

Semi-supervised Audio Classification with Consistency-Based Regularization

Kangkang Lu, Chuan-Sheng Foo, Kah Kuan Teh, Tran Huy Dat, Vijay Ramaseshan Chandrasekhar

Institute for Infocomm Research, A*STAR, Singapore

luk@i2r.a-star.edu.sg, foo_chuan_sheng@i2r.a-star.edu.sg, teh_kah_kuan@i2r.a-star.edu.sg,
hdtran@i2r.a-star.edu.sg, vijay@i2r.a-star.edu.sg

Abstract

Consistency-based semi-supervised learning methods such as the Mean Teacher method are state-of-the-art on image datasets, but have yet to be applied to audio data. Such methods encourage model predictions to be consistent on perturbed input data. In this paper, we incorporate audio-specific perturbations into the Mean Teacher algorithm and demonstrate the effectiveness of the resulting method on audio classification tasks. Specifically, we perturb audio inputs by mixing in other environmental audio clips, and leverage other training examples as sources of noise. Experiments on the Google Speech Command Dataset and an audio classification task show that the method can achieve comparable performance to a purely supervised approach while using only a fraction of the labels.

Index Terms: audio classification, semi-supervised learning, data interpolation, data augmentation

1. Introduction

State-of-the-art deep learning methods typically require large amounts of labeled data to obtain high predictive performance. For audio classification in particular [1, 2], datasets need to include variations caused by the uncontrollable nature of audio sources, thus incurring increased data acquisition and more importantly, data labeling costs. Therefore, building a robust audio classification engine with limited labeled data is an important and practical problem that we are going to address in this paper.

Semi-supervised approaches that utilize small amounts of labeled data together with larger amounts of unlabeled data have previously been explored as a way to mitigate this data labeling burden, as unlabeled audio data is typically far and cheaper to obtain in practice. This approach has been explored for audio event classification [3, 4] and in speech recognition [5]. These works typically adopt a self-training approach – first training a classifier on the small amount of labeled data, then using it to predict labels for the unlabeled data. Confidently predicted labels are then included as part of the training data and the classifier is re-trained on this larger dataset. The process is repeated for several iterations. An alternative approach applied to music instrument recognition first trains a Gaussian mixture model with the labeled data, then continues to improve the model using an iterative EM-algorithm on unlabeled samples [6].

Motivated by the recent successes of consistency-based semi-supervised learning methods in computer vision [7, 8, 9], we investigate their applicability to audio data. We first briefly review this class of methods that encourage model predictions to be consistent on unlabeled examples. Specifically, for unlabeled samples, the model will be trained to have predictions that are consistent with those on perturbed inputs and different model parameters. Besides the classification loss for the labeled samples, a consistency cost will be computed between a student model and a better teacher model on the unlabeled data.

In [7], this consistency loss is between a noisy student model and a clean teacher model. Random noise was added to the latent features in the student model. During optimization, the teacher model tried to denoise the corrupted hidden features and minimize the difference between the features in the student and teacher models. Laine and Aila [8] proposed the π model where the student and teacher models share the same weights but use different data augmentation techniques and dropout. They also introduced temporal ensembling to use the exponential moving average (EMA) of the student model predictions as the outputs of teacher model. Instead of a prediction ensemble, Tarvainen and Valpola [9] used the EMA weights of the student model as the weights of teacher model and achieved state-of-the-art performance on image datasets. However, despite their success on image datasets, to the best of our knowledge these methods have not been applied to audio data.

The specific perturbations applied to the inputs during training play an important role in consistency methods; in some sense they enforce smoothness of the classifier along the data manifold. However, perturbations commonly used for images may not be suitable for audio data. Moreover, some perturbations also are specific to audio data, for instance, mixing with environmental noise. There is a special category of audio in the Google Speech Commands Dataset which is background noise (e.g., “doing the dishes”, “exercise bike”, “running tap”). The background noise is collected to help deal with noisy environments and obtain more robust classifiers. Classifiers are expected to identify the command in spite of the noise.

Besides adding environment noise directly to the training samples, we can also mix different samples together as a form of perturbation. Zhang et al. [10] proposed to train a neural network with convex combination of pairs of samples. This mixup method has been shown to be an effective and simple way of data augmentation and can improve the generalization of networks. Verma et al., [11] further extended mixup to semi-supervised learning. The consistency loss is computed between predictions of interpolated samples from the student model and the interpolated predictions from teacher model. Instead of using small perturbations like noise or data augmentation, the interpolated samples mostly lie in the low density along the decision boundaries. Adding small perturbations to these samples may easily push them to the other side of the decision boundary. By training on these samples, they achieved state-of-the-art SSL performance.

In this work, we apply Mean Teacher to two different audio datasets – the Google Speech Commands Dataset [12] and UrbanSound8K Dataset [13], and show that the consistency-based SSL methods also perform well on audio data. We also show that using environment noise as perturbation further improves performance, using two different sources of such noise. We further explore the effect of mixing samples: through optimizing the consistency of interpolated predictions from the student and

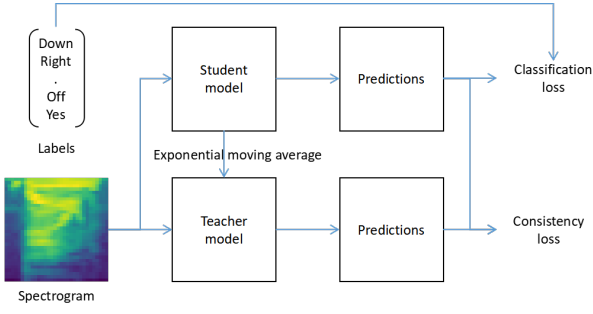


Figure 1: Architecture of Mean Teacher model.

teacher models, the accuracy is further improved. Our work establishes the consistency-based methods as a strong baseline for semi-supervised audio classification.

2. Methods

In this chapter, we will introduce how the Mean Teacher, environment noise and sample mixing are applied to the audio dataset.

2.1. Mean Teacher method

The Mean Teacher method attempts to enforce classifier smoothness in the face of small perturbations on the data and model parameters. In the method, we have a student and teacher model, both assumed to be neural networks that share the same architecture (f). The weights of teacher model (W) is the EMA of the student model's weights (w). Formally, the teacher weights at step t are computed as

$$W_t = \alpha * W_{t-1} + (1 - \alpha) * w_t, \quad (1)$$

where α is the smoothing constant that effectively controls the averaging time period.

Given a labeled set of samples $D_l = \{x, y\}$ and an unlabeled set $D_u = \{X\}$, the classification loss function would be the usual cross-entropy (CE) loss $L_{cls} = CE(f(x, w), y)$. As the student and teacher model have the same model, the consistency loss will be:

$$L_{con} = \|f(x, w) - f(x, W)\|^2 + \|f(X, w) - f(X, W)\|^2. \quad (2)$$

The total loss is the sum of classification and consistency losses

$$L_{total} = L_{cls} + \lambda L_{con}, \quad (3)$$

where λ is a parameter controlling the importance of consistency loss. When applying Mean Teacher, the consistency loss is computed on both labeled and unlabeled data.

A schematic of the method is shown in Figure 1. During training, each batch will contain both labeled and unlabeled samples. We will feed these batches of data into both the student and teacher models but with different perturbations; we discuss these in detail in Section 2.2. These perturbations may include translations, addition of Gaussian noise, horizontal and vertical flips, or training perturbations like dropout. The consistency loss is computed as the mean-square distance between the predictions from the student and teacher models. The classification loss is evaluated as the cross-entropy loss between the

softmax output of the student model and the target label. During optimization, the gradients are not backpropagated to the teacher model, instead, only the weights of the student model are updated.

2.2. Perturbations on audio data

2.2.1. Time and frequency shifts

Random time and frequency shifts are natural perturbations on audio data. Translated to spectrogram images, these perturbations correspond to horizontal translations for time shifts and vertical translations for frequency shifts.

2.2.2. Gaussian noise

Additive Gaussian noise is also used as a possible perturbation as there will be noise when recording and transforming the audio data.

2.2.3. Environment noise as perturbations

In the Google Speech Commands dataset, we have six classes of background noise, one sample per class. Owing to the linearity of the spectrogram, we can simply add the spectrogram of the audio sample with that of the background noise to mimic the situation when the command audio is recorded with the noise. Given the original spectrogram as S_{ori} and the environment noise spectrogram as S_{env} , the final sample (S_{in}) fed into the network would be:

$$S_{in} = (1 - \beta) * S_{ori} + \beta * S_{env}, \quad (4)$$

where β is used to control the environment noise level. Compared to adding random Gaussian noise, real environment noise is a more realistic form of noise, and can potentially enable classifiers to be robust to such noise in the input. As there are only limited number of noise samples in the speech commands dataset, we also use the unrelated UrbanSounds8K urban audio dataset to enrich the quantity and diversity of environment audio data.

2.2.4. Consistency regularization with mixed samples

Another possible source of perturbations are the original samples themselves. Mixed samples can be regarded as the case when different commands are given at the same time. Given two samples as (x_a, x_b) and their corresponding labels (y_a, y_b) , the generated pair of samples would be:

$$x_{com} = (1 - \gamma) * x_a + \gamma * x_b \quad (5)$$

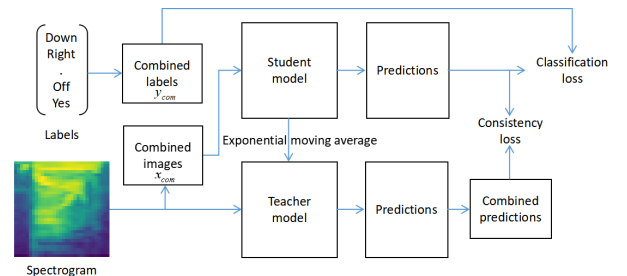


Figure 2: Architecture of Mean Teacher model with mixed samples as input.

$$y_{com} = (1 - \gamma) * y_a + \gamma * y_b \quad (6)$$

The classification loss will be computed as:

$$L_{cls_com} = CE(f(x_{com}, w), y_{com}) \quad (7)$$

For mixed unlabeled data (X_a, X_b), their predictions from the student and teacher models would be:

$$P_{stu} = f(X_{com}, w) \quad (8)$$

$$P_{tch} = (1 - \gamma) * f(X_a, W) + \gamma * f(X_b, W), \quad (9)$$

where $X_{com} = (1 - \gamma) * X_a + \gamma * X_b$. And the consistency loss would be:

$$L_{con_com} = ||P_{stu} - P(tch)||^2 \quad (10)$$

In the case where we utilized mixed samples as perturbations, the training workflow will be different as shown in Figure 2. The samples will be mixed before being fed into the network. For mixed samples with corresponding combined labels, they will contribute to the classification loss. The predictions of original samples from teacher model will be combined as P_{tch} . The consistency loss is the mean-square distance between combined predictions from teacher model and predictions of mixed samples from student model as noted in Equation 10.

3. Experiments

3.1. Datasets

We performed experiments on two datasets: the Google Speech Commands Dataset [12] and the UrbanSound8K Dataset [13]. In both cases, during training, the entire training set is considered as the unlabeled dataset, while labeled samples are randomly selected from the training set in each run.

Google Speech Commands Dataset: This dataset includes one-second long utterances of 30 different speech commands such as “Yes”, “No”, “Up” and “Down”, spoken by many different people. We split the dataset into a 57886-sample training set and a 6835-sample testing set according to the official partition script. We convert the waveform audio into 32×32 mel-spectrogram images for training.

UrbanSound8K Dataset: This is a collection of short audio clips (1 to 4 seconds long) from 10 different classes collected from urban acoustic environments. Classes include sounds of car horns, dogs barking, and jackhammers. There are 8732 audio samples in total; we used the training set including 5434 labeled samples in the version hosted on Kaggle. We split the labeled data into training and test sets. We transform the audio into spectrogram images with the same sampling rate and then zero-pad them to 32×128 for training.

3.2. Experimental setup

For our classifiers (both the supervised baseline and in Mean Teacher), we used the same 13-layer convolutional neural network proposed in [9]. Briefly, the network consists of 3 blocks of 3 convolution layers followed by a single pooling layer; the last block is then connected to a fully-connected layer that computes the prediction logits. The network also includes a dropout layer after each of the first two pooling layers. Batch normalization and weight normalization are applied to the convolution and fully-connected layers. The model is trained using the ADAM optimizer with a batch size of 100 for 80000 iterations.

During the first few epochs of training, the optimization relies more on the labeled data and classification loss. As the

model starts to converge, the consistency loss plays a more important role. As described in [9], a ramp-up on λ in Equation 3 is applied to adapt the importance of the consistency loss as training progresses. Specifically, both the learning rate and λ will be ramped up to their maximum value in the first 40000 steps. The learning rate will also ramp down to 0 in the last 25000 steps for better convergence. For details on the ramp-up and ramp-down procedure, please refer to the experimental setup section in the appendix of [9]. The maximum learning rate is set to 0.003. All experiments are repeated three times with different random seeds, and we report the average classification accuracy along with standard deviations.

3.3. Evaluation of the Mean Teacher method

We first compared the performance of the semi-supervised Mean Teacher method to a purely supervised convolutional neural network with the same architecture (see Section 3.2) that only uses the labeled data during training. For these experiments we only included random time and frequency shifts as well as Gaussian noise as perturbations. We report classification accuracy on the Google Speech Commands Dataset (Table 1) and UrbanSound8K Dataset (Table 2).

On the Google Speech Commands Dataset, Mean Teacher outperforms the supervised baseline when not all labeled data is used. The performance gain increases as less labeled data is used – from 5% when only 1% of labels are used to 0.5% when 25% of labels are used in training. Mean Teacher achieves comparable performance to the supervised baseline using all labels, with only 25% (even 10%) of the labels.

Table 1: *Classification Accuracy (%) of Mean Teacher on the Google Speech Commands Dataset (57886 training samples)*

Labels	Supervised	Mean Teacher
600	86.87 \pm 0.50	92.18 \pm 0.25
3000	93.31 \pm 0.26	95.10 \pm 0.25
6000	94.54 \pm 0.02	95.70 \pm 0.15
15000	95.64 \pm 0.14	96.18 \pm 0.04
57886	96.64 \pm 0.07	96.61 \pm 0.04

The UrbanSound8K dataset is more challenging as it has fewer training samples. We also observe that Mean Teacher achieves performance gains over the supervised baseline across the board. Here, gains are significant even when almost a third of the samples are labeled (3.5% gain when 1500 labels are used). With only 12% of the labels (600 labels), Mean Teacher achieves a large improvement in accuracy of 10% over the supervised baseline. On this dataset it is not possible to achieve comparable accuracy to a supervised baseline using all labels with less labeled data. However, Mean Teacher manages to obtain 90% of the accuracy using only 30% of the labels.

3.4. Incorporating environment noise as perturbations

3.4.1. Background noise from Google Speech Dataset

To add the environment noise, before each iteration, we randomly pick one category of environment noise, crop it to a 32×32 patch and add it to the training batch. We do a grid search for best β in Equation 4 in range of $[0.1, 0.5]$ with step size of 0.1. A randomly selected β in range of $[0.1, 0.5]$ is also evaluated. After validating on 600 labeled samples on the commands dataset, the best performance is achieved with $\beta = 0.1$. From

Table 2: *Classification Accuracy (%) of Mean Teacher on UrbanSound8K (4892 training samples)*

Labels	Supervised	Mean Teacher
60	33.89 ± 0.40	38.31 ± 3.32
300	60.53 ± 1.78	67.65 ± 1.58
600	64.27 ± 5.08	75.20 ± 1.94
1500	81.95 ± 1.97	85.64 ± 1.17
4892	93.43 ± 0.38	93.62 ± 0.09

Table 3, we can see that, with this new perturbation, the accuracy is slightly improved. This lack of significant improvement could be due to the lack of diversity in the background noise samples (6 samples, 40 seconds each).

3.4.2. Urban noise from UrbanSound8K

To see if a more diverse set of environment noise could provide performance benefits, we also considered the Urban Sound Dataset [13] as a source of noise, as it spans more classes and includes more samples. We used the same noise setting as found on the Google Speech Commands Dataset in our experiments.

From the results shown in Table 3, we can see the urban noise brings an impressive increment on performance. Compared with only using the six noise samples in the Google Speech Commands Dataset, urban noise is much more effective. With 600 labels, while the command noise increases accuracy by 0.13%, incorporating urban noise provides a larger boost of 0.98%. With only 15000 (around 25%) labels, performance surpasses that of the supervised baseline using all the labels. We conclude that the quantity and variety of these noise has a great influence on classification performance.

However, these results also raise the question: is the improvement due to the new noise perturbation or simply due to a data augmentation effect. To clarify this, we apply environmental noise to supervised training as a way of data augmentation. The results are shown in Table 4. We can see that, except the improvement in the case where 600 labels are used, there is no clear improvement when using urban noise as part of a data augmentation approach. These results indicate the effectiveness of introducing environmental noise as a new form of perturbation as part of the Mean Teacher method.

3.5. Mixing samples as perturbations

Finally, we evaluate the effect of mixing samples as perturbations. We do not include environment noise in these experiments to isolate the effects of the mixing. Results are shown in Table 5. With mixed samples, instead of introducing more perturbations, we are augmenting the dataset with more training

Table 3: *Classification Accuracy (%) with environmental noise on Google Speech Commands Dataset*

Labels	Mean Teacher	MT+noise	MT+urban noise
600	92.18 ± 0.25	92.31 ± 0.04	93.29 ± 0.59
3000	95.10 ± 0.25	95.07 ± 0.14	95.94 ± 0.36
6000	95.70 ± 0.15	95.82 ± 0.14	96.37 ± 0.13
15000	96.18 ± 0.04	96.22 ± 0.03	96.83 ± 0.22
57886	96.61 ± 0.04	96.63 ± 0.01	97.33 ± 0.18

Table 4: *Classification Accuracy (%) of using environmental noise as data augmentation on Google Speech Commands Dataset*

Labels	Supervised	Urban noise augmentation
600	86.87 ± 0.50	87.19 ± 0.50
3000	93.31 ± 0.26	93.31 ± 0.27
6000	94.54 ± 0.02	94.59 ± 0.20
15000	95.64 ± 0.14	95.73 ± 0.14
57886	96.64 ± 0.07	96.61 ± 0.08

Table 5: *Classification Accuracy (%) incorporating sample mixing on Google Speech Commands Dataset*

Labels	Mean Teacher	MT+sample mixing
600	92.18 ± 0.25	94.30 ± 0.31
3000	95.10 ± 0.25	95.87 ± 0.19
6000	95.70 ± 0.15	96.26 ± 0.14
15000	96.18 ± 0.04	96.61 ± 0.07
57886	96.61 ± 0.04	97.17 ± 0.04

samples near the decision boundary of different target labels. The performance improvements are small (0.5-1%), similar to when we use urban noise as perturbations. We note that the greatest improvement was achieved when only 1% of the labels were used (600 labels), where it provides a 2% boost in accuracy over the vanilla Mean Teacher method.

4. Conclusion

In this paper, we demonstrated the effectiveness of the Mean Teacher method for audio classification on two very different audio classification datasets – speech and urban audio. We introduce the addition of environment noise as a new effective perturbation as part of the Mean Teacher method, showing improvements using background noise from the Google Speech Commands Dataset and more significant gains using noise from the UrbanSound8K dataset. Our results show the importance of using diverse collections of noise as perturbations and that the improvement cannot simply be achieved using a similar data augmentation strategy. Finally, we use mixed samples and labels as a source of perturbations. In all cases, we show that Mean Teacher can achieve comparable results to fully supervised training using only 25% (or even 10%) of the labels. We also observe significant accuracy boosts over the supervised baseline when only 1% of the labels are used. These promising results indicate the feasibility of applying the Mean Teacher method, environment noise and sample mixing on audio classification tasks. We believe that these methods can be adapted and improved for other audio datasets and tasks, and leave this as an interesting direction for future work.

5. References

- [1] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, “Audio Set: An ontology and human-labeled dataset for audio events,” in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2017.
- [2] A. Jansen, J. F. Gemmeke, D. P. Ellis, X. Liu, W. Lawrence, and D. Freedman, “Large-scale audio event discovery in one million

YouTube videos,” in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2017.

- [3] Z. Zhang and B. Schuller, “Semi-supervised learning helps in sound event classification,” in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, 2012.
- [4] W. Han, E. Coutinho, H. Ruan, H. Li, B. Schuller, X. Yu, and X. Zhu, “Semi-supervised active learning for sound classification in hybrid learning environments,” *PLoS ONE*, 2016.
- [5] D. Hakkani-Tur, G. Tur, M. Rahim, and G. Riccardi, “Unsupervised and active learning in automatic speech recognition for call classification,” 2004.
- [6] A. Diment, T. Heittola, and T. Virtanen, “Semi-supervised learning for musical instrument recognition,” in *21st European Signal Processing Conference (EUSIPCO 2013)*, Sep. 2013, pp. 1–5.
- [7] A. Rasmus, M. Berglund, M. Honkala, H. Valpola, and T. Raiko, “Semi-supervised learning with ladder networks,” in *NIPS*, 2015.
- [8] S. Laine and T. Aila, “Temporal ensembling for semi-supervised learning,” *International Conference on Learning Representations (ICLR)*, 2017.
- [9] A. Tarvainen and H. Valpola, “Weight-averaged consistency targets improve semi-supervised deep learning results,” in *NIPS*, 2017.
- [10] Y. N. D. D. L.-P. Hongyi Zhang, Moustapha Cisse, “mixup: Beyond empirical risk minimization,” *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=r1Ddp1-Rb>
- [11] V. Verma, A. Lamb, J. Kannala, Y. Bengio, and D. Lopez-Paz, “Interpolation Consistency Training for Semi-Supervised Learning,” *arXiv e-prints*, p. arXiv:1903.03825, Mar 2019.
- [12] P. Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” 2018.
- [13] J. Salamon, C. Jacoby, and J. P. Bello, “A dataset and taxonomy for urban sound research,” in *22nd ACM International Conference on Multimedia (ACM-MM’14)*, Orlando, FL, USA, Nov. 2014, pp. 1041–1044.