cpu_allocator_impl.cc:82] Allocation of 19932416
exceeds 10% of free system memory.
2022-06-02 08:10:31.121990: W
tensorflow/core/framework/cpu_allocator_impl.cc:82] Allocation of 19932416
exceeds 10% of free system memory.
2022-06-02 08:10:31.547098: W
tensorflow/core/framework/cpu_allocator_impl.cc:82] Allocation of 19932416
exceeds 10% of free system memory.
```

[ ]: