

学習モデル、Data Augmentation の有無、 最適化アルゴリズムの違いによる音響シーン分類性能比較

伊藤 葵

Aoi Ito

法政大学情報科学部デジタルメディア学科 B3 20K1105

aoi.ito.8q@stu.hosei.ac.jp

2023 年 1 月 26 日

1 はじめに

音響シーン・音響イベントの分類・検出に関する技術・研究分野 (Detection and Classification of Acoustic Scenes and Events) では、DCASE Community という組織が音響に関する問題に取り組む DCASE Challenge/Workshop を 2013 年から開催している [1]。この Challenge の一つに、音響シーン分類がある。

このレポートでは、DCASE2013 の Task1 にある音響シーン分類問題のデータセットを縮小したものに対し、異なる学習モデル (CNN, wavelet scattering を用いたアンサンブル学習、CNN とアンサンブル学習の Late Fusion) 間や Data Augmentation の有無、最適化アルゴリズムの違いという 3 つの観点から性能を比較する。

2 手法

この章では、今回の実験で用いたモデルの学習方法と Data Augmentation の方法を説明する。また、このレポートでは、<https://jp.mathworks.com/help/audio/ug/acoustic-scene-recognition-using-late-fusion.html> を参考とした。

2.1 CNN

CNN(Convolutional Neural Network) とは、ニューラルネットワークに畳み込みを追加したものである [2]。本レポートでは、図 1, 2 のネットワークで CNN を行う。最適化アルゴリズムは、sgdm (モーメント項付き確率的勾配降下法) と adam で比較する。



図 1 CNN のネットワーク 1

①の続き

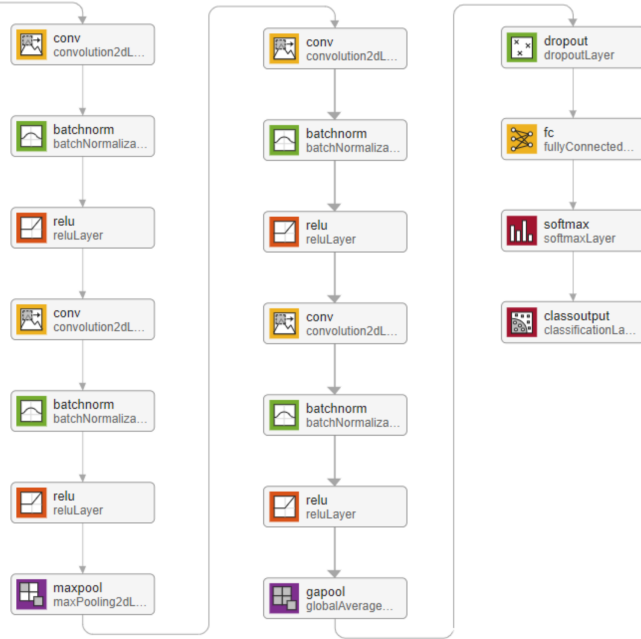


図2 CNN のネットワーク 2

CNN の入力信号 (特徴量) は、対数をとったメルスペクトログラムとする。

2.2 wavelet scattering を用いたアンサンブル学習

wavelet 変換 とは、周波数解析の手法の一つである。

フーリエ解析では、ある入力信号を拡大縮小したサイン波、コサイン波の足し合わせで表現するが、時系列情報が損失する。これに対し、wavelet 変換は時間と周波数を同時に解析する手法である。基準となる小さな波 mother wavelet を様々な縮尺に引き伸ばし、周波数の物差しであるウェーブレットを多数用意する。これらのウェーブレットを、時間軸方向に平行移動させながら入力信号にあてがい、時間と周波数の情報を同時に得るのである [3]。[4] は wavelet scattering が音響シーンの表現に優れていることを示しているため、本レポートでは wavelet scattering を採用した。

アンサンブル学習とは、個々に別々の学習器として学習させたもののそれぞれの結果を融合させることで、未学習のデータに対しての予測能力向上を狙う学習である [5]。複数のモデルを学習し、最後に各モデルの予測結果を多数決原理で決めることで、各モデルの認識精度が低くても性能が向上する確率が高くなるのがメリットである。

2.3 Late Fusion による出力の統合

Late Fusion モデルとは、複数のモデルの出力直後で出力結果を統合したものである。Early Fusion と異なり、統合後にニューラルネットワークの影響を受けない。本レポートでは、CNN の出力と wavelet scattering を用いたアンサンブル学習の出力を統合している。ブロック図で示すと、図 3 の通り [6]。

②

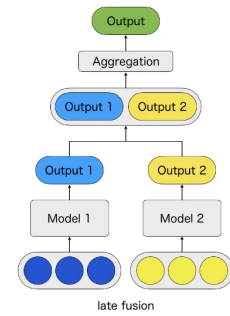


図3 Late Fusion のブロック図

2.4 CNN 用の Data Augmentation

本レポートの実験では、学習用データセットとなる音声ファイル数が全部で 56 個と少ないため、学習用データ数の増加を目的に Data Augmentation を行う。本レポートでは、2 つの異なる音響シーンを 1:1 の割合で混ぜ合わせた。次式の x は音声ファイル、 y は対応する正解ラベルであり、混合比は λ で調整可能である。今回の研究では、混合比 λ は 0.5 とした。実装では、異なるラベルを持つスペクトログラム同士を混合している。

$$\tilde{x} = \lambda x_i + (1 - \lambda) x_j \quad (1)$$

$$\tilde{y} = \lambda y_i + (1 - \lambda) y_j \quad (2)$$

2.5 実験に使用したデータ

実験に使用したデータは、DCASE2013 Task1 のデータセット (https://archive.org/details/dcase2013_scene_classification) を縮小し用いた。音響シーンの種類は、全部で 7 種類 (bus, busystreet, office, park, quietstreet, restaurant, supermarket) とする。各音響シーンは 10 個の wav ファイルが用意されており、音声の長さは 30 秒である。この内、各音響シーン No.01~08 の音声ファイルを学習用データセット、No.09~10 の音声ファイルを検証用データセットとした。

3 実験

3.1 実験条件

実行環境は、Intel(R) Core(TM) i7-10510U CPU @ 1.80GHz 2.30 GHz、RAM 16.0 GB、使用 OS は Windows である。言語は MATLAB を使用した。

実験条件 (学習モデル、DataAugmentation の有無、最適化アルゴリズム) は表 1 の通りである。

3.2 実験方法

実験方法は以下の通りである。

1. 音響シーンデータセットのロード
2. 学習用データセット、検証用データセットへの分類
3. CNN の特徴量抽出
4. Data Augmentation
5. CNN の各層の定義、学習
6. CNN の評価
7. アンサンブル学習
8. アンサンブル学習の評価
9. Late Fusion の適用、評価

手順 1,2 のデータセットの内訳は、2.5 で述べたとおりであ

表 1 実験条件

Model	DataAugmentation	optimization
CNN		sgdm
Ensemble		sgdm
Late Fusion		sgdm
CNN	✓	sgdm
Ensemble	✓	sgdm
Late Fusion	✓	sgdm
CNN		adam
Ensemble		adam
Late Fusion		adam
CNN	✓	adam
Ensemble	✓	adam
Late Fusion	✓	adam

る。手順 3 では、音声の特徴量としてメルスペクトログラムを採用した。手順 6, 8, 9 の評価時には、Average Accuracy だけでなく、混同行列を用いてどの音響シーンの分類を誤っているか可視化した。

3.3 実験結果

各モデルの正解率は表 2 の通り。

表 2 各モデルの正解率 (DA: Data Augmentation)

Model	sgdm	sgdm&DA	adam	adam&DA
CNN	57.14	64.29	71.43	50.00
Ensemble	78.57	78.57	85.71	78.57
Late Fusion	78.57	78.57	85.71	71.43

また、最適化アルゴリズム (sgdm, adam) における各モデル別の混同行列は図 4~13 の通り。尚、アンサンブル学習は各最適化アルゴリズムによる学習方法の違いの影響はないため、DataAugmentation の有無にのみ触れている。

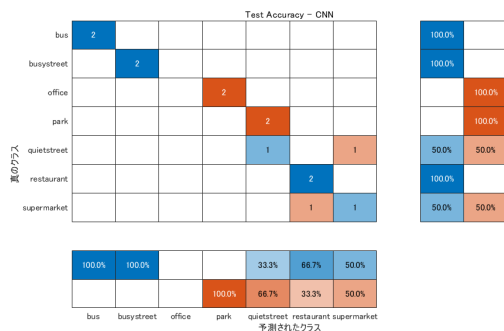


図 4 CNN の混同行列 (sgdm)

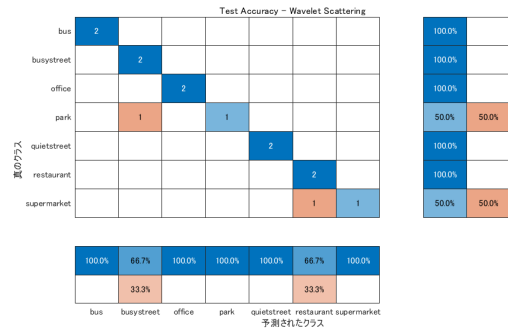


図 5 wavelet scattering を用いたアンサンブル学習の混同行列 (DataAugmentation なし)

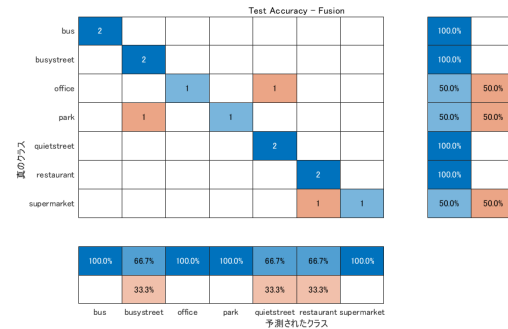


図 6 Late Fusion の混同行列 (sgdm)

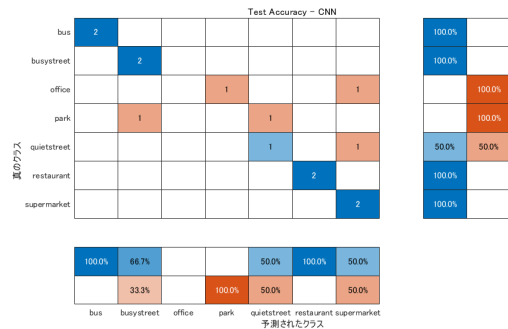


図 7 CNN の混同行列 (sgdm, DataAugmentation)

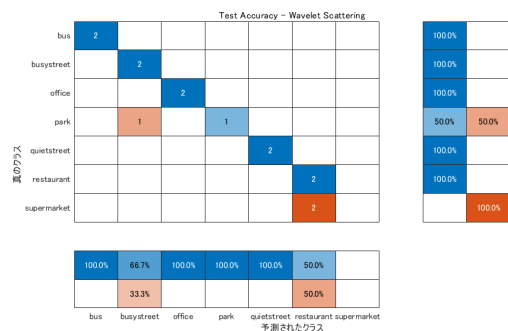


図 8 wavelet scattering を用いたアンサンブル学習の混同行列 (DataAugmentation あり)

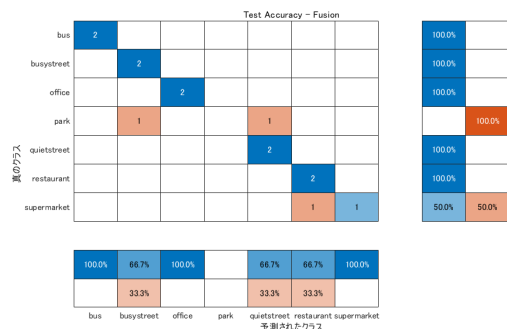


図 9 Late Fusion の混同行列 (sgdm, DataAugmentation)

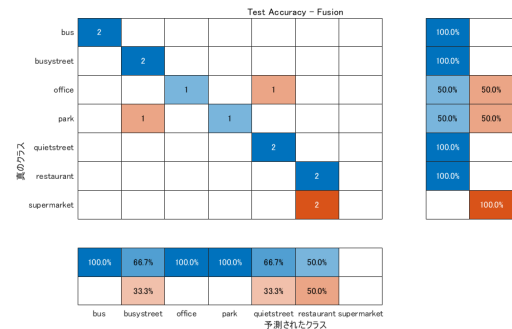


図 13 Late Fusion の混同行列 (adam, DataAugmentation)

最も正解率が低かったモデルは、CNN のみのモデル (Data Augmentation 済、adam 使用) である。最も正解率が高かったモデルは、wavelet scattering を用いたアンサンブル学習のモデルとこのモデルに CNN (adam 使用) の結果を Late Fusion したモデルである。CNN のみのモデルには、正解率が 0% となる音響シーンが存在する。

4 考察

まず、Data Augmentation の有無が認識精度に与える影響について考察する。最も正解率が低かった CNN のみのモデル (Data Augmentation 済、adam 使用) に着目する。一般に、Data Augmentation は、学習データ数の増加や過学習の抑制に効果がある。しかし、今回の実験では Data Augmentation をしても CNN のみ (最適化アルゴリズム: adam) の学習モデルと比較すると、21.43% 正解率が減少した。今回の実験の学習データセットの数は 56 と少なく、かつ人間が聞いても区別がしにくいほど似た音響シーンの組み合わせが多いデータセットであった。そのため、Data Augmentation を異なるラベルを持つスペクトログラム同士で行っても、似たようなデータがさらに増加してしまったため、過学習の促進に繋がってしまった。反面、最適化アルゴリズムに sgdm を使用した際の実験では、CNN のみの学習モデルで 7.15% と正解率が高くなっている。

次に、学習モデルの違いの観点から考察する。CNN だけのモデルに対し、wavelet scattering を用いたアンサンブル学習の出力を CNN の出力に統合した Late Fusion モデルの方が認識性能が向上している。このことから、今回の実験のように、学習データ数が少ない場合は、特に複数のモデルを組み合わせた方が、単体のモデルによる音響シーン分類よりも、はるかに性能向上の効果が期待できると考える。例えば、学習データ数が少なく、似たデータが多かったために Data Augmentation が上手いかなかったと考察した CNN のみのモデル (Data Augmentation 済、adam 使用) を見ても、Late Fusion をすることで、Data Augmentation をしなかった CNN のみ (最適化アルゴリズム: adam) の学習モデルと同等の性能となる。

最後に、混同行列に着目する。図 5,8 を比較すると、Data Augmentation をすることにより、supermarket の正解率が高くなった。これは、Data Augmentation による学習データ数増加が理由だと考察する。特に、supermarket は店内で鳴っている音の種類も多いため、他の音響シーンの音声が入り混じることによって、supermarket の未知のデータに含まれる音に柔軟に対応できるようになった。そして、図 4 ~ 13 の混同行列

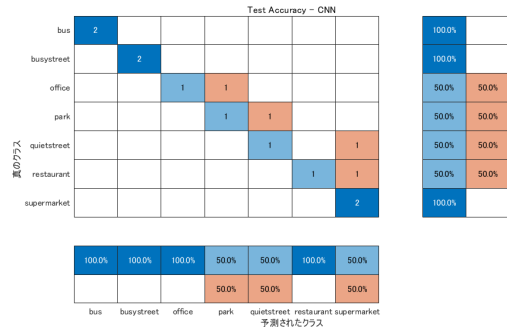


図 10 CNN の混同行列 (adam)

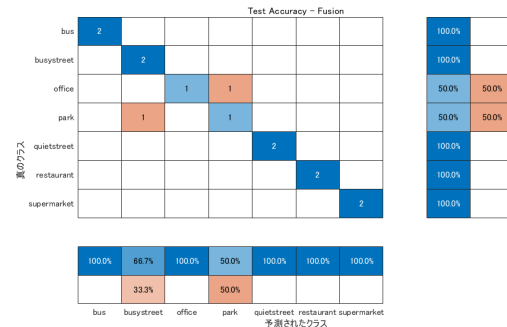


図 11 Late Fusion の混同行列 (adam)

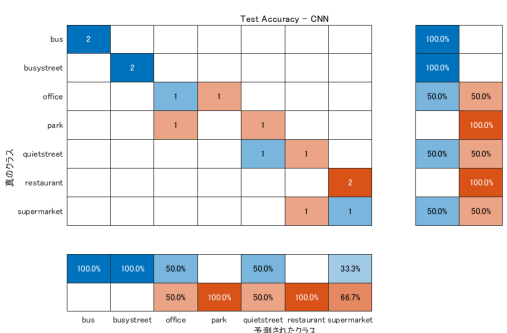


図 12 CNN の混同行列 (adam, DataAugmentation)

を通して、park の認識精度が低い。これに関して、図 4 と 6、図 7 と 9、図 10 と 11、図 12 と 13 を比較すると、複数の学習モデルの統合が認識精度が低い音響シーンの分類精度向上に効果があるといえる。さらに、図 4 と 10、図 6 と 11 といった最適化アルゴリズムが異なる CNN の混同行列を比較すると、adam を用いたモデルの方が park や office といった正解率が 0% と低かった音響シーンの認識精度が向上している。この結果から、音響シーン分類のタスクにおいて、最適化アルゴリズムは adam の方がふさわしい。ただし、adam は sgdm と比較すると汎化性能が劣る [7] といわれるため、AdaBound や AMSBound との比較が今後必要である。

5 おわりに

このレポートでは、DCASE2013 Task1 のデータセットを用いて、異なる学習モデル (CNN, wavelet scattering を用いたアンサンブル学習、CNN とアンサンブル学習の Late Fusion) 間や Data Augmentation の有無、最適化アルゴリズムの違いという 3 つの観点から音響シーン分類の性能を比較した。CNN だけでなく、アンサンブル学習の出力を Late Fusion にて統合し、最適化アルゴリズムには adam を採用することで、音響シーンの分類性能向上が期待できる結果となった。Data Augmentation は、元の学習データ数が少なく似たデータが多い場合は、どのようにデータを組み合わせる新たな学習データを作成するか、方法を検討する必要がある。課題として、CNN に対して、CNN 以外のモデル (LSTM 等) と wavelet scattering を用いたアンサンブル学習の Late Fusion の性能比較がある。

参考文献

- [1] @diesekiefer. "音環境理解のための技術コミュニティ DCASE の紹介". Qiita.2022-12-15. <https://qiita.com/diesekiefer/items/bd0c31c135ede8a9990e>, (参照 2023-01-10).
- [2] @icofog417. "Convolutional Neural Network とは何か". Qiita.2019-10-27. <https://qiita.com/icofog417/items/5fd55fad152231d706c2>, (参照 2023-01-10).
- [3] 三谷 政昭 (2008). "やり直しのための通信数学【オンデマンド版】". CQ 出版社. (参照 2023-01-10).
- [4] Vincent Lostanlen, Joakim Anden(2016). "BINAURAL SCENE CLASSIFICATION WITH WAVELET SCATTERING". Detection and Classification of Acoustic Scenes and Events 2016.2016-09-03. (参照 2023-01-10).
- [5] "アンサンブル学習 (Ensemble learning) 解説と実験". S-Analysis. [https://data-analysis-stats.jp/機械学習/アンサンブル学習 \(Ensemble learning\) 解説と実験](https://data-analysis-stats.jp/機械学習/アンサンブル学習 (Ensemble learning) 解説と実験), (参照 2023-01-10).
- [6] "ニューラルネットワークを用いた複数モーダル の 統 合 に つ い て". AGIRobots.2022-07-09. <https://agirobots.com/multimodal-ai-fusion/>, (参照 2023-01-10).
- [7] @Phoebooooo. "[最新論文] 新しい最適化手法誕生! AdaBound & AMSBound". Qiita. 2019-03-04.

<https://qiita.com/Phoebooooo/items/f610affdcaaae0a28f34>, (参照 2023-01-10).

付録

作成したソースコードを記載する。

最適化アルゴリズムは、ソースコード 333 行目の options にて 'adam' もしくは 'sgdm' を選択する。'sgdm' を選択した場合は、'Momentum',0.9,'も追記する。

```
1 %% Acoustic Scene Recognition Using Late Fusion
2
3 %% Load Acoustic Scene Recognition Data Set
4 % To run the example, you must first download ...
   the data set [1]. The full
5 % data set is approximately 15.5 GB. Depending ...
   on your machine and internet
6 % connection, downloading the data can take ...
   about 4 hours.
7
8 downloadFolder = tempdir;
9 datasetFolder = ...
   fullfile(downloadFolder,'scenes_stereo');
10
11 %%
12 % Read in the development set metadata as a ...
   table. Name the table variables
13 % |FileName|, |AcousticScene|.
14
15 metadata_train = readtable("dataset.txt");
16 head(metadata_train)
17
18 metadata_test = readtable("dataset_test.txt");
19 head(metadata_test)
20
21 train_filePaths = metadata_train.FileName;
22 test_filePaths = metadata_test.FileName;
23
24 %%
25 % Create audio datastores for the train and ...
   test sets. Set the
26 % |Labels| property of the
27 % ...
   <docid:audio_ref#mw-6315b106-9a7b-4a11-a7c6-322c073e343a
28 % |audioDatastore|> to the acoustic scene. Call
29 % ...
   <docid:audio_ref#mw-27293e0c-5066-45c4-b34d-e2814c56921f
30 % |countEachLabel|> to verify an even ...
   distribution of labels in both the
31 % train and test sets.
32
33 adsTrain = audioDatastore(train_filePaths, ...
34   'Labels',categorical(metadata_train.AcousticScene), ...
   ...
35   'IncludeSubfolders',true);
36 display(countEachLabel(adsTrain))
37
38 adsTest = audioDatastore(test_filePaths, ...
39   'Labels',categorical(metadata_test.AcousticScene), ...
   ...
40   'IncludeSubfolders',true);
41 display(countEachLabel(adsTest))
42
43 %%
44 % You can reduce the data set used in this ...
   example to speed up the run time
45 % at the cost of performance. In general, ...
   reducing the data set is a good
46 % practice for development and debugging. Set ...
   |reduceDataset| to |true| to
```

```

47 % reduce the data set.
48 reduceDataset = false;
49 if reduceDataset
50     adsTrain = splitEachLabel(adsTrain,20);
51     adsTest = splitEachLabel(adsTest,10);
52 end
53
54 %%
55 % Call ...
    <docid:audio_ref#mw_4257f192-1344-4f89-8edb-414b0f9efc42>
    |read|> to
56 % get the data and sample rate of a file from ...
    the train set. Audio in the
57 % database has consistent sample rate and ...
    duration. Normalize the audio and
58 % listen to it. Display the corresponding label.
59
60 [data,adsInfo] = read(adsTrain);
61 data = data./max(data,[],'all');
62
63 fs = adsInfo.SampleRate;
64 sound(data,fs)
65
66 fprintf('Acoustic scene = ...
    %s\n',adsTrain.Labels(1))
67
68 %%
69 % Call ...
    <docid:audio_ref#mw_a427bc3f-f73e-448b-896a-65ddbb3a26f1>
    |reset|> to
70 % return the datastore to its initial condition.
71
72 reset(adsTrain)
73
74 %% Feature Extraction for CNN
75 % Each audio clip in the dataset consists of ...
    10 seconds of stereo
76 % (left-right) audio. The feature extraction ...
    pipeline and the CNN
77 % architecture in this example are based on ...
    [3]. Hyperparameters for the
78 % feature extraction, the CNN architecture, ...
    and the training options were
79 % modified from the original paper using a ...
    systematic hyperparameter
80 % optimization workflow.
81 %
82 % First, convert the audio to mid-side ...
    encoding. [3] suggests that mid-side
83 % encoded data provides better spatial ...
    information that the CNN can use to
84 % identify moving sources (such as a train ...
    moving across an acoustic
85 % scene).
86
87 dataMidSide = [sum(data,2),data(:,1)-data(:,2)];
88
89 %%
90 % Divide the signal into one-second segments ...
    with overlap. The final
91 % system uses a probability-weighted average ...
    on the one-second segments to
92 % predict the scene for each 10-second audio ...
    clip in the test set. Dividing
93 % the audio clips into one-second segments ...
    makes the network easier to
94 % train and helps prevent overfitting to ...
    specific acoustic events in the
95 % training set. The overlap helps to ensure ...
    all combinations of features
96 % relative to one another are captured by the ...
    training data. It also
97 % provides the system with additional data ...
    that can be mixed uniquely
    % during augmentation.
98
99
100 segmentLength = 1;
101 segmentOverlap = 0.5;
102
103 [dataBufferedMid,~] = ...
    buffer(dataMidSide(:,1),round(segmentLength*fs),round(segmentLength*fs));
104 [dataBufferedSide,~] = ...
    .buffer(dataMidSide(:,2),round(segmentLength*fs),round(segmentLength*fs));
105 dataBuffered = ...
    zeros(size(dataBufferedMid,1),size(dataBufferedMid,2)+size(dataBufferedSide,2));
106 dataBuffered(:,1:2:end) = dataBufferedMid;
107 dataBuffered(:,2:2:end) = dataBufferedSide;
108
109 %%
110 % Use ...
    <docid:audio_ref#mw_1bb316de-8018-4365-a351-73374473237f>
111 % |melSpectrogram|> to transform the data into ...
    a compact frequency-domain
112 % representation. Define parameters for the ...
    mel spectrogram as suggested by
113 % [3].
114
115 windowLength = 2048;
116 samplesPerHop = 1024;
117 samplesOverlap = windowLength - samplesPerHop;
118 fftLength = 2*windowLength;
119 numBands = 128;
120
121 %%
122 % |melSpectrogram| operates along channels ...
    independently. To optimize
123 % processing time, call |melSpectrogram| with ...
    the entire buffered signal.
124
125 spec = melSpectrogram(dataBuffered,fs, ...
    'Window',hamming(windowLength,'periodic'), ...
    'OverlapLength',samplesOverlap, ...
    'FFTLength',fftLength, ...
    'NumBands',numBands);
126
127
128
129
130
131 %%
132 % Convert the mel spectrogram into the ...
    logarithmic scale.
133
134 spec = log10(spec+eps);
135
136 %%
137 % Reshape the array to dimensions (Number of ...
    bands)-by-(Number of
138 % hops)-by-(Number of channels)-by-(Number of ...
    segments). When you feed an
139 % image into a neural network, the first two ...
    dimensions are the height and
140 % width of the image, the third dimension is ...
    the channels, and the fourth
141 % dimension separates the individual images.
142
143 X = ...
    reshape(spec,size(spec,1),size(spec,2),size(data,2),[]);
144
145 %%
146 % Call |melSpectrogram| without output ...
    arguments to plot the mel
147 % spectrogram of the mid channel for the first ...
    six of the one-second
148 % increments.
149
150 for channel = 1:2:11
151     figure
152     melSpectrogram(dataBuffered(:,channel),fs, ...
    'Window',hamming(windowLength,'periodic'), ...

```



```

154         'OverlapLength',samplesOverlap, ...
155         'FFTLength',fftLength, ...
156         'NumBands',numBands);
157     title(sprintf('Segment %d',ceil(channel/2)))
158 end
159
160 %%
161 % The helper function ...
162 %   |HelperSegmentedMelSpectrograms| performs ...
163 %   the feature
164 % extraction steps outlined above.
165
166 %%
167 % To speed up processing, extract mel ...
168 % spectrograms of all audio files in
169 % the datastores using ...
170 % <docid:matlab.ref#bvaomuj-1 |tall|> ...
171 % arrays. Unlike
172 % in-memory arrays, tall arrays remain ...
173 % unevaluated until you request that
174 % the calculations be performed using the ...
175 % <docid:matlab.ref#bvaolov
176 % |gather|> function. This deferred evaluation ...
177 % enables you to work quickly
178 % with large data sets. When you eventually ...
179 % request the output using
180 % |gather|, MATLAB combines the queued ...
181 % calculations where possible and
182 % takes the minimum number of passes through ...
183 % the data. If you have Parallel
184 % Computing Toolbox(TM), you can use tall ...
185 % arrays in your local MATLAB
186 % session, or on a local parallel pool. You ...
187 % can also run tall array
188 % calculations on a cluster if you have ...
189 % MATLAB(R) Parallel Server(TM)
190 % installed.
191 %
192 % If you do not have Parallel Computing ...
193 % Toolbox(TM), the code in this
194 % example still runs.
195
196 train_set_tall = tall(adsTrain);
197 xTrain = ...
198     cellfun(@ (x) HelperSegmentedMelSpectrograms(x,fs, ...
199         'SegmentLength',segmentLength, ...
200         'SegmentOverlap',segmentOverlap, ...
201         'WindowLength',windowLength, ...
202         'HopLength',samplesPerHop, ...
203         'NumBands',numBands, ...
204         'FFTLength',fftLength), ...
205         train_set_tall, ...
206         'UniformOutput',false);
207
208 test_set_tall = tall(adsTest);
209 xTest = ...
210     cellfun(@ (x) HelperSegmentedMelSpectrograms(x,fs, ...
211         'SegmentLength',segmentLength, ...
212         'SegmentOverlap',segmentOverlap, ...
213         'WindowLength',windowLength, ...
214         'HopLength',samplesPerHop, ...
215         'NumBands',numBands, ...
216         'FFTLength',fftLength), ...
217         test_set_tall, ...
218         'UniformOutput',false);
219
220 xTest = cat(4,xTest{:});
221
222 %%
223 % Replicate the labels of the training set so ...
224 % that they are in one-to-one
225 % correspondence with the segments.
226
227 numSegmentsPer10seconds = size(dataBuffered,2)/2;
228 yTrain = ...
229     repmat(adsTrain.Labels,1,numSegmentsPer10seconds)';
230 yTrain = yTrain(:);
231
232 %% Data Augmentation for CNN
233 % The DCASE 2017 dataset contains a relatively ...
234 % small number of acoustic
235 % recordings for the task, and the development ...
236 % set and evaluation set were
237 % recorded at different specific locations. As ...
238 % a result, it is easy to
239 % overfit to the data during training. One ...
240 % popular method to reduce
241 % overfitting is mixup. In mixup, you ...
242 % augment your dataset by mixing the
243 % features of two different classes. When you ...
244 % mix the features, you mix the
245 % labels in equal proportion. That is:
246
247 
$$\begin{array}{l} \tilde{\mathbf{x}} = \lambda \mathbf{x}_i + (1-\lambda) \mathbf{x}_j \\ \tilde{\mathbf{y}} = \lambda \mathbf{y}_i + (1-\lambda) \mathbf{y}_j \end{array}$$

248
249 % Mixup was reformulated by [2] as labels ...
250 % drawn from a probability
251 % distribution instead of mixed labels. The ...
252 % implementation of mixup in this
253 % example is a simplified version of mixup: ...
254 % each spectrogram is mixed with
255 % a spectrogram of a different label with ...
256 % lambda set to 0.5. The
257 % original and mixed datasets are combined for ...
258 % training.
259
260 % {
261 xTrainExtra = xTrain;
262 yTrainExtra = yTrain;
263 lambda = 0.5;
264 for i = 1:size(xTrain,4)
265
266     % Find all available spectrograms with ...
267     % different labels.
268     availableSpectrograms = ...
269         find(yTrain~=yTrain(i));
270
271     % Randomly choose one of the available ...
272     % spectrograms with a different label.
273     numAvailableSpectrograms = ...
274         numel(availableSpectrograms);
275     idx = randi([1,numAvailableSpectrograms]);
276
277     % Mix.
278     xTrainExtra(:, :, :, i) = ...
279         lambda*xTrain(:, :, :, i) + ...
280         (1-lambda)*xTrain(:, :, :, availableSpectrograms(idx));
281
282     % Specify the label as randomly set by lambda.
283     if rand > lambda
284         yTrainExtra(i) = ...
285             yTrain(availableSpectrograms(idx));
286     end
287 end
288 xTrain = cat(4,xTrain,xTrainExtra);

```

```

258 yTrain = [yTrain;yTrainExtra];
259
260 %%
261 % Call |summary| to display the distribution ...
    of labels for the augmented
262 % training set.
263
264 summary(yTrain)
265 %}
266
267 %% Define and Train CNN
268 %
269 % Define the CNN architecture. This ...
    architecture is based on [1] and
270 % modified through trial and error. See
271 % ...
    <docid:nnet_ug#mw.25a9e4c2-614f-48b9-97e2-bbdd7aaf936f ...
    List of Deep
272 % Learning Layers> to learn more about deep ...
    learning layers available in
273 % MATLAB(R).
274
275 imgSize = ...
    [size(xTrain,1),size(xTrain,2),size(xTrain,3)];
276 numF = 32;
277 layers = [ ...
278     imageInputLayer(imgSize)
279
280     batchNormalizationLayer
281
282     convolution2dLayer(3,numF,'Padding','same')
283     batchNormalizationLayer
284     reluLayer
285     convolution2dLayer(3,numF,'Padding','same')
286     batchNormalizationLayer
287     reluLayer
288
289     maxPooling2dLayer(3,'Stride',2,'Padding','same')
290
291     convolution2dLayer(3,2*numF,'Padding','same')
292     batchNormalizationLayer
293     reluLayer
294     convolution2dLayer(3,2*numF,'Padding','same')
295     batchNormalizationLayer
296     reluLayer
297
298     maxPooling2dLayer(3,'Stride',2,'Padding','same')
299
300     convolution2dLayer(3,4*numF,'Padding','same')
301     batchNormalizationLayer
302     reluLayer
303     convolution2dLayer(3,4*numF,'Padding','same')
304     batchNormalizationLayer
305     reluLayer
306
307     maxPooling2dLayer(3,'Stride',2,'Padding','same')
308
309     convolution2dLayer(3,8*numF,'Padding','same')
310     batchNormalizationLayer
311     reluLayer
312     convolution2dLayer(3,8*numF,'Padding','same')
313     batchNormalizationLayer
314     reluLayer
315
316     globalAveragePooling2dLayer
317
318     dropoutLayer(0.5)
319
320     fullyConnectedLayer(7)
321     softmaxLayer
322     classificationLayer];
323
324

```

```

325 %%
326 % Define <docid:nnet_ref#bu59f0q ...
    |trainOptions|> for the CNN. These
327 % options are based on [3] and modified ...
    through a systematic hyperparameter
328 % optimization workflow.
329
330 miniBatchSize = 128;
331 tuneme = 128;
332 lr = 0.05*miniBatchSize/tuneme;
333 options = trainingOptions('sgdm', ...
334     'InitialLearnRate',lr, ...
335     'MiniBatchSize',miniBatchSize, ...
336     'L2Regularization',0.005, ...
337     'MaxEpochs',8, ...
338     'Momentum',0.9, ...
339     'Shuffle','every-epoch', ...
340     'Plots','training-progress', ...
341     'Verbose',false, ...
342     'LearnRateSchedule','piecewise', ...
343     'LearnRateDropPeriod',2, ...
344     'LearnRateDropFactor',0.2);
345
346 %%
347 % Call <docid:nnet_ref#bu6sn4c |trainNetwork|> ...
    to train the network.
348
349 trainedNet = ...
    trainNetwork(xTrain,yTrain,layers,options);
350
351 %% Evaluate CNN
352 % Call ...
    <docid:nnet_ref#mw.0a51db93-cccf-4b2f-ae4c-6724cbf5ec46 ...
    |predict|>
353 % to predict responses from the trained ...
    network using the held-out test
354 % set.
355
356 cnnResponsesPerSegment = ...
    predict(trainedNet,xTest);
357
358 %%
359 % Average the responses over each 10-second ...
    audio clip.
360
361 classes = trainedNet.Layers(end).Classes;
362 numFiles = numel(adsTest.Files);
363
364 counter = 1;
365 cnnResponses = zeros(numFiles,numel(classes));
366 for channel = 1:numFiles
367     cnnResponses(channel,:) = ...
        sum(cnnResponsesPerSegment(counter:counter+numSegmentsPer10seconds-1),2);
368     counter = counter + numSegmentsPer10seconds;
369 end
370
371 %%
372 % For each 10-second audio clip, choose the ...
    maximum of the predictions,
373 % then map it to the corresponding predicted ...
    location.
374
375 [~,classIdx] = max(cnnResponses,[],2);
376 cnnPredictedLabels = classes(classIdx);
377
378 %%
379 % Call ...
    <docid:nnet_ref#mw.b571feea-6af0-489f-b52c-a4b5141ac550
380 % |confusionchart|> to visualize the accuracy ...
    on the test set. Return the
381 % average accuracy to the Command Window.
382
383 figure('Units','normalized','Position',[0.2 ...

```



```

0.2 0.5 0.5])
384 cm = ...
    confusionchart(adsTest.Labels,cnnPredictedLabels,'title','Test
    Accuracy - CNN');
385 cm.ColumnSummary = 'column-normalized';
386 cm.RowSummary = 'row-normalized';
387
388 fprintf('Average accuracy of CNN = ...
    %0.2f\n',mean(adsTest.Labels==cnnPredictedLabels)*100)
389
390
391 %% Feature Extraction for Ensemble Classifier
392 % Wavelet scattering has been shown in [4] to ...
    provide a good representation
393 % of acoustic scenes. Define a
394 % ...
    <docid:wavelet.ref#mw.a0dd2386-66f5-4c0b-b96d-7ae0ba327cd1corresponding predicted location. Call
395 % |waveletScattering|> object. The invariance ...
    scale and quality factors
396 % were determined through trial and error.
397
398 sf = ...
    waveletScattering('SignalLength',size(data,1), ...
    ...
    'SamplingFrequency',fs, ...
    'InvarianceScale',0.75, ...
    'QualityFactors',[4 1]);
399
400
401
402 %%
403
404 % Convert the audio signal to mono, and then call
405 % ...
    <docid:wavelet.ref#mw.edb9d918-ed3b-469f-8cb0-5e6cbad25cm7
406 % |featureMatrix|> to return the scattering ...
    coefficients for the scattering
407 % decomposition framework, |sf|.
408
409 dataMono = mean(data,2);
410 scatteringCoefficients = ...
    featureMatrix(sf,dataMono,'Transform','log');
411
412 %%
413 % Average the scattering coefficients over the ...
    10-second audio clip.
414
415 featureVector = mean(scatteringCoefficients,2);
416 fprintf('Number of wavelet features per ...
    10-second clip = %d\n',numel(featureVector))
417
418 %%
419 % The helper function ...
    |HelperWaveletFeatureVector| performs the ...
    above steps.
420 % Use a <docid:matlab.ref#bvaomuj-1 |tall|> ...
    array with
421 % <docid:matlab.ref#bsz9tpz |cellfun|> and ...
    |HelperWaveletFeatureVector| to
422 % parallelize the feature extraction. Extract ...
    wavelet feature vectors for
423 % the train and test sets.
424
425 scatteringTrain = ...
    cellfun(@ (x) HelperWaveletFeatureVector(x,sf),trainSet,tall,'UniformOutput',false);
426 xTrain = gather(scatteringTrain);
427 xTrain = cell2mat(xTrain)';
428
429 scatteringTest = ...
    cellfun(@ (x) HelperWaveletFeatureVector(x,sf),testSet,tall,'UniformOutput',false);
430 xTest = gather(scatteringTest);
431 xTest = cell2mat(xTest)';
432
433 %% Define and Train Ensemble Classifier
434 % Use |fitcensemble| to create a trained ...
    classification ensemble model
435 % (|ClassificationEnsemble|).
436
437 subspaceDimension = min(150,size(xTrain,2) - 1);
438 numLearningCycles = 30;
439 classificationEnsemble = ...
    fitcensemble(xTrain,adsTrain.Labels, ...
    'Method','Subspace', ...
    'NumLearningCycles',numLearningCycles, ...
    'Learners','discriminant', ...
    'NPredToSample',subspaceDimension, ...
    'ClassNames',removecats(unique(adsTrain.Labels)));
440
441 %% Evaluate Ensemble Classifier
442 % For each 10-second audio clip, call ...
    |predict| to return the labels and
443 % the weights, then map it to the ...
    corresponding predicted location. Call
444 % ...
    <docid:nnet.ref#mw.b571feea-6af0-489f-b52c-a4b5141ac550 ...
    |confusionchart|>
445
446 % to visualize the accuracy on the test set. ...
    Print the average.
447
448 [waveletPredictedLabels,waveletResponses] = ...
    predict(classificationEnsemble,xTest);
449
450 figure('Units','normalized','Position',[0.2 ...
    0.2 0.5 0.5])
451 cm = ...
    confusionchart(adsTest.Labels,waveletPredictedLabels,'title',
    'Accuracy - Wavelet Scattering');
452 cm.ColumnSummary = 'column-normalized';
453 cm.RowSummary = 'row-normalized';
454
455 fprintf('Average accuracy of classifier = ...
    %0.2f\n',mean(adsTest.Labels==waveletPredictedLabels)*100)
456
457
458 %% Apply Late Fusion
459 % For each 10-second clip, calling predict on ...
    the wavelet classifier and
460 % the CNN returns a vector indicating the ...
    relative confidence in their
461 % decision. Multiply the |waveletResponses| ...
    with the |cnnResponses| to
462 % create a late fusion system.
463
464 fused = waveletResponses .* cnnResponses;
465 [~,classIdx] = max(fused,[],2);
466
467 predictedLabels = classes(classIdx);
468
469 %% Evaluate Late Fusion
470 % Call |confusionchart| to visualize the fused ...
    classification accuracy.
471 % Print the average accuracy to the Command ...
    Window.
472
473 figure('Units','normalized','Position',[0.2 ...
    0.2 0.5 0.5])
474 cm = ...
    confusionchart(adsTest.Labels,predictedLabels,'title','Test
    Accuracy - Fusion');
475 cm.ColumnSummary = 'column-normalized';
476 cm.RowSummary = 'row-normalized';
477
478 fprintf('Average accuracy of fused models = ...
    %0.2f\n',mean(adsTest.Labels==predictedLabels)*100)
479
480
481 function X = ...
    HelperSegmentedMelSpectrograms(x,fs,varargin)
482 % Copyright 2019 The MathWorks, Inc.
483
484 p = inputParser;

```

```

487 addParameter(p, 'WindowLength', 1024);
488 addParameter(p, 'HopLength', 512);
489 addParameter(p, 'NumBands', 128);
490 addParameter(p, 'SegmentLength', 1);
491 addParameter(p, 'SegmentOverlap', 0);
492 addParameter(p, 'FFTLength', 1024);
493 parse(p, varargin{:})
494 params = p.Results;
495
496 x = [sum(x,2), x(:,1)-x(:,2)];
497 x = x./max(max(x));
498
499 [xb_m,~] = ...
    buffer(x(:,1), round(params.SegmentLength*fs), round(params.SegmentOverlap*fs), 'nodelay');
500 [xb_s,~] = ...
    buffer(x(:,2), round(params.SegmentLength*fs), round(params.SegmentOverlap*fs), 'nodelay');
501 xb = ...
    zeros(size(xb_m,1), size(xb_m,2)+size(xb_s,2));
502 xb(:,1:2:end) = xb_m;
503 xb(:,2:2:end) = xb_s;
504
505 spec = melSpectrogram(xb, fs, ...
506     'Window', hamming(params.WindowLength, 'periodic'), ...
507     ...
508     'OverlapLength', params.WindowLength - ...
509     params.HopLength, ...
510     'FFTLength', params.FFTLength, ...
511     'NumBands', params.NumBands, ...
512     'FrequencyRange', [0, floor(fs/2)]);
513 spec = log10(spec+eps);
514
515 X = ...
516     reshape(spec, size(spec,1), size(spec,2), size(x,2), []);
517 end
518 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
519 %HelperWaveletFeatureVector
520 function features = ...
521     HelperWaveletFeatureVector(x, sf)
522 x = mean(x,2);
523 features = featureMatrix(sf, x, 'Transform', 'log');
524 features = mean(features,2);
525 end
526 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
527 %%
528 % Copyright 2019 The MathWorks, Inc.

```