

COMBINING MULTIPLE DISTRIBUTIONS BASED ON SUB-CLUSTER ADACOS FOR ANOMALOUS SOUND DETECTION UNDER DOMAIN SHIFTED CONDITIONS

Kevin Wilkinghoff

Fraunhofer Institute for Communication, Information Processing and Ergonomics FKIE
 Fraunhoferstraße 20, 53343 Wachtberg, Germany
 kevin.wilkinghoff@fkie.fraunhofer.de

ABSTRACT

Systems based on sub-cluster AdaCos yield state-of-the-art performance on the DCASE 2020 dataset for anomalous sound detection. In contrast to the previous year, the dataset belonging to task 2 “Unsupervised Anomalous Sound Detection for Machine Condition Monitoring under Domain Shifted Conditions” of the DCASE challenge 2021 contains not only source domains with 1000 normal training samples for each machine but also so-called target domains with different acoustic conditions for which only 3 normal training samples are available. To address this additional problem, a novel anomalous sound detection system based on sub-cluster AdaCos for the DCASE challenge 2021 is presented. The system is trained to extract embeddings whose distributions are estimated in different ways for source and target domains, and utilize the resulting negative log-likelihoods as anomaly scores. In experimental evaluations, it is shown that the presented system significantly outperforms both baseline systems on the source and target domains of the development set. On the evaluation set of the challenge, the proposed system is ranked third among all 27 teams’ submissions.

Index Terms— anomalous sound detection, machine listening, representation learning, angular margin loss, domain shift

1. INTRODUCTION

The goal of semi-supervised anomalous sound detection is to decide whether a given audio sample resembles the training data i.e. is *normal* or substantially differs from the training data and thus is *anomalous*. Basically, one can distinguish two major strategies for anomalous sound detection: The first approach is based on training autoencoders to encode normal data into a lower-dimensional space and then reconstruct it again [1, 2]. The underlying assumption is that normal data can be reconstructed well after training while anomalous data cannot, leading to a higher reconstruction error. Thus, the reconstruction error can be used as an anomaly score. The second approach is to train neural networks to discriminate among classes as for example machine types and utilize the trained neural network to extract representation of the data, so-called embeddings, as features [3, 4, 5, 6, 7]. Here, the assumption is that the information needed to discriminate among the classes and thus is contained in the embeddings is also sufficient to distinguish normal from anomalous data. Angular margin losses such as ArcFace [8] or AdaCos [9], which ensure a margin between the classes, have been shown to outperform standard softmax losses in this context. To our knowledge, the best performing system on the anomalous sound detection dataset belonging to task 2 of the DCASE challenge 2020 [10] uses an extension of AdaCos, called sub-cluster AdaCos [11].

This loss learns more than a single mean value for each class to estimate less restrictive distributions of the embeddings than standard AdaCos and utilizes Gaussian mixture models (GMMs) to estimate these distributions for the normal data instead of comparing embeddings to the learned mean values by using the cosine similarity. This superior performance is the reason why this work focuses entirely on a system based on the sub-cluster AdaCos loss.

The system presented in this paper is designed for and submitted to task 2 “Unsupervised Anomalous Sound Detection for Machine Condition Monitoring under Domain Shifted Conditions” of the DCASE challenge 2021 [12]. The dataset of this task consists of audio recordings with a length of 10 seconds and a sampling rate of 16 kHz belonging to the machine types “ToyCar” and “ToyTrain” from ToyADMOS2 [13] and the machines types “fan”, “gearbox”, “pump”, “slide rail” and “valve” from MIMII DUE [14]. The organizers of the challenge also provided two baseline systems: An autoencoder, which is the same as the baseline system of the previous edition of the task, and a discriminatively trained MobileNetV2-based baseline that predicts the section, which is a subset of the data within one machine type, a given audio sample belongs to.

In contrast to the DCASE challenge 2020, there are several differences for this year’s task: First and foremost, the dataset is split into source domains for which about 1000 normal training samples are provided for each of the 6 sections per machine type and so-called target domains for the same sections with different acoustic conditions than the source domains for which only 3 normal training samples are available. For both domains, the same number of test samples is provided, about 100 normal samples and 100 anomalous samples. Furthermore, the dataset is split into a development set consisting of half of the sections and an evaluation set consisting of the other half of the sections. Another difference between the datasets is, that the sections do not directly correspond to specific products of a machine type but the same products can appear in different sections or different products can appear in the same sections. Both of these changes make the task much more challenging than before. Last but not least, the DCASE 2020 dataset consists of slightly different machine types, namely “ToyCar” and “ToyConveyor” from the ToyADMOS dataset [15] and the machine types “fan”, “pump”, “slide rail” and “valve” from the MIMII dataset [16].

The goal of this work is to investigate how to utilize the sub-cluster AdaCos loss for the DCASE 2021 anomalous sound detection dataset with its novel challenges. To this end, a system based on the sub-cluster AdaCos loss is presented. As a second contribution, different ways to compute anomaly scores for the source and target domains are proposed. Furthermore, it is shown how to decide whether samples are normal or anomalous only based on these

layer name	structure	output size
input	-	313×128
2D convolution	7×7 , stride= 2	$157 \times 64 \times 16$
residual block	$\begin{pmatrix} 3 \times 3 \\ 3 \times 3 \end{pmatrix} \times 2$, stride= 1	$78 \times 31 \times 16$
residual block	$\begin{pmatrix} 3 \times 3 \\ 3 \times 3 \end{pmatrix} \times 2$, stride= 1	$39 \times 16 \times 32$
residual block	$\begin{pmatrix} 3 \times 3 \\ 3 \times 3 \end{pmatrix} \times 2$, stride= 1	$20 \times 8 \times 64$
residual block	$\begin{pmatrix} 3 \times 3 \\ 3 \times 3 \end{pmatrix} \times 2$, stride= 1	$10 \times 4 \times 128$
max pooling	10×1 , stride= 1	4×128
flatten	-	512
dense (representation)	linear	128
sub-cluster AdaCos	-	42
sub-cluster AdaCos	-	199

Table 1: Modified ResNet architecture used for all experiments.

scores from normal data. In experimental evaluations, it is shown that the proposed system significantly outperforms both baseline systems on the source and target domains of the development set.

2. METHODOLOGY

2.1. Sub-Cluster AdaCos loss

The standard AdaCos [9] loss is an angular margin loss that uses an adaptive scale parameter and thus does not require any hyperparameter to be set. For detecting anomalous sounds, this loss has been extended to the sub-cluster AdaCos loss [11] that learns multiple mean values for each class instead of a single one. This loss has been shown to significantly outperform the standard AdaCos loss on the DCASE 2020 dataset and thus is the only loss used for training the proposed system. The probability of sample $x_i \in \mathbb{R}^D$ belonging to class j of the $C \in \mathbb{N}$ classes is given by

$$\hat{P}_{i,j} := \sum_{l \in \mathcal{M}^{(j)}} \frac{\exp(\hat{s} \cdot \cos \theta_{i,l})}{\sum_{k=1}^{CS} \exp(\hat{s} \cdot \cos \theta_{i,k})} \quad (1)$$

where $\mathcal{M}^{(j)}$ denotes all sub-clusters belonging to class j , $S \in \mathbb{N}$ denotes the number of sub-clusters, \hat{s} denotes the adaptive scale parameter as defined in [11], and $\theta_{i,k} \in [0, \pi]$ denotes the angle defined through the cosine similarity $\cos \theta_{i,k} = \langle x_i, W_k \rangle / \|x_i\| \|W_k\|$ for a learned class center $W_k \in \mathbb{R}^D$.

2.2. Data preprocessing

To compute input features for the neural network, log-Mel spectrograms with 128 Mel-bins, a window size of 1024 and a hop size of 512 are extracted from all raw waveforms with a sampling rate of 16 kHz resulting in features of size 313×128 . These features are then standardized by subtracting the temporal mean and dividing by the temporal standard deviation estimated from all training files.

2.3. Neural network architecture

The network architecture used throughout this work is the same as used in [11] and can be found in Tab. 2.3. It consists of several residual blocks [17] whose output is further processed by max-pooling over time, flattening and using a final linear dense layer

to obtain the embeddings of size $D = 128$. In each residual block, batch normalization [18] is applied and LeakyReLU [19] with $\alpha = 0.1$ is used as a non-linear transfer function.

To train the neural network, two sub-cluster AdaCos losses with equal weights are minimized using Adam [20]. One is for classifying jointly among the sections and machine types and the other one for classifying among the different attribute information given in the filenames. When training, all normal data contained in the training set and the additional training set has been used resulting in a total of 42 sections and 199 different types of attribute information. Furthermore, mixup [21] is used during training to avoid overfitting of the model to the training data. The network is implemented in Tensorflow [22] and trained for 400 epochs with a batch size of 64.

2.4. Calculating anomaly scores

Throughout this work, all anomaly scores are computed by training Gaussian mixture models (GMMs) on the embeddings and utilizing negatively weighted log-likelihoods as scores. In [11], it has been shown that using GMMs to estimate the underlying distribution of the embeddings outperforms other backends such as using cosine similarity to the class means. Unless stated otherwise, all GMMs are realized using scikit-learn [23], initialized with the learned mean values of the sub-cluster AdaCos loss, and have a regularized covariance matrix by adding 10^{-3} to the diagonal. To calculate the anomaly scores, two different strategies for the source and target domain are used.

For the source domain, one GMM is trained for all normal data of the source domain belonging to a single section. Another GMM is trained for all normal data of the source domain labeled with different attribute information. Let $x \in \mathbb{R}^D$ denote an embedding, $s(x) \in \mathcal{S}$ denote its section and $a(s(x)) \subset \mathcal{A}$ denote all attribute information that are present in this section. Then, the anomaly score $Z_{\text{source}}(x)$ for x is the given by

$$Z_{\text{source}}(x) := - \max_k \log P(x|s(x), k) - \max_k \max_{a \in a(s(x))} \log P(x|a, k) \quad (2)$$

where $P(\cdot|s, k)$, $P(\cdot|a, k)$ denote the weighted likelihoods of component k of the GMMs trained for section $s \in \mathcal{S}$ and target information $a \in \mathcal{A}$, respectively.

For the target domain, the same GMMs trained on the normal data of the target domain belonging to single sections are used. Furthermore, another GMM with three components is trained on the three target samples and thus the cosine distance to the closest normal target sample is also utilized. Using the above notation, the anomaly score $Z_{\text{target}}(x)$ for embedding x is given by

$$Z_{\text{target}}(x) := - \max_k \log P(x|s(x), k) - \max_{k=1,2,3} \log P(x|X_{\text{target}}(s(x)), k) \quad (3)$$

where $X_{\text{target}}(s(x)) \subset \mathbb{R}^D$ denotes the normal training samples of the target domain belonging to the section of x .

In [11], it has also been shown that a simple representation derived from the input features leads to surprisingly good performance for the machine type ‘‘valve’’ on the DCASE 2020 dataset. This is the reason why an additional term based on the temporal maximum of the log-Mel spectrogram, denoted by $t_{\text{max}}(x) \in \mathbb{R}^D$ for embedding x , is introduced when calculating the anomaly score for the source domain of the machine type ‘‘valve’’. To this end, a

GMM with a single Gaussian component is trained and the altered anomaly score $\tilde{Z}_{\text{source}}(x)$ for embedding x belonging to machine type “valve” is given by

$$\tilde{Z}_{\text{source}}(x) := Z_{\text{source}}(x) - \max_{a \in \mathcal{A}(s(x))} \log P_{t_{\max}}(t_{\max}(x)|a) \quad (4)$$

where $P_{t_{\max}}(\cdot|a)$ denotes the weighted likelihoods of the single Gaussian trained on the temporal maxima of the log-Mel spectrograms belonging to target information $a \in \mathcal{A}$.

2.5. Ensembling strategy

As done in [11], the proposed neural network for extracting the embeddings is trained with a different number of sub-clusters ranging from 2^0 to 2^4 . The same value is used for both sub-cluster AdaCos losses. Thus, there are 5 differently trained versions of each network to extract embeddings. Furthermore, after each 100 epochs of training, the embeddings are extracted and GMMs are trained to calculate the anomaly detection scores. Then, all of these scores are summed-up resulting in 4 subsystems for each network with a specified number of sub-clusters and hence an ensemble consisting of a total of $4 \times 5 = 20$ models.

In addition to that, the described ensembling procedure is repeated by using only a single sub-cluster AdaCos loss classifying among the sections and machine types only and thus removing the second sub-cluster AdaCos loss. This led to slightly better performance for some machine types and to slightly worse performance for other machine types. To obtain anomaly scores for each machine type, the single system is used that led to better performance for the given machine type. More concretely, for the machine types “ToyCar”, “ToyTrain”, “pump” and “slide rail” the anomaly scores obtained by using the model trained on both losses are used and for “fan”, “gearbox” and “valve” the anomaly scores obtained with the models trained on only a single loss are used.

2.6. Setting decision thresholds

Next, it is described how decision thresholds for deciding whether a given test sample is normal or anomalous solely based on the anomaly score are obtained. To this end, the 90th percentile of the anomaly scores of all normal training samples belonging to a given section and a given domain is calculated. Then, all anomaly scores of test samples belonging to the same section and domain that are above this threshold are marked as *anomalous*. For the source domain, $Z_{\text{source}}(x)$ as defined in Eq. (1) is used but for the target domain, only the first term of $Z_{\text{target}}(x)$ is used. The reason is that the likelihoods from the second term belonging to the training data are inappropriately high since the three means of the corresponding GMM are initialized as the three training samples. Hence, when also using the second term of Eq. (2) the decision threshold would also be overestimated and thus all test data samples belonging to the target domain would be considered anomalous.

3. EXPERIMENTAL RESULTS

3.1. Performance on the development set

The experimental results obtained with the proposed system on the development set compared to both baseline systems can be found in Tab. 2. It can be seen that the proposed system significantly outperforms both baseline systems, which both have roughly

machine type	dataset split	section	domain	baseline				proposed system	
				autoencoder		MobileNetV2		AUC	pAUC
ToyCar	0	source		67.63%	51.87%	66.56%	<u>66.47%</u>	<u>75.03%</u>	60.32%
ToyCar	1	source		61.97%	51.82%	71.58%	66.44%	89.30%	<u>73.26%</u>
ToyCar	2	source		74.36%	55.56%	40.37%	47.48%	<u>92.27%</u>	<u>75.32%</u>
ToyCar	0	target		54.50%	50.52%	61.32%	52.61%	<u>92.41%</u>	<u>78.11%</u>
ToyCar	1	target		64.12%	51.14%	72.48%	<u>63.99%</u>	<u>78.34%</u>	61.47%
ToyCar	2	target		56.57%	52.61%	45.17%	48.85%	<u>56.78%</u>	<u>57.26%</u>
ToyCar		harmonic mean		62.49%	52.36%	56.04%	56.37%	<u>78.37%</u>	<u>66.64%</u>
ToyTrain	0	source		72.67%	69.38%	69.84%	54.43%	<u>95.62%</u>	<u>88.74%</u>
ToyTrain	1	source		72.65%	62.52%	64.79%	54.09%	<u>90.67%</u>	<u>76.26%</u>
ToyTrain	2	source		69.91%	47.48%	69.28%	<u>47.66%</u>	<u>87.78%</u>	47.37%
ToyTrain	0	target		56.07%	50.62%	46.28%	<u>51.27%</u>	<u>71.32%</u>	48.47%
ToyTrain	1	target		51.13%	48.60%	<u>53.38%</u>	49.60%	<u>50.89%</u>	<u>50.89%</u>
ToyTrain	2	target		55.57%	50.79%	51.42%	53.40%	<u>92.98%</u>	<u>76.58%</u>
ToyTrain		harmonic mean		61.71%	53.81%	57.46%	51.61%	<u>77.17%</u>	<u>60.71%</u>
fan	0	source		66.69%	57.08%	43.62%	50.45%	<u>72.17%</u>	<u>62.95%</u>
fan	1	source		67.43%	50.72%	78.33%	78.37%	<u>89.24%</u>	<u>84.68%</u>
fan	2	source		64.21%	53.12%	74.21%	<u>76.80%</u>	<u>83.03%</u>	74.58%
fan	0	target		<u>69.70%</u>	55.13%	53.34%	<u>56.01%</u>	55.30%	48.47%
fan	1	target		49.99%	48.49%	78.12%	66.41%	<u>87.68%</u>	74.58%
fan	2	target		66.19%	56.93%	60.35%	60.97%	<u>72.43%</u>	<u>70.63%</u>
fan		harmonic mean		63.24%	53.38%	61.56%	63.02%	<u>74.66%</u>	<u>67.34%</u>
gearbox	0	source		56.03%	51.59%	81.35%	70.46%	<u>85.80%</u>	<u>74.56%</u>
gearbox	1	source		72.77%	52.30%	60.74%	<u>53.88%</u>	<u>85.37%</u>	52.54%
gearbox	2	source		58.96%	51.82%	<u>71.58%</u>	<u>62.23%</u>	61.39%	48.23%
gearbox	0	target		74.29%	55.67%	75.02%	64.77%	<u>81.93%</u>	<u>69.25%</u>
gearbox	1	target		72.12%	51.78%	56.27%	<u>53.30%</u>	<u>86.02%</u>	51.29%
gearbox	2	target		<u>66.41%</u>	53.66%	64.45%	<u>55.58%</u>	65.26%	49.92%
gearbox		harmonic mean		65.97%	52.76%	66.70%	<u>59.16%</u>	<u>76.13%</u>	56.00%
pump	0	source		67.48%	61.83%	64.09%	62.40%	<u>77.49%</u>	<u>63.47%</u>
pump	1	source		82.38%	58.29%	86.27%	66.66%	<u>98.26%</u>	<u>91.21%</u>
pump	2	source		63.93%	55.44%	53.70%	50.98%	<u>78.56%</u>	<u>63.68%</u>
pump	0	target		58.01%	51.53%	<u>59.09%</u>	<u>53.96%</u>	<u>57.00%</u>	51.74%
pump	1	target		47.35%	49.65%	71.86%	62.69%	<u>88.82%</u>	<u>62.42%</u>
pump	2	target		62.78%	51.67%	50.16%	51.69%	<u>72.88%</u>	<u>57.26%</u>
pump		harmonic mean		61.92%	54.41%	61.89%	57.37%	<u>76.59%</u>	<u>63.00%</u>
slide rail	0	source		74.09%	52.45%	61.51%	53.97%	<u>96.28%</u>	<u>83.16%</u>
slide rail	1	source		82.16%	60.29%	79.97%	55.62%	<u>93.73%</u>	<u>69.16%</u>
slide rail	2	source		78.34%	65.16%	79.86%	71.88%	<u>84.01%</u>	<u>77.30%</u>
slide rail	0	target		67.22%	57.32%	51.96%	51.96%	<u>82.06%</u>	<u>60.63%</u>
slide rail	1	target		<u>66.94%</u>	<u>53.08%</u>	46.83%	52.02%	62.92%	49.76%
slide rail	2	target		46.20%	50.10%	55.61%	55.71%	<u>72.39%</u>	<u>55.72%</u>
slide rail		harmonic mean		66.74%	55.94%	59.26%	56.00%	<u>80.16%</u>	<u>63.86%</u>
valve	0	source		50.34%	50.82%	58.34%	54.97%	<u>81.47%</u>	<u>63.00%</u>
valve	1	source		53.52%	49.33%	53.57%	50.09%	<u>90.09%</u>	<u>64.37%</u>
valve	2	source		59.91%	51.96%	56.13%	51.69%	<u>98.07%</u>	<u>91.47%</u>
valve	0	target		47.12%	48.68%	52.19%	51.54%	<u>65.86%</u>	<u>64.74%</u>
valve	1	target		56.39%	53.88%	68.59%	57.83%	<u>81.38%</u>	<u>58.58%</u>
valve	2	target		55.16%	48.97%	53.58%	50.86%	<u>77.42%</u>	<u>56.95%</u>
valve		harmonic mean		53.41%	50.54%	56.51%	52.64%	<u>81.13%</u>	<u>64.92%</u>
all		harmonic mean		61.93%	53.27%	59.72%	56.37%	<u>77.69%</u>	<u>62.99%</u>

Table 2: AUCs and pAUCs (with $p = 0.1$) obtained with the baseline systems and the proposed system on the development set. High-est AUCs and pAUCs in each row are underlined.

the same overall performance, on the source and target domains. However, the improvement in terms of AUC is much greater than for pAUC. For nearly all dataset splits the proposed system has a higher AUC than both baseline systems. But for some dataset splits the MobileNetV2-based baseline system has a higher pAUC than the proposed system. For the machine type “gearbox” the harmonic mean of all pAUCs belonging to the proposed system is even slightly worse than the harmonic mean of the MobileNetV2-based baseline system.

Next, the performances of all distributions and losses used for the subsystems of the ensemble have been evaluated and compared on the development set. The results can be found in Tab. 3 and Tab. 4, respectively. For the distributions it can be seen that in most cases the distribution conditioned on the sections has a lower AUC but a higher pAUC than the distribution conditioned on the file endings and the distribution conditioned on the target samples. One exception to this is the machine type “gearbox” for which $P(\cdot|s, k)$ performs best and “ToyCar” and “ToyConveyor”

dataset split		$P(\cdot s, k)$		distribution $P(\cdot a, k)$		$P(\cdot X_{\text{target}}(s), k)$		proposed ensemble	
machine type	domain	AUC	pAUC	AUC	pAUC	AUC	pAUC	AUC	pAUC
ToyCar	source	82.49%	65.82%	<u>86.22%</u>	<u>69.39%</u>	-	-	85.04%	69.25%
ToyCar	target	65.02%	58.45%	-	-	<u>76.02%</u>	64.13%	73.21%	<u>64.55%</u>
ToyTrain	source	90.54%	65.60%	<u>91.55%</u>	<u>66.11%</u>	-	-	91.30%	66.01%
ToyTrain	target	54.03%	53.21%	-	-	<u>68.13%</u>	<u>57.29%</u>	66.81%	56.14%
fan	source	<u>81.05%</u>	<u>73.34%</u>	80.55%	72.12%	-	-	80.98%	73.20%
fan	target	67.52%	<u>65.71%</u>	-	-	69.17%	60.20%	<u>69.30%</u>	62.20%
gearbox	source	75.20%	58.63%	75.66%	55.65%	-	-	75.88%	56.46%
gearbox	target	<u>81.02%</u>	<u>57.93%</u>	-	-	75.68%	55.34%	76.64%	55.59%
pump	source	83.63%	70.60%	<u>83.83%</u>	<u>70.83%</u>	-	-	83.70%	70.66%
pump	target	69.90%	<u>60.11%</u>	-	-	70.16%	55.51%	<u>70.52%</u>	56.81%
slide rail	source	91.10%	<u>76.23%</u>	90.97%	75.59%	-	-	91.11%	76.20%
slide rail	target	65.90%	<u>56.17%</u>	-	-	<u>71.63%</u>	-	-	54.94%
valve	source	83.85%	<u>73.84%</u>	<u>89.34%</u>	68.89%	-	-	89.33%	70.77%
valve	target	71.66%	<u>61.64%</u>	-	-	73.31%	59.16%	<u>74.26%</u>	60.00%
all	source	83.67%	<u>68.66%</u>	<u>85.09%</u>	67.81%	-	-	85.00%	68.38%
all	target	67.00%	<u>58.80%</u>	-	-	<u>71.90%</u>	57.93%	71.63%	58.41%

Table 3: Harmonic means of AUCs and pAUCs (with $p = 0.1$) obtained with different distributions on the development set. Highest AUCs and pAUCs in each row are underlined.

dataset split		sub-cluster AdaCos losses for				proposed ensemble	
machine type	domain	sections		sections and file endings		AUC	pAUC
		AUC	pAUC	AUC	pAUC		
ToyCar	source	74.49%	60.67%	<u>88.58%</u>	<u>71.58%</u>	85.04%	69.25%
ToyCar	target	67.44%	61.55%	75.53%	65.84%	73.21%	64.55%
ToyTrain	source	86.88%	62.00%	91.15%	<u>67.03%</u>	<u>91.30%</u>	66.01%
ToyTrain	target	63.73%	54.18%	<u>68.10%</u>	<u>56.43%</u>	66.81%	56.14%
fan	source	<u>81.66%</u>	73.16%	80.22%	72.60%	80.98%	<u>73.20%</u>
fan	target	69.11%	<u>62.55%</u>	69.00%	61.16%	<u>69.30%</u>	62.20%
gearbox	source	74.97%	<u>57.04%</u>	75.34%	56.16%	<u>75.88%</u>	56.46%
gearbox	target	<u>78.53%</u>	<u>59.53%</u>	72.02%	51.90%	76.64%	55.59%
pump	source	83.44%	69.75%	<u>83.75%</u>	<u>71.35%</u>	83.70%	70.66%
pump	target	68.77%	56.31%	<u>71.30%</u>	56.34%	70.52%	<u>56.81%</u>
slide rail	source	90.44%	74.91%	90.99%	<u>76.43%</u>	<u>91.11%</u>	76.20%
slide rail	target	68.84%	54.44%	<u>73.13%</u>	<u>55.41%</u>	71.56%	54.94%
valve	source	88.90%	<u>72.11%</u>	<u>89.51%</u>	68.95%	89.33%	70.77%
valve	target	<u>75.41%</u>	<u>61.12%</u>	72.36%	56.55%	74.26%	60.00%
all	source	82.54%	66.44%	<u>85.26%</u>	<u>68.58%</u>	85.00%	68.38%
all	target	69.96%	58.34%	71.56%	57.37%	<u>71.63%</u>	<u>58.41%</u>

Table 4: Harmonic means of AUCs and pAUCs (with $p = 0.1$) obtained with different losses on the development set. Highest AUCs and pAUCs in each row are underlined.

for which $P(\cdot|s, k)$ yields significantly lower AUCs and pAUCs. The ensemble obtained by taking the mean of the distributions performs relatively close to the best performing distribution for each machine type. For the two losses the same three machine types stated before have the largest differences in performance. When only using the sections for the loss, the AUCs and pAUCs are higher for the machine type “gearbox” but lower for machine types “ToyCar” and “ToyConveyor” than when also using the file endings for the loss. For each of the other machine types, no loss is preferable since both perform better than the other one for some classes but the differences in performance are relatively small compared to the three machine types mentioned before. Again, the ensemble yields results close to the best performing loss for each machine type and thus seems to combine the strengths of both losses.

3.2. Performance on the evaluation set

The experimental results obtained on the evaluation set compared to both baseline systems can be found in Tab. 5. Overall, it can be seen that the proposed system significantly outperforms both baseline

dataset split		baseline				proposed system	
machine type	domain	autoencoder		MobileNetV2		AUC	pAUC
		AUC	pAUC	AUC	pAUC		
ToyCar	source	<u>76.33%</u>	51.26%	34.32%	53.49%	67.07%	<u>63.05%</u>
ToyCar	target	58.02%	53.42%	56.62%	58.89%	<u>72.83%</u>	<u>63.77%</u>
ToyTrain	source	69.89%	55.49%	47.30%	52.49%	<u>70.87%</u>	<u>56.19%</u>
ToyTrain	target	<u>67.18%</u>	<u>59.78%</u>	39.27%	48.75%	48.38%	52.39%
fan	source	66.58%	51.36%	70.88%	57.76%	89.07%	<u>69.85%</u>
fan	target	55.74%	49.68%	59.96%	58.53%	<u>88.89%</u>	<u>70.55%</u>
gearbox	source	<u>67.81%</u>	<u>55.71%</u>	53.16%	53.47%	61.19%	50.97%
gearbox	target	<u>63.32%</u>	<u>58.06%</u>	49.27%	49.83%	54.68%	49.40%
pump	source	62.75%	51.18%	67.12%	60.77%	<u>70.89%</u>	<u>65.52%</u>
pump	target	54.43%	50.79%	68.85%	59.79%	<u>79.20%</u>	<u>67.81%</u>
slide rail	source	64.13%	50.91%	73.06%	60.47%	<u>88.06%</u>	<u>64.38%</u>
slide rail	target	51.65%	51.92%	72.78%	60.94%	<u>85.66%</u>	<u>69.69%</u>
valve	source	51.56%	50.89%	54.71%	53.03%	<u>73.19%</u>	<u>55.97%</u>
valve	target	52.19%	49.27%	51.64%	50.10%	<u>54.90%</u>	<u>51.47%</u>
all	source	64.76%	52.32%	53.82%	55.73%	<u>73.13%</u>	<u>60.21%</u>
all	target	57.03%	53.01%	54.80%	54.80%	<u>65.76%</u>	<u>59.47%</u>
all	both	56.38%		54.77%		64.20%	

Table 5: Harmonic means of AUCs and pAUCs (with $p = 0.1$) obtained with the baseline systems and the proposed system on the evaluation set. Highest AUCs and pAUCs in each row are underlined.

systems. However, for some machine types as for example “gearbox” or the target domains of “ToyTrain”, the autoencoder-based baseline system has a much better performance than the proposed system. This shows that not all information needed to detect anomalies is captured sufficiently well in the discriminatively trained embeddings but an autoencoder is able to identify these anomalous structures. Hence, also utilizing an autoencoder structure for training the embeddings may be helpful to further improve the results.

4. CONCLUSIONS AND FUTURE WORK

In this work, a system for detecting anomalous sounds based on the sub-cluster AdaCos loss function for domain shifted conditions has been presented. The proposed system consists of multiple discriminatively trained neural networks for extracting embeddings from log-Mel spectrograms and utilizes multiple GMMs for estimating distributions of the normal embeddings. These estimated distributions are then used to calculate log-likelihoods for test data that are combined into actual anomaly scores. To decide whether a given test sample is anomalous or normal, individual decision thresholds for each section are computed by taking the 90th percentile from the log-likelihoods of the corresponding normal training samples. In experimental evaluations conducted on the dataset of task 2 of the DCASE challenge 2021, it has been shown that the proposed system significantly outperforms both baseline systems of the challenge in terms of AUC and pAUC on source and target domains of the development and evaluation set. Moreover, the system ranked third among all 27 teams that participated in this challenge.

For the near future, it is planned to reduce the size of the ensemble and thus simplify the presented system. This can be achieved by identifying subsystems that only result in marginal gains and therefore can be removed from the ensemble. Furthermore, the performance can probably be improved by incorporating an autoencoder into the proposed system. Instead of building an even larger ensemble, it seems promising to use an additional reconstruction loss and therefore enforce the embeddings to capture more information about the structure of the data as for example done in [24, 25].

5. REFERENCES

- [1] S. Kapka, “ID-conditioned auto-encoder for unsupervised anomaly detection,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2020, pp. 71–75.
- [2] K. Miyazaki, T. Komatsu, T. Hayashi, S. Watanabe, T. Toda, and K. Takeda, “Conformer-based sound event detection with semi-supervised learning and data augmentation,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2020, pp. 100–104.
- [3] T. Inoue, P. Vinayavekchin, S. Morikuni, S. Wang, T. Hoang Trong, D. Wood, M. Tatsubori, and R. Tachibana, “Detection of anomalous sounds for machine condition monitoring using classification confidence,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2020, pp. 66–70.
- [4] R. Giri, S. V. Tenneti, F. Cheng, K. Helwani, U. Isik, and A. Krishnaswamy, “Self-supervised classification for detecting anomalous sounds,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2020, pp. 46–50.
- [5] Q. Zhou, “ArcFace based sound mobilenets for DCASE 2020 task 2,” DCASE2020 Challenge, Tech. Rep., 2020.
- [6] J. A. Lopez, H. Lu, P. Lopez-Meyer, L. Nachman, G. Stemmer, and J. Huang, “A speaker recognition approach to anomaly detection,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2020, pp. 96–99.
- [7] K. Wilkinghoff, “Using look, listen, and learn embeddings for detecting anomalous sounds in machine condition monitoring,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2020, pp. 215–219.
- [8] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, “ArcFace: Additive angular margin loss for deep face recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2019, pp. 4690–4699.
- [9] X. Zhang, R. Zhao, Y. Qiao, X. Wang, and H. Li, “AdaCos: Adaptively scaling cosine logits for effectively learning deep face representations,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2019, pp. 10 823–10 832.
- [10] Y. Koizumi, Y. Kawaguchi, K. Imoto, T. Nakamura, Y. Nikaido, R. Tanabe, H. Purohit, K. Suefusa, T. Endo, M. Yasuda, and N. Harada, “Description and discussion on DCASE2020 challenge task2: Unsupervised anomalous sound detection for machine condition monitoring,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, 2020, pp. 81–85.
- [11] K. Wilkinghoff, “Sub-cluster AdaCos: Learning representations for anomalous sound detection,” in *International Joint Conference on Neural Networks (IJCNN)*, 2021.
- [12] Y. Kawaguchi, K. Imoto, Y. Koizumi, N. Harada, D. Niizumi, K. Dohi, R. Tanabe, H. Purohit, and T. Endo, “Description and discussion on DCASE 2021 challenge task 2: Unsupervised anomalous sound detection for machine condition monitoring under domain shifted conditions,” in *arXiv e-prints: 2106.04492, 1–5*, 2021.
- [13] N. Harada, D. Niizumi, D. Takeuchi, Y. Ohishi, M. Yasuda, and S. Saito, “ToyADMOS2: Another dataset of miniature-machine operating sounds for anomalous sound detection under domain shift conditions,” *arXiv preprint arXiv:2106.02369*, 2021.
- [14] R. Tanabe, H. Purohit, K. Dohi, T. Endo, Y. Nikaido, T. Nakamura, and Y. Kawaguchi, “MIMII DUE: Sound dataset for malfunctioning industrial machine investigation and inspection with domain shifts due to changes in operational and environmental conditions,” in *arXiv e-prints: 2006.05822, 1–4*, 2021.
- [15] Y. Koizumi, S. Saito, H. Uematsu, N. Harada, and K. Imoto, “ToyADMOS: A dataset of miniature-machine operating sounds for anomalous sound detection,” in *Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2019, pp. 313–317.
- [16] H. Purohit, R. Tanabe, T. Ichige, T. Endo, Y. Nikaido, K. Suefusa, and Y. Kawaguchi, “MIMII Dataset: Sound dataset for malfunctioning industrial machine investigation and inspection,” in *Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*. New York University, 2019, pp. 209–213.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016, pp. 770–778.
- [18] S. Ioffe and C. Szegedy, “Batch normalization: accelerating deep network training by reducing internal covariate shift,” in *32nd International Conference on Machine Learning (ICML)*, vol. 37, 2015, pp. 448–456.
- [19] A. L. Maas, A. Y. Hannun, and A. Y. Ng, “Rectifier nonlinearities improve neural network acoustic models,” in *30th International Conference on Machine Learning (ICML)*, 2013.
- [20] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations (ICLR)*, Y. Bengio and Y. LeCun, Eds., 2015.
- [21] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, “Mixup: Beyond empirical risk minimization,” in *International Conference on Learning Representations (ICLR)*, 2018.
- [22] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, *et al.*, “Tensorflow: A system for large-scale machine learning,” in *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2016, pp. 265–283.
- [23] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [24] X. Cai, H. Dinkel, Z. Yan, Y. Wang, J. Zhang, and Y. Wang, “The small rice camera ready submission to the DCASE2021,” DCASE2021 Challenge, Tech. Rep., 2021.
- [25] H. Narita and A. Tamamori, “Unsupervised anomalous sound detection using intermediate representation of trained models and metric learning based variational autoencoder,” DCASE2021 Challenge, Tech. Rep., 2021.