

ZERO-SHOT AUDIO CLASSIFICATION WITH FACTORED LINEAR AND NONLINEAR ACOUSTIC-SEMANTIC PROJECTIONS

Huang Xie*, Okko Räsänen*[†], Tuomas Virtanen*

* Unit of Computing Sciences, Tampere University, Finland

[†] Dept. Signal Processing and Acoustics, Aalto University, Finland

ABSTRACT

In this paper, we study zero-shot learning in audio classification through factored linear and nonlinear acoustic-semantic projections between audio instances and sound classes. Zero-shot learning in audio classification refers to classification problems that aim at recognizing audio instances of sound classes, which have no available training data but only semantic side information. In this paper, we address zero-shot learning by employing factored linear and nonlinear acoustic-semantic projections. We develop factored linear projections by applying rank decomposition to a bilinear model, and use nonlinear activation functions, such as tanh, to model the non-linearity between acoustic embeddings and semantic embeddings. Compared with the prior bilinear model, experimental results show that the proposed projection methods are effective for improving classification performance of zero-shot learning in audio classification.

Index Terms— audio classification, zero-shot learning, acoustic-semantic projection

1. INTRODUCTION

Supervised learning has been well-studied for tackling audio classification problems, such as acoustic scene classification [1] and environmental sound classification [2, 3]. To obtain classifiers with satisfactory performance, existing supervised learning techniques require large amounts of annotated training data from target sound classes, which is labor-intensive and costly to acquire. Moreover, with the increasing diversity of observed sound classes, it becomes even more challenging for humans to collect sufficient annotated training data for all possible sound classes. Recent work [4, 5, 6, 7] in the audio recognition literature that deal with the lack of adequate training data mainly apply *data augmentation* [8], *meta learning* [9] and *few-shot learning* [10] methods. However, a certain amount of representative training data from target classes is still indispensable to make these methods work. Furthermore, to classify instances from novel classes, a supervised learning classifier would require retraining for the novel classes, which can be time-consuming and requires exhaustive parameter tuning.

In this paper, we consider the extreme case of audio classification where target sound classes have no available training samples but only class side information (e.g., textual descriptions). This problem is generally referred to as the *zero-shot learning* [11], which has been increasingly studied in the context of image classification. In contrast to conventional supervised learning, zero-shot learning uses training data from only predefined classes (i.e., seen classes) to obtain classifiers that can be generalized to novel classes (i.e., unseen classes).

There is only limited work that has been done for zero-shot learning in audio classification. Due to the lack of training data from unseen classes, class side information is used as a compensation for exploring the relationship between seen classes and unseen classes to make zero-shot learning possible. Prior work [12, 13, 14] tackled zero-shot learning by leveraging semantic side information of sound classes, such as textual labels, with a two-phase learning process. First, intermediate-level representations were learned for audio instances and sound classes, respectively. Audio instances were embedded into a low-dimensional acoustic space with feature learning techniques, such as *VGGish* [15]. Sound classes were represented by word embeddings in a semantic space, which were extracted from their semantic side information with pre-trained language models, such as *Word2Vec* [16]. Then, an acoustic-semantic projection was learned to associate acoustic embeddings with semantic embeddings. Islam et al. [12] employed a two-layer fully-connected neural network to model a nonlinear projection of acoustic embeddings onto semantic embeddings. In our previous work [13, 14], a bilinear model was used to learn a bidirectional linear projection between acoustic embeddings and semantic embeddings. Another prior work [17] was conducted by integrating the acoustic embedding learning phase with the acoustic-semantic projection learning phase to optimize them holistically. Thus, a nonlinear acoustic-semantic projection was inherently built into their model.

In this paper, we extend our previous work [13, 14] by introducing matrix decomposition and nonlinear activation functions (e.g., tanh) to the bilinear model. We develop factored linear and nonlinear acoustic-semantic projections for zero-shot learning in audio classification. Experimental results show that the proposed projection methods are effective

tive for improving classification performance of zero-shot learning in audio classification.

The remainder of this paper is organized as follows. In Section 2, we introduce the concept of zero-shot learning in audio classification. Then, we present the proposed factored linear and nonlinear projections in Section 3. We describe the training algorithm in Section 4, and discuss the experimental results in Section 5. Finally, we conclude this paper in Section 6.

2. ZERO-SHOT LEARNING

In this section, we introduce the concept of zero-shot learning via semantic side information in audio classification. We denote the audio sample space by X , the set of seen sound classes by Y , and the set of unseen sound classes by Z . Note that Y and Z are disjoint. Let $\theta(x) \in \mathbb{R}^{d_a}$ be the acoustic embedding of an audio instance $x \in X$ in an acoustic space, and $\phi(y) \in \mathbb{R}^{d_s}$, $\phi(z) \in \mathbb{R}^{d_s}$ be the semantic embeddings of sound classes $y \in Y$ and $z \in Z$ in a semantic space, respectively. We are given the training data $S_{tr} = \{(x_n, y_n) \in X \times Y | n = 1, \dots, N\}$, where x_n is an annotated audio sample belonging to a seen sound class y_n . In zero-shot learning, our goal is to learn an audio classifier $f : X \rightarrow Z$, which can predict the correct sound class for an audio instance $x \in X$, given its acoustic embedding $\theta(x)$ and the semantic embeddings $\phi(z)$ of every candidate sound class $z \in Z$.

In prior work [12, 13, 14, 17], this is done via learning an acoustic-semantic projection $T : \mathbb{R}^{d_a} \rightarrow \mathbb{R}^{d_s}$ of acoustic embeddings onto semantic embeddings. Given an audio instance $x \in X$ belonging to a sound class $z_x \in Z$, it is generally assumed that the projected acoustic embedding $T(\theta(x))$ in the semantic space should be closer to the semantic embedding $\phi(z_x)$ of its correct sound class z_x rather than those of other sound classes. A similarity scoring function $F : \mathbb{R}^{d_s} \times \mathbb{R}^{d_s} \rightarrow \mathbb{R}$, as known as the *compatibility function*, is then defined to measure how similar/compatible a projected acoustic embedding and a semantic embedding are. To measure the similarity of two embedding vectors, the popular choices of F could be Euclidean distance [12], cosine similarity [17], dot product [13, 14], etc. Therefore, the classifier $f : X \rightarrow Z$ is formulated as¹

$$z_x = f(x) = \operatorname{argmax}_{z \in Z} F(T(\theta(x)), \phi(z)). \quad (1)$$

During the training stage, the projection T is trained with audio samples $(x_n, y_n) \in S_{tr}$ such that

$$y_n = f(x_n) = \operatorname{argmax}_{y \in Y} F(T(\theta(x_n)), \phi(y)). \quad (2)$$

For prediction, an audio instance will be classified into the sound class, the semantic embedding of which is most compatible with its projected acoustic embedding.

¹Note that the *argmin* operation is used for other compatibility functions, such as Euclidean distance, etc.

3. FACTORED LINEAR AND NONLINEAR PROJECTIONS

In this section, we introduce our approach for developing factored linear and nonlinear acoustic-semantic projections for zero-shot learning in audio classification. First, we present the bilinear model used in our previous work [13, 14], on which we build factored linear and nonlinear acoustic-semantic projections. Then, we describe a factored linear projection obtained by applying matrix decomposition to the bilinear model. After that, we introduce the nonlinear projections derived from the factored linear projection by introducing nonlinear activation functions, such as \tanh .

3.1. Bilinear Model

Inspired by prior work [18, 19] in computer vision, we employed a bilinear model in [13, 14] to learn a bidirectional linear projection between acoustic embeddings and semantic embeddings. The audio classifier $f : X \rightarrow Z$ was then written in a bilinear form as

$$f(x) = \operatorname{argmax}_{z \in Z} \theta(x)' W \phi(z), \quad (3)$$

where W was the learned projection matrix. Considering a projection of acoustic embeddings onto semantic embeddings, T could be formulated as

$$T(\theta(x)) = W' \theta(x). \quad (4)$$

Thus, in (3), the dot product was inherently defined as the compatibility function F

$$F(T(\theta(x)), \phi(z)) = T(\theta(x))' \phi(z). \quad (5)$$

3.2. Factored Linear Projection

In the case where d_a and d_s are large, it would be valuable to decompose W into a product of two low-rank matrices $U_{d_a \times r}$ and $V_{r \times d_s}$ to reduce the effective number of parameters in the bilinear model. Thus, we consider the rank decomposition $W = UV$ in (4) to build a factored linear projection

$$T(\theta(x)) = (UV)' \theta(x) = V' U' \theta(x). \quad (6)$$

3.3. Nonlinear Projection

Based on the linear projection (6), we consider introducing nonlinear activation functions (e.g., \tanh) into it to model the potential nonlinear relationship between acoustic embeddings and semantic embeddings. We apply a nonlinear activation function t to $U' \theta(x)$. Therefore, a nonlinear projection is formulated as

$$T(\theta(x)) = V' t(U' \theta(x)). \quad (7)$$

We notice the fact that (7) also defines a two-layer fully-connected neural network² with input $\theta(x)$, layer weights U'

²Note that the bias parameters of each layer are ignored from (7) and later formulas in this paper.

Class Fold	Sound Class	Audio Sample
Fold0	104	23007
Fold1	104	22889
Fold2	104	22762
Fold3	104	22739
Fold4	105	21377

Table 1. Class folds in the selected subset of AudioSet.

and V' . To describe deeper fully-connected neural networks, we can simply add more activation functions and matrices between U' and V' . For example, a three-layer fully-connected neural network is formulated as

$$T(\theta(x)) = V' t(Q t(U' \theta(x))). \quad (8)$$

where the matrix Q denotes the weight parameters of the second fully-connected layer.

4. TRAINING ALGORITHM

In this section, we introduce the algorithm for learning an acoustic-semantic projection T with training data S_{tr} . Given an audio sample $(x_n, y_n) \in S_{tr}$, we consider the task of sorting sound classes $y \in Y$ in descending order according to their compatibility values $F(T(\theta(x_n)), \phi(y))$. Our objective is to optimize T so that the correct class y_n would be ranked at top of the sorted class list, i.e., having the maximal compatibility value for x_n .

Let r_y be the position index of a sound class y in the sorted class list. We define $r_y = 0$ when y is sorted at the first position. By applying a ranking error function [20], we transform position index r into loss $\beta(r)$:

$$\beta(r) = \sum_{i=1}^r \alpha_i, \quad (9)$$

with $\alpha_1 \geq \alpha_2 \geq \dots \geq 0$ and $\beta(0) = 0$. Specifically, α_i denotes a penalty to a class losing a position from $i - 1$ to i . In this paper, we follow previous work [18, 20] and choose $\alpha_i = 1/i$.

To learn an acoustic-semantic projection T with audio samples $(x_n, y_n) \in S_{tr}$, we minimize the weighted approximate-rank pairwise objective [21]

$$\frac{1}{N} \sum_{n=1}^N \frac{\beta(r_{y_n})}{r_{y_n}} \sum_{y \in Y} \max\{0, l(x_n, y_n, y)\}, \quad (10)$$

with the convention $0/0 = 0$ when y_n is top-ranked. In this paper, we define the hinge loss $l(x_n, y_n, y)$ as

$$l(x_n, y_n, y) = \Delta(y_n, y) + F(T(\theta(x_n)), \phi(y)) - F(T(\theta(x_n)), \phi(y_n)), \quad (11)$$

where $\Delta(y_n, y) = 0$ if $y_n = y$ and 1 otherwise. The objective (10) is convex and can be optimized through stochastic gradient descent. To prevent over-fitting, we regularize (10) with L2 norms of parameter matrices in T .

5. EXPERIMENTS

In this section, we evaluate the proposed method with AudioSet [22] and report the effectiveness of factored linear and nonlinear projections.

5.1. Dataset

AudioSet [22] is a large unbalanced audio dataset, which contains over two million weakly labeled audio clips covering 527 sound classes. Most of these audio clips are multi-label. In this work, we focus on zero-shot learning in single-label classification problems. We follow the same experimental setup as in [14]. An audio subset containing 112,774 single-label 10-second audio clips and 521 sound classes is selected from AudioSet. We randomly split the selected subset into five disjoint class folds. The number of sound classes and audio samples in each class fold is shown in Table 1.

5.2. Acoustic Embeddings

A pre-trained VGGish [15] was used to generate acoustic embeddings from audio clips in [13]. In zero-shot learning, it is generally assumed that unseen classes are unknown at training stage. Since a pre-trained VGGish may have already embedded knowledge about unseen sound classes, using it to generate acoustic embedding can lead to a biased evaluation of zero-shot learning in audio classification. Therefore, we train VGGish from scratch with audio data excluding unseen sound classes in this work.

Following [13, 14], an 10-second audio clip is first split into ten one-second audio segments without overlapping. Then, a 128-dimensional embedding vector is generated for each audio segments with the trained VGGish. To obtain the clip-level acoustic embedding for an audio clip, we take the average of the 128-dimensional embedding vectors extracted from its one-second audio segments.

5.3. Class Semantic Embeddings

In AudioSet [22], a sound class is described by one or several textual labels (i.e., words and phrases) and an additional short description (i.e., sentences). In this work, we consider only textual labels as class semantic side information. We adopt Word2Vec [16] as a word embedding model for generating semantic embeddings from these textual labels. For the sake of simplicity, we use a publicly available pre-trained Word2Vec³, which embeds roughly three million English words and phrases. It outputs a 300-dimensional semantic word vector for a single word or a phrase. To represent a sound class with semantic word vectors, we calculate the average of these word vectors extracted from its textual labels.

³Word2Vec: <https://code.google.com/archive/p/word2vec>.

5.4. Experimental Setup

In the following experiments, we first train an VGGish for generating acoustic embeddings with class folds “Fold0” and “Fold1”. Then, we conduct zero-shot learning in audio classification with “Fold2” for training, “Fold3” for parameter validation and “Fold4” for test, respectively.

VGGish Training. Audio samples of class folds “Fold0” and “Fold1” are first randomly split into training/validation partitions with a class-specific proportion of 75/25. Then, we train an VGGish from scratch by feeding log mel spectrogram extracted from audio clips into it. After training, the trained VGGish achieves a classification accuracy (TOP-1) of 27.4% on the validation partition.

Zero-Shot Learning. We conduct zero-shot learning in audio classification with the proposed factored linear and nonlinear projections, respectively. We use the bilinear model as the baseline method. For the factored linear projection (6), we experiment with low-rank decomposition and full-rank decomposition of W to investigate the effect of the rank r and the L2 norm regularization. For the nonlinear projection (7), we experiment on three widely used activation functions, i.e., ReLU, sigmoid and tanh. We implement (7) by two-layer fully-connected neural networks, which are denoted by $FC2_{relu}$, $FC2_{sigmoid}$ and $FC2_{tanh}$, respectively. Similarly, we implement (8) by an three-layer fully-connected neural network with tanh, which is denoted by $FC3_{tanh}$. To prevent randomness, each projection method is evaluated twenty times with random initialization. The averages and standard deviations of their TOP-1 accuracies are reported in Table 2.

5.5. Result and Analysis

As a baseline method, the bilinear model achieves an averaged TOP-1 of 5.7% with a standard deviation of 1.1%. For factored linear projections with either low-rank decomposition or full-rank decomposition of W , we obtain similar results (roughly $6.4 \pm 0.6\%$). We conclude that, with the L2 norm regularization, the rank r has a limited influence on classification performance. Here, we report the result from a factored linear projection with the full-rank decomposition (i.e., $r=128$) of W , which has an averaged TOP-1 of 6.3% with a standard deviation of 0.8%. Unpaired t-test with $\alpha=0.05$ is used to measure the statistical significance among different methods. We find the results of the factored linear projection are significantly different from those of the bilinear model ($t(38)=2.09$, $p=0.04$). It shows that classification performance is improved by applying rank decomposition to W with L2 norm regularization. For nonlinear projections, we set $r=128$. Classification performance is impaired with $FC2_{relu}$ ($5.5 \pm 0.9\%$) while it is improved with $FC2_{sigmoid}$ ($7.0 \pm 0.5\%$) and $FC2_{tanh}$ ($7.2 \pm 0.6\%$). Particularly, the averaged TOP-1 of $FC2_{tanh}$ is significantly better than those of the bilinear model ($t(38)=5.60$, $p=4.50e-6$) and the factored linear projection ($t(38)=3.88$, $p=4.59e-4$). Compared with

Acoustic-Semantic Projection		TOP-1 (%) (avg \pm std)
Bilinear (baseline)		5.7 ± 1.1
Factored Linear		6.3 ± 0.8
Nonlinear	$FC2_{relu}$	5.5 ± 0.9
	$FC2_{sigmoid}$	7.0 ± 0.5
	$FC2_{tanh}$	7.2 ± 0.6
	$FC3_{tanh}$	6.0 ± 0.6

Table 2. Zero-Shot learning in audio classification with different acoustic-semantic projections.

sigmoid and tanh, the ReLU function introduces non-linearity by simply dropping negative values from its inputs. In a two-layer fully-connected neural network, this can lead to a poor acoustic-semantic projection and results in an impaired performance. For $FC2_{sigmoid}$ and $FC2_{tanh}$, we think both of them capture non-linearity between acoustic embeddings and semantic embeddings, which is useful for improving classification performance. However, compared with $FC2_{tanh}$, classification performance is impaired with $FC3_{tanh}$ ($6.0 \pm 0.6\%$). It seems that it would not be helpful for improving classification performance in zero-shot learning by simply introducing more nonlinear layers in a fully-connected neural network.

6. CONCLUSION

In this paper, we present an approach for zero-shot learning in audio classification with factored linear and nonlinear acoustic-semantic projections. We develop factored linear and nonlinear projections by applying rank decomposition and nonlinear activation functions to a bilinear model. We evaluate our proposed approach with a large unbalanced audio dataset. The experimental results show that both factored linear and nonlinear projections are effective for zero-shot learning in audio classification. With a factored linear projection, it achieves an averaged TOP-1 accuracy of 6.3%, which is better than the prior bilinear model (5.7%). By introducing nonlinear activation functions into it, classification performance can be further improved. Classification performance achieves an averaged TOP-1 accuracy of 7.2% by using a nonlinear projection with the tanh activation function.

7. ACKNOWLEDGEMENT

The research leading to these results has received funding from the European Research Council under the European Unions H2020 Framework Programme through ERC Grant Agreement 637422 EVERYSOUND. OR was funded by Academy of Finland grant no. 314602.

8. REFERENCES

- [1] D. Barchiesi, D. D. Giannoulis, D. Stowell, and M. D. Plumbley, “Acoustic Scene Classification: Classifying environments from the sounds they produce,” *IEEE Signal Process. Mag.*, vol. 32, no. 3, pp. 16–34, 2015.
- [2] K. J. Piczak, “Environmental sound classification with convolutional neural networks,” in *Proc. 25th IEEE Int. Workshop Mach. Learn. Signal Process. (MLSP)*, 2015, pp. 1–6.
- [3] K. J. Piczak, “ESC: Dataset for Environmental Sound Classification,” in *Proc. 23rd ACM Int. Conf. Multimedia (MM)*, 2015, pp. 1015–1018.
- [4] J. Salamon and J. P. Bello, “Deep Convolutional Neural Networks and Data Augmentation for Environmental Sound Classification,” *IEEE Signal Process. Lett.*, vol. 24, no. 3, pp. 279–283, 2017.
- [5] N. R. Koluguri, M. Kumar, S. H. Kim, C. Lord, and S. Narayanan, “Meta-Learning for Robust Child-Adult Classification from Speech,” in *Proc. IEEE Int. Conf. Acoustic., Speech and Signal Process. (ICASSP)*, 2020, pp. 8094–8098.
- [6] B. Shi, M. Sun, K. C. Puvvada, C.-C. Kao, S. Matsoukas, and C. Wang, “Few-Shot Acoustic Event Detection Via Meta Learning,” in *Proc. IEEE Int. Conf. Acoustic., Speech and Signal Process. (ICASSP)*, 2020, pp. 76–80.
- [7] S.-Y. Chou, K.-H. Cheng, J.-S. R. Jang, and Y.-H. Yang, “Learning to Match Transient Sound Events Using Attentional Similarity for Few-shot Sound Recognition,” in *Proc. IEEE Int. Conf. Acoustic., Speech and Signal Process. (ICASSP)*, 2019, pp. 26–30.
- [8] M. Brian, E. J. Humphrey, and J. P. Bello, “A software framework for musical data augmentation,” in *Proc. 16th Int. Soc. Music Inf. Retrieval Conf. (ISMIR)*, 2015, pp. 248–254.
- [9] C. Finn, P. Abbeel, and S. Levine, “Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks,” in *Proc. 34th Int. Conf. Mach. Learn. (ICML)*, 2017, vol. 70, pp. 1126–1135.
- [10] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni, “Generalizing from a Few Examples: A Survey on Few-Shot Learning,” *ACM Comput. Surv.*, vol. 53, no. 3, June 2020.
- [11] M. Palatucci, D. Pomerleau, G. Hinton, and T. M. Mitchell, “Zero-Shot Learning with Semantic Output Codes,” in *Proc. 22nd Int. Conf. Neural Inf. Process. Syst. (NIPS)*, 2009, pp. 1410–1418.
- [12] M. T. Islam and S. Nirjon, “SoundSemantics: Exploiting Semantic Knowledge in Text for Embedded Acoustic Event Classification,” in *Proc. 18th Int. Conf. Inf. Process. Sensor Networks (IPSN)*, 2019, pp. 217–228.
- [13] H. Xie and T. Virtanen, “Zero-Shot Audio Classification Based on Class Label Embeddings,” in *Proc. IEEE Workshop Appl. Signal Process. Audio and Acoustic. (WASPAA)*, 2019, pp. 264–267.
- [14] H. Xie and T. Virtanen, “Zero-Shot Audio Classification via Semantic Embeddings,” 2020, Submitted.
- [15] S. Hershey, S. Chaudhuri, D. P.W. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold, M. Slaney, R. J. Weiss, and K. Wilson, “CNN Architectures for Large-Scale Audio Classification,” in *Proc. IEEE Int. Conf. Acoustic., Speech and Signal Process. (ICASSP)*, 2017, pp. 131–135.
- [16] T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” in *Proc. 1st Int. Conf. Learn. Representations (ICLR)*, 2013.
- [17] J. Choi, J. Lee, J. Park, and J. Nam, “Zero-shot Learning for Audio-based Music Classification and Tagging,” in *Proc. 20th Int. Soc. Music Inf. Retrieval Conf. (ISMIR)*, 2019, pp. 67–74.
- [18] Z. Akata, F. Perronnin, Z. Harchaoui, and C. Schmid, “Label-Embedding for Image Classification,” *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 1425–1438, 2016.
- [19] Y. Xian, Z. Akata, G. Sharma, Q. Nguyen, M. Hein, and B. Schiele, “Latent Embeddings for Zero-Shot Classification,” in *Proc. IEEE Conf. Comput. Vision and Pattern Recognit. (CVPR)*, 2016, pp. 69–77.
- [20] N. Usunier, D. Buffoni, and P. Gallinari, “Ranking with Ordered Weighted Pairwise Classification,” in *Proc. 26th Annu. Int. Conf. Mach. Learn. (ICML)*, 2009, pp. 1057–1064.
- [21] J. Weston, S. Bengio, and N. Usunier, “WSABIE: Scaling Up to Large Vocabulary Image Annotation,” in *Proc. 22nd Int. Joint Conf. Artif. Intell. (IJCAI)*, 2011, pp. 2764–2770.
- [22] J. F. Gemmeke, D. P. W. 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 *Proc. IEEE Int. Conf. Acoustic., Speech and Signal Process. (ICASSP)*, 2017, pp. 776–780.