

Machine Learning Engineer Nanodegree

Capstone Report

Kannappan Natarasan
May 12th 2018

I . Definition

Project Overview

Able to understand and classify a sound file by an application fascinated me, since I'm planning to build an App which can apply these techniques on machines in factory thus those can be monitored remotely and any mechanical malfunction could be predicted earlier.

Thanks to Deep learning content in MLND, which gave me nice exposure to handle data in array format (the assignment was classifying dog images). For the capstone project when I've been exploring Kaggle challenges, sound classifier challenge attracted me, since I already got exposure to identify and classify image files which are 3D arrays. Moving to sound may give me opportunity to understand handling audio files.

Picking a challenge from Kaggle gives me the luxury of clean dataset, so that I can focus more on solving prediction problem rather spending too much time on data preparation. Since my background is from data warehousing and Big data, I never felt data processing is a problem. I was able to download required dataset and information about labels from the following link.
<https://www.kaggle.com/c/freesound-audio-tagging>

Though machine learning has been exist for the past 30 years (earlier it has been known as data mining), it popped up to common usage, due to following reasons.

- i) Maturity of internet and continuously expanded network bandwidth enabled collecting data from various ends feasible
- ii) Continuous improvement in hardware (Disk, RAM and CPU) enabled to store and process these data.
- iii) Especially improvement in CPU, GPU and TPU power enabled to run machine learning algorithms on larger dataset possible. GPU power enabled to build Neural networks upto multiple layers which is called deep learning.

When I was exploring similar project in web found following links, Giants like Amazon Google Apple went to advanced level like able to process human speech in real time. I feel starting from here would at least put me in first step in the journey of voice processing.

<https://medium.com/@ageitgey/machine-learning-is-fun-part-6-how-to-do-speech-recognition-with-deep-learning-28293c162f7a>
<https://www.iotforall.com/tensorflow-sound-classification-machine-learning-applications/>
<https://www.analyticsvidhya.com/blog/2017/08/audio-voice-processing-deep-learning/>

Problem statement

The training dataset contains around 9000 audio files classified into 41 labels . Audio files in training set are in wav format, so finding python packages which can parse this format and convert into array format would help to prepare dataset format which could be fed into keras deep learning algorithms.

Finding necessary preprocessing steps would help to execute these algorithms efficiently. Suitable benchmark method would help to assess the model built. Defining proper metrics would be helpful to measure each learning iteration and tuning the model to improve accuracy.

Developing right CNN architecture would be vital to train the model and derive a classifier.

Metrics

The prediction model built in this project would be evaluated by Mean Average Precision@3 (MAP@3). This is a measurement method suggested in Kaggle challenge.

$$MAP@3 = \frac{1}{U} \sum_{u=1}^U \sum_{k=1}^{\min(n,3)} P(k)$$

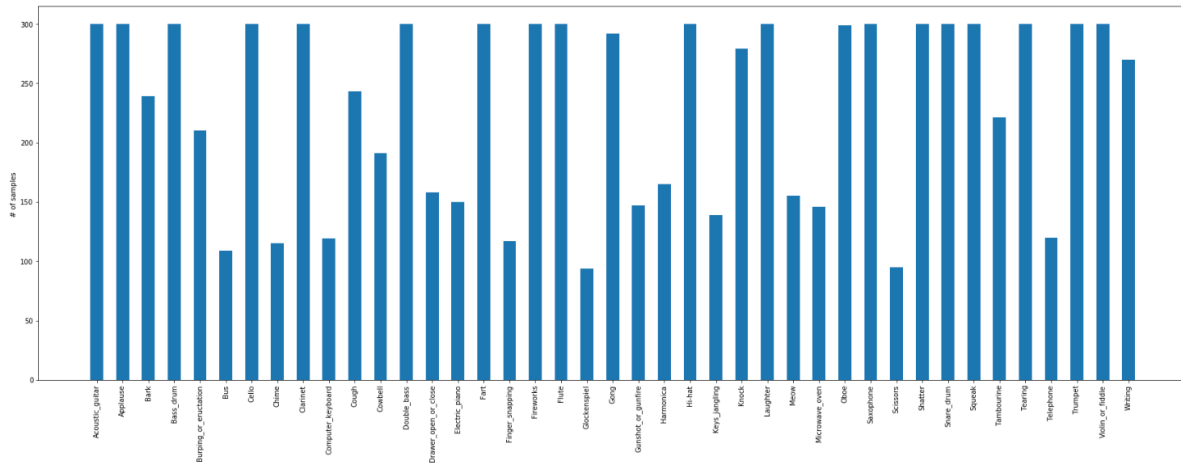
U : the number of scored audio files in the test data,
P(k) : the precision at cutoff k
n : the number predictions per audio file.

I've decided to use this method since it measures accuracy of prediction not only based on exact predicted value , but also considers next 2 predictions in pipeline, i.e , imagine an audio file from validation set, say it is Guitar, and predicted value is Double-bass. Other metrics system would immediately mark as false . But using MAP@3 would look for next two predictions , if the second prediction is Guitar , it would add some positive weightage for the accuracy measure, which would obviously less than if the first prediction is Guitar. For the App in my mind would be benefited if it can measure next two predictions.

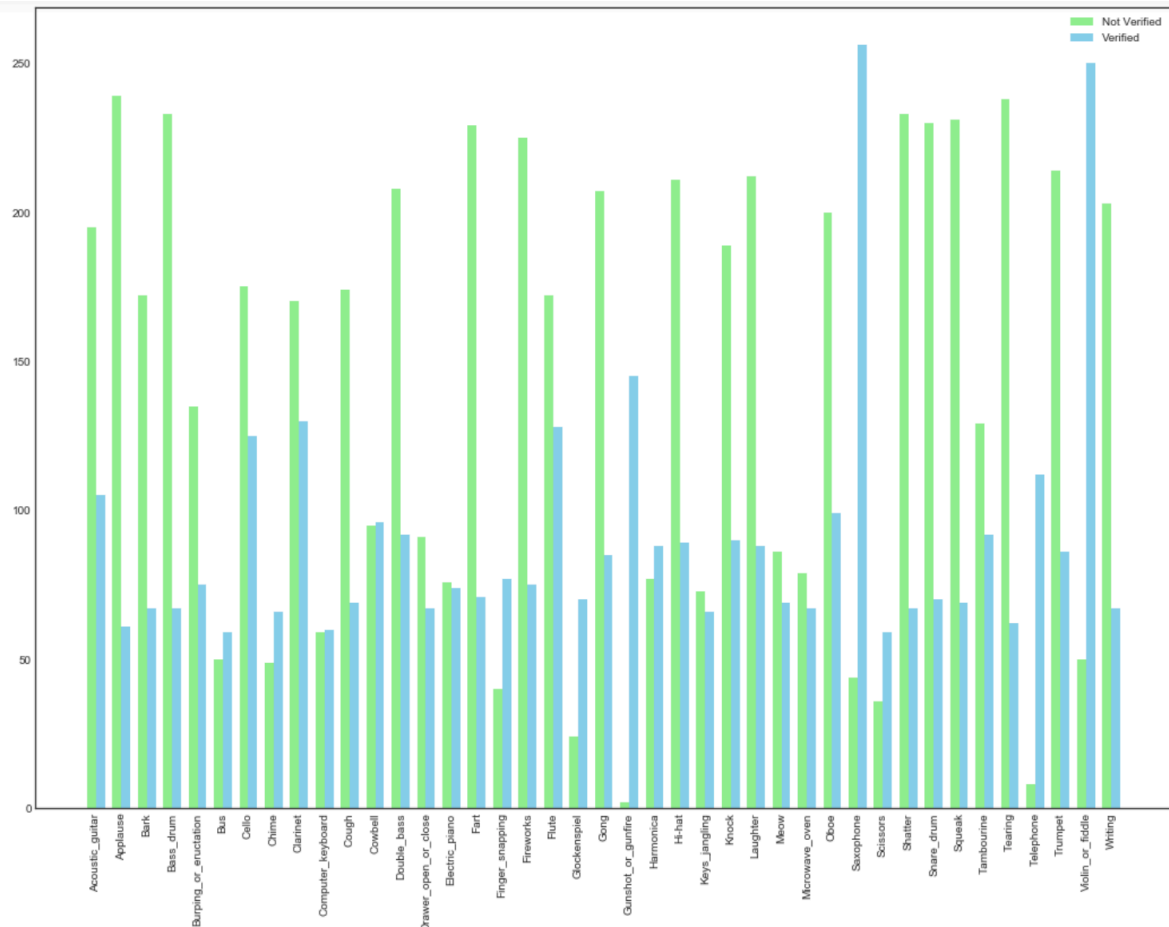
II . Analysis

Data Exploration

Training set used in this problem is made up of 9473 audio files in .wav format and classified into 41 labels. Labels like *Acoustic guitar*, *Applause*, *Bass drum*, *Cello*, *Clarinet*, *Double bass*, etc has 300 files each. Labels like *Bus*, *Scissors*, *Glockenspiel* has audio files close to 100 audio files.



On top of above following graph shows on each label , number of files which are human verified.



Properties of Single Audio file

To understand a single audio file python packages like `scipy` was useful, which revealed that following two properties vital to understand sound files.

i) **Frame rate** (frames per second):

Decides quality of sound , higher the frame rate could be clear sound, found that framerate of all audio files in this data is constant 44100, refer following snippet

```

sample_rat=[]
sample_len=[]

for i in trainFilesList['fname']:
    sample_rate,samples=wavfile.read(home+'data/audio_train/'+i)
    sample_rat.append(sample_rate)
    sample_len.append(len(samples))

print('Sample rate on all files is ',set(sample_rat))

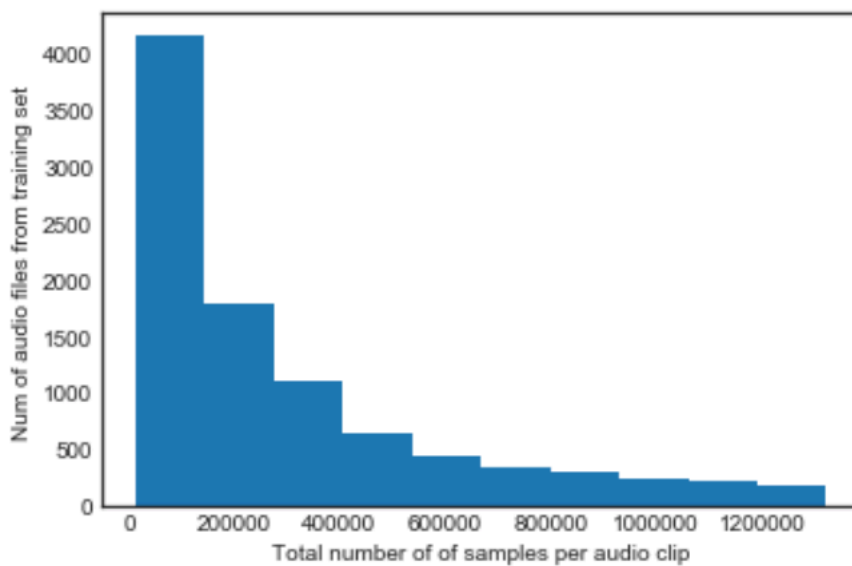
```

Sample rate on all files is {44100}

ii) **Number of samples/frames per audio file :**

Size of the audio file, more samples in a file means the length of audio is longer.

Following histogram provides number of number of audio files with bucket of number of audio samples (length of audio file) in each file.

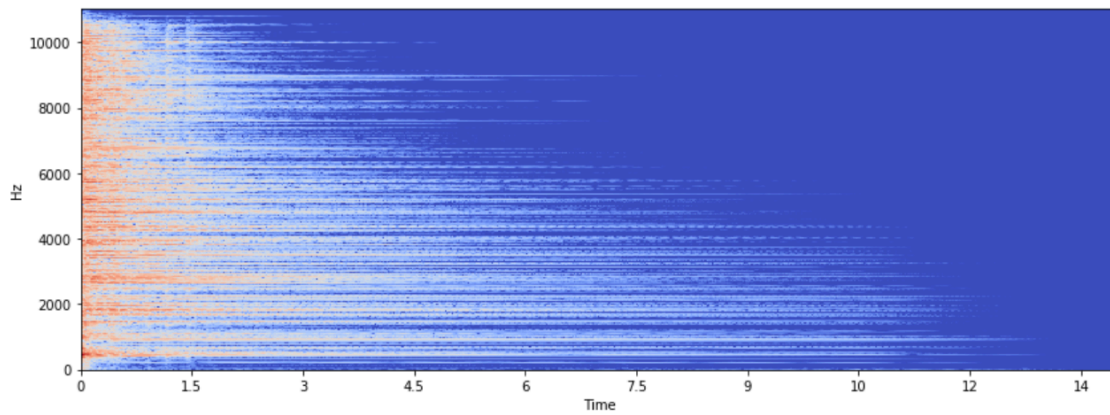
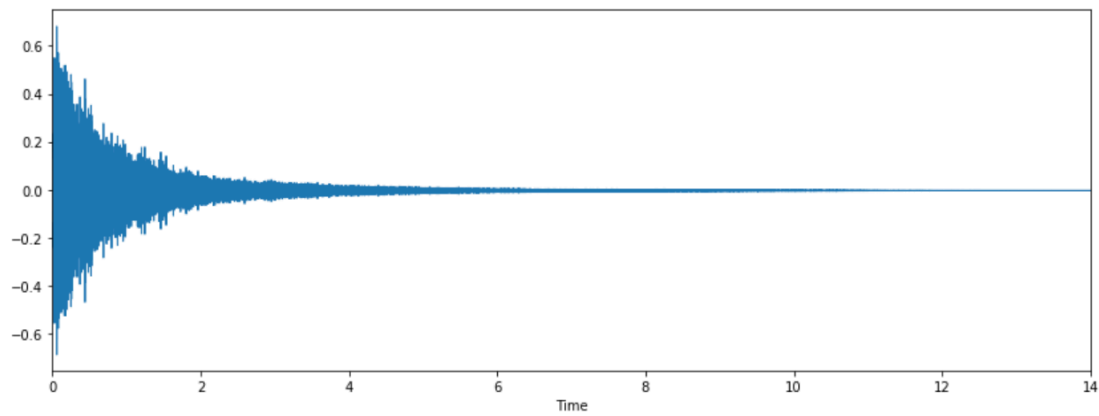


Close to 4k files have 200k audio samples, this is half of the dataset. In remaining half close to 2k files has samples between 200k and 400k . Less than 500 files have samples as high as 1200k.

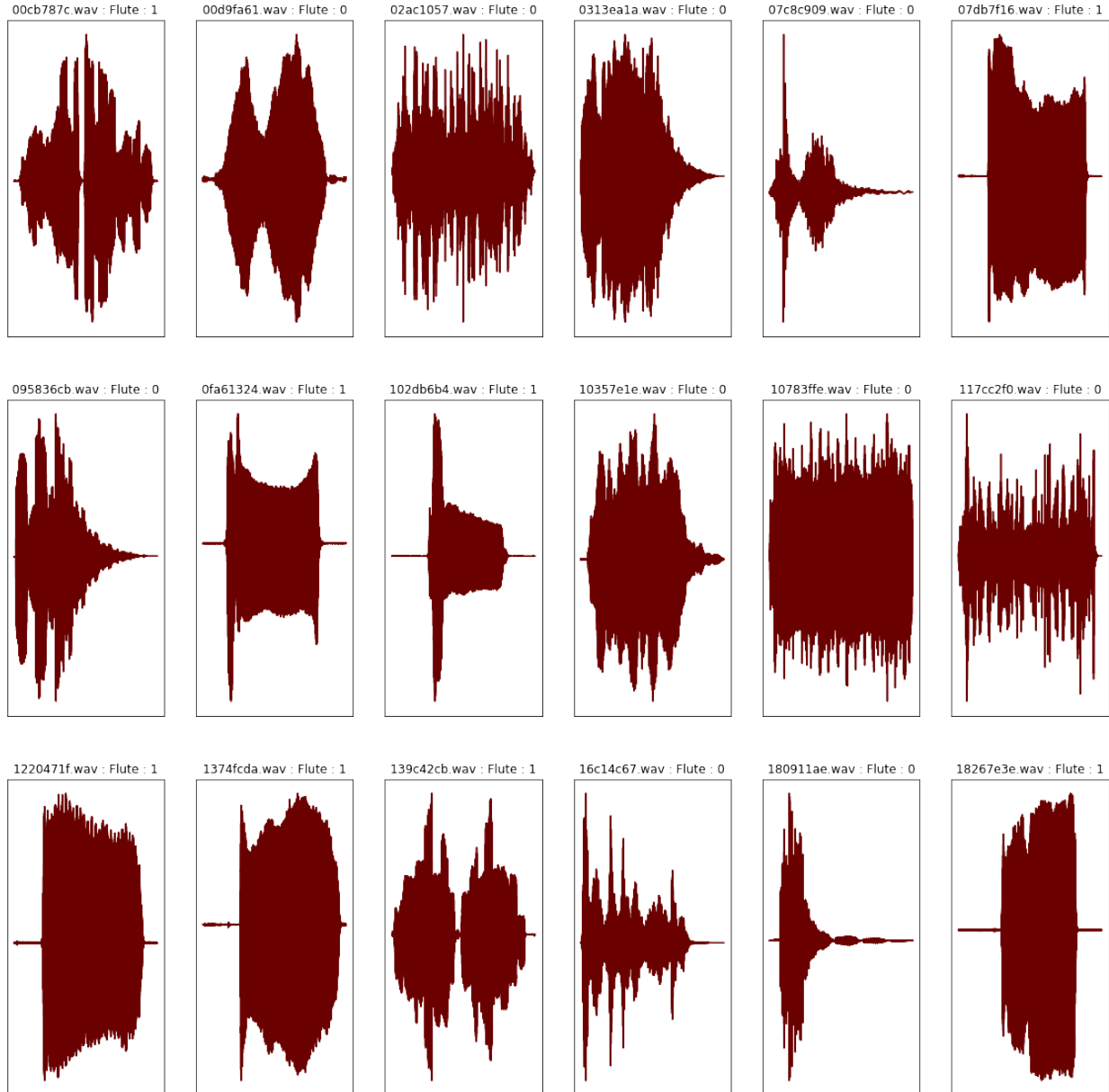
Exploratory Visualization

Visualizing individual audio files would be interesting to study possible features which could automatically extracted and used by CNN architecture. Package librosa helped to visualize individual files. Sound features within a second could be understood by following visuals.

First two seconds are intense, After two seconds sound intensity deteriorates gradually and becomes zero. Following two graphs depicts this effect of same audio file.



Followings are smaller size visuals of same kind of sound files (Flue files) , which would help to see the pattern of same kind of audio files.



Algorithms and Techniques

I've used Convolutional Neural Networks to solve this problem through keras implementation . Since there are similarities between sound files and image files (except sound is single dimension array but image could be 2D array for gray and 3D array for color image) and CNN is proven solution through MNIST and CIFAR 100, I've decided to use CNN for this problem.

Convolutional Neural Network

Epochs : 20 (An epoch is an iteration over the entire X and y data provided. If suppose epoch=20 , then entire dataset is iterated on training for 20 times. Between epochs weights kept updated based on reduction in loss).

Filter : 16 (Number of filter decides how many patterns to detect)

Slider size: 9 (This decides number of steps to be moved by a filter)

Padding : valid (valid is no padding , When the slider moves, action to be taken in the edge is decided by this value)

Layers : This could be Input layer, Convolution layer, Pooling layer, Output layer

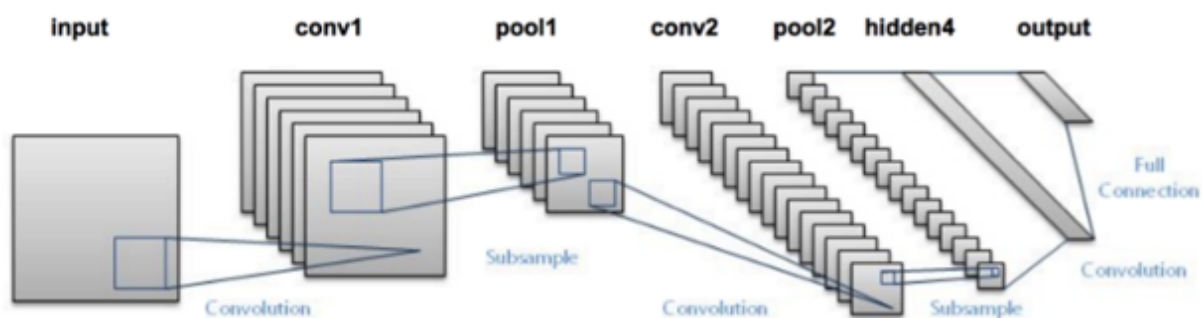
Activation function: Adam , smoothens the output between 0 to 1

Loss function : categorical_crossentropy – since this problem is to categorize between 41 labels (This helps to decide set value of weights for every pattern)

Layer (type)	Output Shape	Param #
=====		
input_4 (InputLayer)	(None, 32000, 1)	0
conv1d_25 (Conv1D)	(None, 31992, 16)	160
conv1d_26 (Conv1D)	(None, 31984, 16)	2320
max_pooling1d_10 (MaxPooling)	(None, 1999, 16)	0
dropout_13 (Dropout)	(None, 1999, 16)	0
conv1d_27 (Conv1D)	(None, 1997, 32)	1568
conv1d_28 (Conv1D)	(None, 1995, 32)	3104
max_pooling1d_11 (MaxPooling)	(None, 498, 32)	0
dropout_14 (Dropout)	(None, 498, 32)	0
conv1d_29 (Conv1D)	(None, 496, 32)	3104
conv1d_30 (Conv1D)	(None, 494, 32)	3104
max_pooling1d_12 (MaxPooling)	(None, 123, 32)	0
dropout_15 (Dropout)	(None, 123, 32)	0
conv1d_31 (Conv1D)	(None, 121, 256)	24832
conv1d_32 (Conv1D)	(None, 119, 256)	196864
global_max_pooling1d_4 (Glob	(None, 256)	0
dropout_16 (Dropout)	(None, 256)	0
dense_10 (Dense)	(None, 64)	16448
dense_11 (Dense)	(None, 1028)	66820

dense_12 (Dense)	(None, 41)	42189
=====		
Total params: 360,513		
Trainable params: 360,513		
Non-trainable params: 0		
=====		

CNN is made up of multiple layers, starting with input layer, pooling layer, drop out layer. The final layer would converge into a dense layer made up of all labels. Input layer is a matrix, then convolution layer, a small window slide over entire matrix and outputs next hidden layer. Convolution layer uses randomly assigned weights to get next layer. Using activation function output from convolution layer would be transformed between 0 to 1. A filter with specific pattern traces for same pattern in all sample arrays. Weights at filter is assigned randomly, between iterations these random weights are updated based reduction in loss. CNN decides what kinds patterns it want to detect based on loss function. Stride decides amount the filter moves for a particular layer..



Source: <https://blog.dataiku.com/deep-learning-with-dss>

Input layer:

This is initial layer the window has to be equal to input file array size .

Conv1D: This is Convolution layer number of filters size of slide window , padding strategy and activation methods are set here .

Max Pooling : This layer helps to consolidate many neurons to single, this can be selected from average pooling or max pooling.

Dropout : This layer is to dropout un interesting features and deepen the network only with most impacting features.

Dense : This are final layers to predict labels from derived features.

Benchmark model

As this is Kaggle challenge, I've observed the leaderboard for this competition, best performing algorithms who are in #1 rank scored 95% accuracy. By comparing this I've set an accuracy benchmark of 70%.

For the application where I'm going to use this model , a mean accuracy of 70% MAP@3 would be great. As the app I'm planning to build going to collect values from various other sensors, I can handle this 70% accuracy with other values.

III . Methodology

Data PreProcessing

Since dataset is downloaded from Kaggle, it can be assumed that the data is clean, when looking for any outliers in features like framerate, number of samples per audio I haven't seen anything significant to be removed.

However CNN would expect every data point would be same shape array for learning . But in the dataset number of samples per audio file or length of audio is not constant. I think this can happen for image files also, since image files could vary based on format , resolution and size would easily end up differing sized arrays. For this dataset I've decided to standardize size of audio that I would use for training as well. Audio length is parameterized and for files with smaller than the standard length, say 2 seconds, I would pad them with zero then audio files are normalized using following formula,

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Implementation

When I've started building CNN model, set input shape more than 2 seconds of audio, but the data set doesn't have more than 2 seconds across all samples so learning process was failing. Later decided to go with 2 seconds which also accelerated learning process. Since it is a classification problem. I've used categorical_crossentropy loss function. Optimizer I've decided to go from Adam.

Decided to go with epoch 20 and fold 2 arbitrarily, but this took lot of time to complete, at least 3 hours. I have to leave it overnight for a full run to complete. Set checkpoint to save best only , since final best model is enough , but saving all would help if overfitting occurs , but

when I predicted files from test set (whose labels unknown) observed accuracy 53% which is close to training accuracy.

Input layer has to be in the shape of input files, for this project being framerate is 16000 and I've taken 2 seconds of audio for training shape of Input layer would be (32000,1).

Refinement

By executing different learning flows by modifying size of audio file, number of folds for stratifiedKfold and max epochs shows that , Increasing number folds, audio duration doesn't reduce validation loss.

of Audio samples : 200

config = Config(sampling_rate=16000, audio_duration=2, n_folds=4, max_epochs=4)

```
Fold: 0
#####
Epoch 1/4
3/3 [=====] - 14s 5s/step - loss: 3.7132 - acc: 0.0461 - val_loss: 3.713
1 - val_acc: 0.0000e+00

Epoch 00001: val_loss improved from inf to 3.71309, saving model to models/best_0.h5
Epoch 2/4
3/3 [=====] - 9s 3s/step - loss: 3.7120 - acc: 0.0345 - val_loss: 3.7126
- val_acc: 0.0462

Epoch 00002: val_loss improved from 3.71309 to 3.71264, saving model to models/best_0.h5
Epoch 3/4
3/3 [=====] - 9s 3s/step - loss: 3.7106 - acc: 0.0664 - val_loss: 3.7122
- val_acc: 0.0462

Epoch 00003: val_loss improved from 3.71264 to 3.71218, saving model to models/best_0.h5
Epoch 4/4
3/3 [=====] - 10s 3s/step - loss: 3.7096 - acc: 0.0664 - val_loss: 3.711
7 - val_acc: 0.0462

Epoch 00004: val_loss improved from 3.71218 to 3.71166, saving model to models/best_0.h5
Fold: 1
#####
Epoch 1/4
3/3 [=====] - 14s 5s/step - loss: 3.7129 - acc: 0.0380 - val_loss: 3.713
0 - val_acc: 0.0339

Epoch 00001: val_loss improved from inf to 3.71301, saving model to models/best_1.h5
Epoch 2/4
3/3 [=====] - 11s 4s/step - loss: 3.7118 - acc: 0.0551 - val_loss: 3.712
6 - val_acc: 0.0339

Epoch 00002: val_loss improved from 3.71301 to 3.71257, saving model to models/best_1.h5
Epoch 3/4
3/3 [=====] - 10s 3s/step - loss: 3.7103 - acc: 0.0704 - val_loss: 3.712
1 - val_acc: 0.0339

Epoch 00003: val_loss improved from 3.71257 to 3.71205, saving model to models/best_1.h5
Epoch 4/4
3/3 [=====] - 11s 4s/step - loss: 3.7085 - acc: 0.0647 - val_loss: 3.711
5 - val_acc: 0.0339

Epoch 00004: val_loss improved from 3.71205 to 3.71146, saving model to models/best_1.h5
```

```

Fold: 2
#####
Epoch 1/4
3/3 [=====] - 15s 5s/step - loss: 3.7134 - acc: 0.0220 - val_loss: 3.713
1 - val_acc: 0.0465

Epoch 00001: val_loss improved from inf to 3.71307, saving model to models/best_2.h5
Epoch 2/4
3/3 [=====] - 10s 3s/step - loss: 3.7129 - acc: 0.0628 - val_loss: 3.712
5 - val_acc: 0.0465

Epoch 00002: val_loss improved from 3.71307 to 3.71255, saving model to models/best_2.h5
Epoch 3/4
3/3 [=====] - 11s 4s/step - loss: 3.7119 - acc: 0.0691 - val_loss: 3.711
9 - val_acc: 0.0465

Epoch 00003: val_loss improved from 3.71255 to 3.71195, saving model to models/best_2.h5
Epoch 4/4
3/3 [=====] - 11s 4s/step - loss: 3.7112 - acc: 0.0746 - val_loss: 3.711
3 - val_acc: 0.0465

Epoch 00004: val_loss improved from 3.71195 to 3.71125, saving model to models/best_2.h5
Fold: 3
#####
Epoch 1/4
3/3 [=====] - 14s 5s/step - loss: 3.7137 - acc: 0.0428 - val_loss: 3.712
9 - val_acc: 0.0606

Epoch 00001: val_loss improved from inf to 3.71286, saving model to models/best_3.h5
Epoch 2/4
3/3 [=====] - 9s 3s/step - loss: 3.7131 - acc: 0.0374 - val_loss: 3.7123
- val_acc: 0.0909

Epoch 00002: val_loss improved from 3.71286 to 3.71235, saving model to models/best_3.h5
Epoch 3/4
3/3 [=====] - 11s 4s/step - loss: 3.7125 - acc: 0.0399 - val_loss: 3.711
8 - val_acc: 0.0909

Epoch 00003: val_loss improved from 3.71235 to 3.71177, saving model to models/best_3.h5
Epoch 4/4
3/3 [=====] - 9s 3s/step - loss: 3.7120 - acc: 0.0453 - val_loss: 3.7111
- val_acc: 0.0606

Epoch 00004: val_loss improved from 3.71177 to 3.71108, saving model to models/best_3.h5

```

of Audio samples : 200

config = Config(sampling_rate=16000, audio_duration=4, n_folds=2, max_epochs=4)

```

Fold: 0
#####
Epoch 1/4
2/2 [=====] - 20s 10s/step - loss: 3.7138 - acc: 0.0155 - val_loss: 3.71
34 - val_acc: 0.0648

Epoch 00001: val_loss improved from inf to 3.71344, saving model to models/best_0.h5
Epoch 2/4
2/2 [=====] - 15s 7s/step - loss: 3.7133 - acc: 0.0576 - val_loss: 3.713
1 - val_acc: 0.0648

Epoch 00002: val_loss improved from 3.71344 to 3.71310, saving model to models/best_0.h5
Epoch 3/4
2/2 [=====] - 15s 8s/step - loss: 3.7127 - acc: 0.0420 - val_loss: 3.712
8 - val_acc: 0.0556

Epoch 00003: val_loss improved from 3.71310 to 3.71279, saving model to models/best_0.h5
Epoch 4/4
2/2 [=====] - 15s 8s/step - loss: 3.7122 - acc: 0.0353 - val_loss: 3.712
5 - val_acc: 0.0185

```

```

Epoch 00004: val_loss improved from 3.71279 to 3.71248, saving model to models/best_0.h5
Fold: 1
#####
Epoch 1/4
2/2 [=====] - 20s 10s/step - loss: 3.7140 - acc: 0.0083 - val_loss: 3.71
33 - val_acc: 0.0326

Epoch 00001: val_loss improved from inf to 3.71329, saving model to models/best_1.h5
Epoch 2/4
2/2 [=====] - 17s 8s/step - loss: 3.7130 - acc: 0.0545 - val_loss: 3.712
9 - val_acc: 0.0326

Epoch 00002: val_loss improved from 3.71329 to 3.71285, saving model to models/best_1.h5
Epoch 3/4
2/2 [=====] - 17s 8s/step - loss: 3.7122 - acc: 0.0616 - val_loss: 3.712
4 - val_acc: 0.0326

Epoch 00003: val_loss improved from 3.71285 to 3.71243, saving model to models/best_1.h5
Epoch 4/4
2/2 [=====] - 17s 9s/step - loss: 3.7117 - acc: 0.0438 - val_loss: 3.712
0 - val_acc: 0.0326

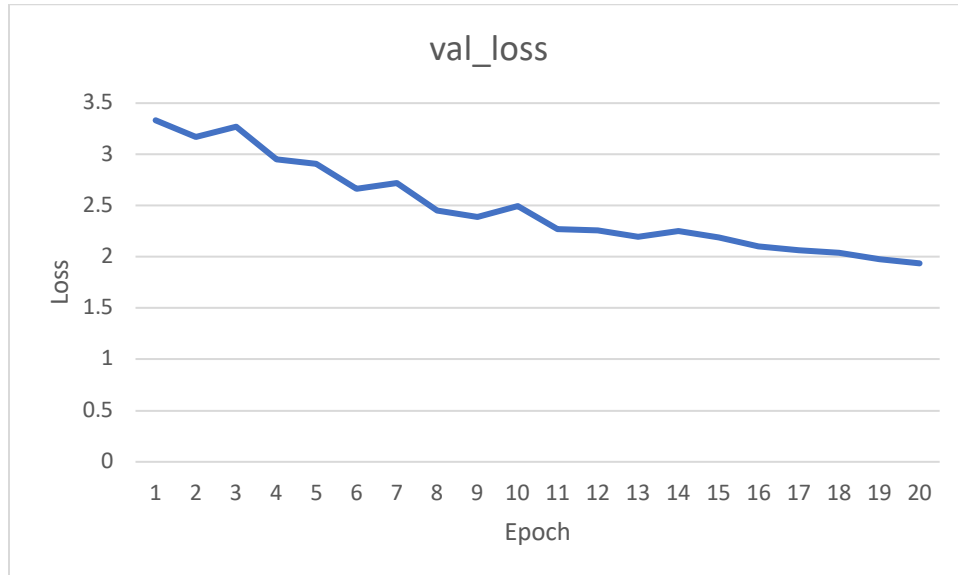
Epoch 00004: val_loss improved from 3.71243 to 3.71198, saving model to models/best_1.h5

```

However during complete run , Validation loss keeps reducing with epoch. So number of epochs has to be increased till validation loss start increasing (i.e, until overfitting start occurs) .

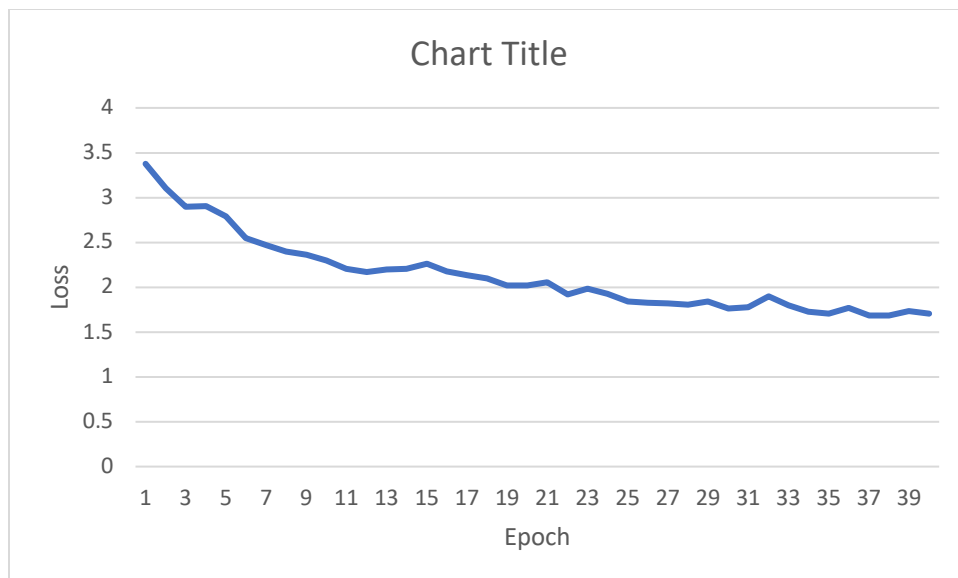
of Audio samples : 9400

config = Config(sampling_rate=16000, audio_duration=2, n_folds=2, max_epochs=20)



of Audio samples : 9400

config = Config(sampling_rate=16000, audio_duration=2, n_folds=2, max_epochs=40)



Log files used for above graphs are stored at `logs/training-history.txt`

IV . Results



Model Evaluation and Validation

Model has been evaluated with Average Precision@3 , since APK@3 considers next 3 predicted values , I prefer this than other accuracy measures .

Epoch	Accuracy on Training set	Accuracy on Test set
20	52%	53%
40	66%	66%

Since accuracy on test set is closer to accuracy on training set the model generalizes well.

<https://www.kaggle.com/c/freesound-audio-tagging/leaderboard>

156	new	Kannappan Natarasan		0.531	1	9h
147	new	Kannappan Natarasan		0.660	2	2h

Here is the sample output

Predicted	Actual	Predicted-next-3-values
Acoustic_guitar	Flute	[Acoustic_guitar Laughter Cello]
Tearing	Tearing	[Tearing Keys_jangling Computer_keyboard]
Fireworks	Writing	[Fireworks Writing Knock]

Cello	Clarinet	[Cello Acoustic_guitar Gong]
Cough	Cough	[Cough Laughter Bark]
Acoustic_guitar	Double_bass	[Acoustic_guitar Cello Double_bass]
Cello	Cello	[Cello Double_bass Acoustic_guitar]
Hi-hat	Tambourine	[Hi-hat Tambourine Shatter]
Harmonica	Harmonica	[Harmonica Violin_or_fiddle Gong]
Cello	Clarinet	[Cello Clarinet Flute]

Justification

After multiple iterations the test accuracy is reached to 66%, which is very close to my benchmark 70%. With the explorations so far I'm confident that I can increase accuracy, by using GPU or TPU instances from cloud and increase learning iterations with shorter time with less money spent.

The constraints are lack of computing power and number of audio samples .

V . Conclusion

Free-form visualization (listening)

After built the model tested for accuracy , I wanted to test how the prediction works, since the model is to predict sound files “Free-form listening” could be appropriate title for this section. When I tried to predict a random sound file (*37965099.wav*) from test set , the model predicted it as clarinet, next two predictions are Flute and Saxophone.

Free form visualization

This section would provide room to select a data set randomly and predict it using the model built, then gives a provision to play it


```

In [280]: 1 test=pd.read_csv(home+'doc/sample_submission.csv')

In [320]: 1 j=2000
          2 myaudio=home+'data/audio_test/'+test['fname'][j]

In [321]: 1 pred3=getPred3(predictions_test[j])
          2 print('Predicted      :',LABELS[predictions_test[j].argmax()])
          3 print('Predicted next 3 vals:',end='')
          4 print('[',end='')
          5 for i in range(3):
          6     print(LABELS[pred3[i]],end=' ')
          7 print(']')

Predicted      : Clarinet
Predicted next 3 vals:[Clarinet Flute Saxophone ]

In [322]: 1 ipd.Audio(myaudio)
Out[322]: 

```

A mini video demonstration is attached to following link, which depicts above steps
<https://github.com/knatarasan/freesound/blob/master/doc/Visualization.mp4>

But when I played the sound (refer the video from github to listen the sound) it sounds close to clarinet, since the model works with 66% accuracy it makes sense to understand that the model predicted it as clarinet.

Reflection

To recap steps I've followed.

- Prepared dataset using packages like librosa, pyaudio and scipy which are very helpful to process sound files and convert them to single dimension array.
- Prepared plots to using numpy and matplotlib packages, which helped to understand overall dataset. Using the package librosa plotted single audio file.
- Standardized dataset by taking equal length of data sample only, for the smaller sound files padded with 0, and normalized the data set.
- Built CNN architecture using Convolution1D, MaxPool, Dropout and Dense layers.
- Executed learning process by splitting dataset into multiple folds.
- Verified model through few sample predictions .
- Able to derive a model with AP@3 accuracy of 66% .

Most difficult thing in this project is execution time, A complete run with 9400 audio files on 40 epoch would take 4 hours. But it wasn't easy to decide the optimum configuration (audio file size, activation function, optimizer , number of filters etc) with out few test runs. I've tried few sampling techniques like testing a configuration with 2000 files , which also took significant time. So worked with 200 files. Having better understanding on how sound is stored and replayed and ability to cut the piece of sound which has most relevant features for prediction could have solved two problems, a) how to test a configuration quickly b) Increasing model efficiency by making the learning algorithm to focus on sound clip which has most relevant features .

This is very interesting project, which helped me to explore the details of CNN, how to build the architecture, select filter size, slider padding, activation function , optimizer. Optimum epoch count could be derived by running the learning multiple times and observations from change in validation loss.

Improvement

As per current trend on reduction in validation loss, yet the model is in under fit stage, If I've enough computing power I could easily reached the sweet spot by increasing epochs.

Also exploring mechanics of sound would help to cut window of sound file which has the relevant pattern ,i.e, the current selection picks first 2 seconds of audio, suppose the actual recording starts after 2 seconds , required features would be missed.

Reference:

<https://adeshpande3.github.io/A-Beginner%27s-Guide-To-Understanding-Convolutional-Neural-Networks/>