Training General-Purpose Audio Tagging Networks with Noisy Labels and Iterative Self-Verification

DCASE-2018 Challange Task 2



Matthias Dorfer and Gerhard Widmer Institute of Computational Perception





Audio Signal Pre-Processing

- Normalize to a dB-level of -0.1
- Clip silence in beginning and end of signal
- re-sampled to 32 kHz





Audio Signal Pre-Processing

- Normalize to a dB-level of -0.1
- Clip silence in beginning and end of signal
- re-sampled to 32 kHz





Spectrogram Parameters

Two spectrogram types to capture different aspects of the audio

Version – 1

- STFT hop-size: 192
- 1024-sample hann windows
- Perceptual weighting
- Mel-scaled filterbank (128 bins)



Version – 2

- STFT hop-size: 128
- 1024-sample hann windows
- Logarithm of the power spectrogram
- Log-scaled filterbank (128 bins)

| Contraction and an and and | | | |
|--|---|--|--|
| A CONTRACTOR OF A CONTRACTOR O | | | - HARD BEER STOLEN. |
| | | | and the second se |
| | | | and the second |
| the second se | | | 1 Contraction of the local division of the l |
| Contraction of the second seco | | | CONTRACTOR OF THE OWNER. |
| the second se | | | and the second s |
| | | | - CONTRACTOR OF CONTRACTOR |
| The second | | | Contraction of the second s |
| The second se | | | ALC: NOT AND DESCRIPTION OF |
| And the second sec | | | and the second se |
| | | | and set in a little |
| The second | | | |
| Contraction of the second s | | | |
| | | | the second se |
| Carlo Car | our loss of the loss of the loss of the loss of the | - 18C | And the second of the local division of the second |
| 18-19, menerosconosconosco - 191 - en esperante | P | Section of the sectio | |
| · · · · · · · · · · · · · · · · · · · | | | A REAL PROPERTY AND ADDRESS OF THE OWNER OWNER OF THE OWNER OWNE |
| at the faith of the second | | | A CONTRACT OF A CONTRACT. |
| | and the second se | | |
| the second | The Physical State State of the State of the | | The second s |
| the second | | | |



Spectrogram Length Distribution



JOHANNES KEPLER UNIVERSITY LINZ

Spectrogram Length Distribution



UNIVERSITY LINZ

Dealing with Spectrogram Lengths

Fix length to 3000 frames

- Repeat a given excerpt in case it is too short
- Clip at 3000 frames in case it is too long





Dealing with Spectrogram Lengths

Fix length to 3000 frames

- Repeat a given excerpt in case it is too short
- Clip at 3000 frames in case it is too long





Network Architecture

- Fully Convolutional Neural Network
 - VGG-Style (3 x 3 convolutions & 2 x 2 max-pooling)
 - Global Average Pooling over 41 feature maps
- Why?
 - Less parameters in classification layer
 - Deals with varying spectrogram length (Nice to have for application time)



Training Procedure

- ADAM: 500 epochs with initial lerning rate 0.001
- Linear learning rate decay starting from epoch 100





Training Procedure

- ADAM: 500 epochs with initial lerning rate 0.001
- Linear learning rate decay starting from epoch 100
- Spectrogram Excerpt Sub-Sampling (384 frame excerpts)





Training Procedure

- ADAM: 500 epochs with initial lerning rate 0.001
- Linear learning rate decay starting from epoch 100
- Spectrogram Excerpt Sub-Sampling (384 frame excerpts)
- Mixup Data Augmentation (α=0.3)





4-Fold Iterative Self-Verification

- Address the noisy labels in the development dataset.
- Central Idea: Gradually shift unverified labels into the verified, trusted training set for fine-tuning the models



Is this really class *Knock*?



4-Fold Cross-Validation Setup

- Crucial component for self-verification
- Parts of the data (to be verified) must not be presented to the verification network for training
- Prediction would be worthless
- Stratified sub-folds! (keep label distribution)





4-Fold Cross-Validation Setup

- Crucial component for self-verification
- Parts of the data (to be verified) must not be presented to the verification network for training

"The test set is composed of ~1.6k samples with manually-verified annotations and with a similar category distribution than that of the train set."









Matthias Dorfer



Matthias Dorfer

JOHANNES KEPLE

17



Verification Conditions

- 1) Automatic annotation and avg. prediction agree $(y = y_p)$
- 2) Average target class posteriors exceeds 0.95
- 3) Count of 40 self-verified examples per class is not reached

| Manually Verified | Automatically Verified | Unverified |
|----------------------|---------------------------|------------|
|----------------------|---------------------------|------------|







Experimental Setup

- We evaluate the model of the iteration with highest verified validation set score.
- Test set comprises 1600 unseen audio clips
- Evaluation Measures
 - Mean Average Precison (MAP@3)
 - F-Score for Individual Classes



Experimental Results





Experimental Results





Experimental Results

Private Kaggle Leaderboard



| | MAP@3 |
|---------|--------|
| Public | 0.9563 |
| Private | 0.9518 |



Matthias Dorfer

Live Machine Listening Demo on unseen sounds ...

Live Demo

Summary and Conclusions

Proposed Approach:

- Iterative Self-verification Loop
- Fully Convolutional Neural Network (VGG, Global Average Pooling, 2nd place Task 1A)
- Improvement from 93.87% to 96.01% (nice but we can't expect miracles)
- Reminder for how important **the right ML setup** is
- Audio (Signal) Pre-Processing is still key
- <u>https://cpjku.github.io/dcase_task2/</u>



