

# DBT: A DETECTION BOOSTER TRAINING METHOD FOR IMPROVING THE ACCURACY OF CLASSIFIERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1 Deep learning models owe their success at large, to the availability of a large  
2 amount of annotated data. They try to extract features from the data that contain  
3 useful information needed to improve their performance on target applications.  
4 Most works focus on directly optimizing the target loss functions to improve the  
5 accuracy by allowing the model to implicitly learn representations from the data.  
6 There has not been much work on using background/noise data to estimate the  
7 statistics of in-domain data to improve the feature representation of deep neural  
8 networks. In this paper, we probe this direction by deriving a relationship between  
9 the estimation of unknown parameters of the probability density function (pdf)  
10 of input data and classification accuracy. Using this relationship, we show that  
11 having a better estimate of the unknown parameters using background and in-  
12 domain data provides better features which leads to better accuracy. Based on  
13 this result, we introduce a simple but effective detection booster training (DBT)  
14 method that applies a detection loss function on the early layers of a neural network  
15 to discriminate in-domain data points from noise/background data, to improve  
16 the classifier accuracy. The background/noise data comes from the same family  
17 of pdfs of input data but with different parameter sets (e.g., mean, variance). In  
18 addition, we also show that our proposed DBT method improves the accuracy even  
19 with limited labeled in-domain training samples as compared to normal training.  
20 We conduct experiments on face recognition, image classification, and speaker  
21 classification problems and show that our method achieves superior performance  
22 over strong baselines across various datasets and model architectures.

## 23 1 INTRODUCTION

24 Modern pattern recognition systems achieve outstanding accuracies on a vast domain of challenging  
25 computer vision, natural language, and speech recognition benchmarks (Russakovsky et al. (2015);  
26 Lin et al. (2014); Everingham et al. (2015); Panayotov et al. (2015)). The success of deep learning  
27 approaches relies on the availability of a large amount of annotated data and on extracting useful  
28 features from them for different applications. Learning rich feature representations from the available  
29 data is a challenging problem in deep learning. A related line of work includes learning deep latent  
30 space embedding through deep generative models (Kingma & Welling (2014); Goodfellow et al.  
31 (2014); Berthelot et al. (2019) or using self-supervised learning methods (Noroozi & Favaro (2016);  
32 Gidaris et al. (2018); Zhang et al. (2016b)) or through transfer learning approaches (Yosinski et al.  
33 (2014); Oquab et al. (2014); Razavian et al. (2014)).

34 In this paper, we propose to use a different approach to improve the feature representations of deep  
35 neural nets and eventually improve their accuracy by estimating the unknown parameters of the  
36 probability density function (pdf) of input data. Parameter estimation or Point estimation methods  
37 are well studied in the field of statistical inference (Lehmann & Casella (1998)). The insights from  
38 the theory of point estimation can help us to develop better deep model architectures for improving  
39 the model's performance. We make use of this theory to derive a correlation between the estimation  
40 of unknown parameters of pdf and classifier outputs. However, directly estimating the unknown  
41 pdf parameters for practical problems such as image classification is not feasible since it can sum  
42 up to millions of parameters. In order to overcome this bottleneck, we assume that the input data  
43 points are sampled from a family of pdfs instead of a single pdf and propose to use a detection  
44 based training approach to better estimate the unknowns using in-domain and background/noise data.  
45 One alternative is that we can use generative models for this task, however, they mimic the general

46 distribution of training data conditioned on random latent vectors and hence cannot be directly applied  
 47 for estimating the unknown parameters of a family of pdfs. Our proposed detection method involves  
 48 a binary class discriminator that separates the target data points from noise or background data. The  
 49 noise or background data is assumed to come from the same family of distribution of in-domain  
 50 data but with different moments (Please refer to the appendix for more details about the family of  
 51 distributions and its extension to a general structure). In image classification, this typically represents  
 52 the background patches from input data that fall under the same distribution family. In speech domain,  
 53 it can be random noise or the silence intervals in speech data. Collecting such background data to  
 54 improve the feature representations is much simpler as compared to using labeled training data since  
 55 it is time-consuming and expensive to collect labeled data. Since the background patches in images  
 56 or noise in speech signals are used for binary classification in our method, we refer to such data  
 57 as the noise of an auxiliary binary classification problem denoted by auxiliary binary classification  
 58 (ABC)-noise dataset. An advantage of using ABC-noise data during training is that it can implicitly  
 59 add robustness to deep neural networks against the background or noisy data.

60 Since ABC-noise data can be collected in large quantities for free and using that data in our approach  
 61 improves the classification benchmarks, we investigate whether this data can act as a substitute for  
 62 labeled data. We conduct empirical analysis and show that using only a fraction of labeled training  
 63 data together with ABC-noise data in our DBT method, indeed improves the accuracy as compared  
 64 to normal training.

65 To summarize, our contributions are threefold. First, we present a detailed theoretical analysis on  
 66 the relation between the estimation of unknown parameters of pdf of data and classification outputs.  
 67 Second, based on the theoretical analysis, we present a simple booster training method to improve  
 68 classification accuracy which also doubles up as an augmented training method when only limited  
 69 labeled data is available. Third, we consistently achieve improved performances over strong baselines  
 70 on face recognition, image classification, and speaker recognition problems using our proposed  
 71 method, showing its generalization across different domains and model architectures.

## 72 2 RELATED WORK

73 **Notations and Preliminary:** In this paper, vectors, matrices, functions, and sets are denoted by bold  
 74 lower case, bold uppercase, lower case, and calligraphic characters, respectively. Consider a datapoint  
 75 denoted by  $\mathbf{x}$ . We assume that  $\mathbf{x}$  belongs to a family of probability density functions (pdf's) defined  
 76 as  $\mathcal{P} = \{p(\mathbf{x}, \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ , where  $\Theta$  is the possible set of parameters of the pdf. In general,  $\boldsymbol{\theta}$  is a real  
 77 vector in higher dimensions. For example, in a mixture of Gaussians,  $\boldsymbol{\theta}$  is a vector containing the  
 78 component weights, the component means, and the component covariance matrices. In this paper, we  
 79 assume that  $\boldsymbol{\theta}$  is an unknown deterministic function (There are other approaches such as bayesian  
 80 that consider  $\boldsymbol{\theta}$  as a random vector). In general, although the structure of the family of pdfs is itself  
 81 unknown, defining a family of pdfs such as  $\mathcal{P}$  can help us to develop theorems and use those results  
 82 to derive a new method. For the family of distribution  $\mathcal{P}$ , we can define the following classification  
 83 problem

$$\{ \mathcal{C}_1 : \boldsymbol{\theta} \in \Theta_1, \mathcal{C}_2 : \boldsymbol{\theta} \in \Theta_2, \dots, \mathcal{C}_n : \boldsymbol{\theta} \in \Theta_n \} \quad (1)$$

84 where set of  $\Theta_i$ 's is a partition of  $\Theta$ . The notation of (1) means that, class  $\mathcal{C}_i$  deals with a set of  
 85 data points whose pdf is  $p(\mathbf{x}, \boldsymbol{\theta}_i)$  where  $\boldsymbol{\theta}_i \in \Theta_i$ . A wide range of classification problems can be  
 86 defined using (1) e.g., ((Lehmann & Casella, 2006, Chapter 3)) and ((Duda et al., 2012, Chapter 4)).  
 87 The problem of estimating  $\boldsymbol{\theta}$  comes under the category of parametric estimation or point estimation  
 88 (Lehmann & Casella (1998)). Estimating the unknown parameters of a given pdf  $p(\mathbf{x}, \boldsymbol{\theta})$ , have been  
 89 extensively studied in the field of point estimation methods (Lindgren (2017); Lee et al. (2018);  
 90 Lehmann & Casella (2006)). An important estimator in this field is the minimum variance unbiased  
 91 estimator and it is governed by the Cramer Rao bound. The Cramer Rao bound provides the lower  
 92 bound of the variance of an unbiased estimator (Bobrovsky et al. (1987)). Let the estimation of  
 93  $\boldsymbol{\theta}$  be denoted by  $\hat{\boldsymbol{\theta}}$ , and assume that  $\hat{\boldsymbol{\theta}}$  is an unbiased estimator, i.e.,  $E(\hat{\boldsymbol{\theta}}) = \boldsymbol{\theta}$ . Its covariance  
 94 matrix denoted by  $\Sigma_{\hat{\boldsymbol{\theta}}}$  satisfies  $\Sigma_{\hat{\boldsymbol{\theta}}} - \mathbf{I}^{-1}(\boldsymbol{\theta}) \succeq \mathbf{0}$ , where  $\mathbf{A} \succeq \mathbf{0}$  implies that  $\mathbf{A}$  is a non-negative  
 95 definite matrix ((Lehmann & Casella, 1998, chapter 5)) and  $\mathbf{I}(\boldsymbol{\theta}) := -E(\partial^2 \log(p(\mathbf{x}, \boldsymbol{\theta}))/\partial \boldsymbol{\theta}^2)$   
 96 is called the Fisher information matrix. For an arbitrary differentiable function  $g(\cdot)$ , an efficient  
 97 estimator of  $\mathbf{g}(\boldsymbol{\theta})$  is an unbiased estimator when its covariance matrix equals to  $\mathbf{I}_{\mathbf{g}}^{-1}(\boldsymbol{\theta})$ , where  $\mathbf{I}_{\mathbf{g}}^{-1}(\boldsymbol{\theta})$   
 98 is the fisher information matrix of  $\mathbf{g}(\boldsymbol{\theta})$ , i.e., the efficient estimator achieves the lowest possible

99 variance among all unbiased estimators. The efficient estimator can be achieved using factorization of  
 100  $\partial \log(p(\mathbf{x}, \boldsymbol{\theta})) / \partial \mathbf{g}(\boldsymbol{\theta}) = I_{\mathbf{g}}(\boldsymbol{\theta})(\widehat{\mathbf{g}}(\mathbf{x}) - \mathbf{g}(\boldsymbol{\theta}))$ , if it exists (Rao (1992); Lehmann & Casella (1998)).  
 101 Based on these results, we derive a relationship between the efficient estimation of unknowns and  
 102 maximum likelihood classifier of (1) and use auxiliary binary classifiers to apply that result in our  
 103 proposed DBT method.

104 **Parameter Estimations:** Independent component analysis (Hyvärinen (1999)) decomposes a multi-  
 105 variate signal into independent non-Gaussian signals. ICA can extract non-Gaussian features from  
 106 Gaussian noise. Additionally, there is a class of classifiers called generalized likelihood ratio functions  
 107 that replaces the estimation of unknown parameters into the likelihood functions. This approach  
 108 provides a huge improvement in the field of parametric classifiers, where the family of pdf of data  
 109 is given (Zeitouni et al. (1992), Conte et al. (2001), Lehmann & Casella (2006)). Noise-contrastive  
 110 estimation (NCE) (Gutmann & Hyvärinen (2010)) involves training a generative model that allows  
 111 a model to discriminate data from a fixed noise distribution. Then, this trained model can be used  
 112 for training a sequence of models of increasing quality. This can be seen as an informal competition  
 113 mechanism similar in spirit to the formal competition used in the adversarial networks game. In  
 114 Bachman et al. (2019), a feature selection is proposed by maximizing the mutual information of the  
 115 difference between features extracted from multiple views of a shared context. In that work, it is  
 116 shown that the best results is given by using a mutual information bound based on NCE. The key  
 117 difference between our method and NCE is that, we do not construct a generative model for noise.  
 118 Instead of estimating the pdf of noise in NCE, we estimate the parameters of pdf of in-domain dataset  
 119 using an auxiliary class that has many common parameters in its pdf. Moreover, we show that the  
 120 estimation of that parameters are sufficient statistic for a classifier. We assume that the noise dataset is  
 121 not pure and it has some similarity with the in-domain dataset, where it can help the feature selection  
 122 layers to select relevant (in-domain) features, e.g., see Fig. 3. Further, in our approach, we do not  
 123 construct the pdf of noise or in-domain data, instead we estimate its parameters directly, which is  
 124 more efficient in terms of training, computation and also dimensionality reduction.

125 Auxiliary classifiers were introduced in inception networks (Szegedy et al. (2015)) and used in (Lee  
 126 et al. (2015); S. et al. (2016)) for training very deep networks to prevent vanishing gradient problems.  
 127 Further, auxiliary classifiers were also proposed for early exit schemes (Teerapittayanon et al. (2016))  
 128 and self-distillation methods (Zhang et al. (2019a;b)). Such auxiliary classifiers tackle different  
 129 problems by predicting the same target as the final classification layer. In contrast, our proposed DBT  
 130 method involves auxiliary binary classifiers that detect noise, interference, and/or background data  
 131 from in-domain data points for improving the target classification accuracy.

### 132 3 ESTIMATION OF PARAMETERS OF PDF AND CLASSIFICATION

133 For (1), we define a deterministic discriminative function of  $\Theta_i$ , denoted by  $t_i(\cdot)$  such that the  
 134 following conditions are satisfied:

- 135 •  $t_i(\cdot)$  maps  $\Theta$  to real numbers such that  $t_i(\boldsymbol{\theta}) > 0$ , if  $\boldsymbol{\theta} \in \Theta_i$  and  $t_i(\boldsymbol{\theta}) \leq 0$  for  $\boldsymbol{\theta} \notin \Theta_i$ .
- 136 •  $t_i(\cdot)$  is a differentiable function almost everywhere and  $\int_{\Theta} |t_i(\boldsymbol{\theta})| d\mu_i(\boldsymbol{\theta}) < \infty$ , where  $\mu_i$  denotes  
 137 the Lebesgue measure.

138 The following theorem shows the relationship of  $t_i(\cdot)$  and the log-likelihood ratio of class  $\mathcal{C}_i$  versus  
 139 other classes. The proofs of Theorems 1, 2 and 3 are provided in the appendix.

140 **Theorem 1** *Assume that the pdf  $p(\mathbf{x}, \boldsymbol{\theta})$  is differentiable with respect to  $\boldsymbol{\theta}$  almost everywhere. If the  
 141 efficient minimum variance and unbiased estimation of a deterministic discriminative function of  $\Theta_i$   
 142 exists, then the log likelihood ratio of class  $i$  against the rest of classes is an increasing function of  
 143 the minimum variance and unbiased estimation of  $\Theta_i$ .*

144 Directly from this theorem, it follows that the optimal classifier using the maximum likelihood for (1)  
 145 is given as follows  $d(\mathbf{x}) = \arg \max_{i \in \{1, \dots, n\}} k_i(t_i(\mathbf{x}))$ , where  $k_i$ 's are some increasing functions and  
 146  $t_i(\cdot)$ 's are the deterministic discriminative function of  $\Theta_i$ 's such that the efficient minimum variance  
 147 and unbiased estimation for them exists. Based on this result, a set of minimum variance and unbiased  
 148 estimation of deterministic discriminative functions of  $\Theta_i$ 's leads us to the maximum likelihood  
 149 classifier. One approach is to directly estimate the deterministic discriminative functions, instead of  
 150 maximizing the likelihood function. However, finding deterministic discriminative functions that  
 151 have efficient minimum variance and unbiased estimation may not be feasible in practical problems,

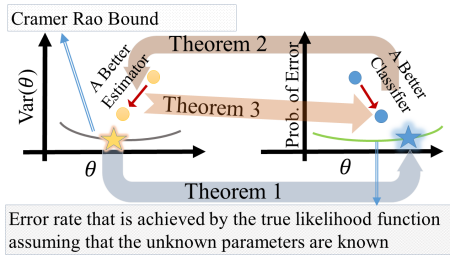


Figure 1: Visualizing Theorems 1,2 and 3

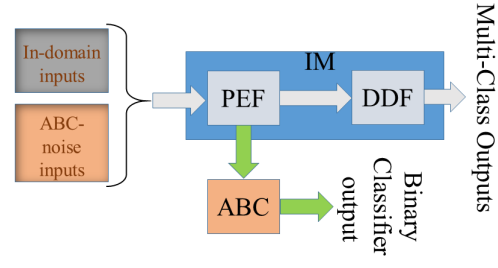


Figure 2: A general schema of our proposed DBT method with PEF, DDF and ABC blocks

152 especially when the dimension of  $\theta$  increases. Theorems 2 and 3 study the same relationship between  
 153 the estimation of unknown parameters and the accuracy of classifiers for sub-optimal estimators and  
 154 classifiers.

155 **Theorem 2** Consider the output of two classifiers for the  $i$ th class as follows:  $r_j(\mathbf{x}) = i$  if  $h_j(\mathbf{x}) > \tau$   
 156 and  $r_j(\mathbf{x}) = \text{other classes}$  if  $h_j(\mathbf{x}) < \tau$ , where  $j \in \{1, 2\}$ . where  $h_j(\mathbf{x})$  is the estimation of a  
 157 deterministic discriminative function and  $\tau$  is a classification threshold. Assume that the cumulative  
 158 distribution function of  $h_j(\mathbf{x})$ 's have bounded inflection points, and also, the probability of true  
 159 positive of  $r_j(\mathbf{x})$  is an increasing function of  $d(\theta)$ , which is the deterministic discriminative function  
 160 of class  $i$ , for all  $i$ . Further assume that for each  $\tau$  the probability of false positive of  $r_1(\mathbf{x})$  is less  
 161 than the probability of false positive of  $r_2(\mathbf{x})$  and the probability of true positive of  $r_1(\mathbf{x})$  is greater  
 162 than the probability of true positive of  $r_2(\mathbf{x})$ . Then, there exists a  $h_{\min}$  such that for all  $d(\theta) > h_{\min}$   
 163 and all  $\theta$  we have  $\Pr(|h_1(\mathbf{x}) - d(\theta)| < \epsilon) > \Pr(|h_2(\mathbf{x}) - d(\theta)| < \epsilon)$ .

164 Theorem 2 shows that a better classifier leads to a better estimation of  $d(\theta)$ . In the next theorem, we  
 165 show the dual property of this result.

166 **Theorem 3** Let  $\Theta_m$  be a Borel set with positive Lebesgue measure in (1) for all  $m \in \{1, \dots, n\}$ .  
 167 Assume that  $r_1(\cdot)$  and  $r_2(\cdot)$  are given as follows  $r_1(\mathbf{x}) = m$ , if  $\hat{\theta}_1 \in \Theta_m$  and  $r_2(\mathbf{x}) = m$ , if  $\hat{\theta}_2 \in \Theta_m$ .  
 168 Also, assume that  $\Pr(\|\hat{\theta}_1 - \theta\| \leq \epsilon) \geq \Pr(\|\hat{\theta}_2 - \theta\| \leq \epsilon)$ , for all  $\theta \in \Theta = \cup_{m=1}^n \Theta_m$  and  $\epsilon > 0$ ,  
 169 then the probability of classification error  $r_1(\cdot)$  is less than  $r_2(\cdot)$  where  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two different  
 170 estimators of  $\theta \in \Theta = \cup_{m=0}^{M-1} \Theta_m$ .

171 Theorem 3 proves that a more accurate estimator leads to a classifier that has a lower probability  
 172 of classification error. From Theorem 1, we can infer that a sufficient statistic for developing the  
 173 maximum likelihood classification is  $\hat{t}_i(\mathbf{x})$ , which is the efficient minimum variance and unbiased  
 174 estimation of the deterministic discriminative functions of  $\Theta_i$ 's denoted by  $t_i(\theta)$ . In other words, the  
 175 maximum likelihood classifier is a function of  $\mathbf{x}$  only via the efficient minimum variance and unbiased  
 176 estimation  $t_i(\theta)$ . We can estimate  $t_i(\theta)$  by replacing the estimation  $\theta$  in  $t_i(\cdot)$ , i.e.,  $\widehat{t}_i(\hat{\theta}) \approx t_i(\hat{\theta})$ ,  
 177 where  $\hat{\theta}$  is a function of  $\mathbf{x}$ . From the above theorems, we conclude that improving the estimation  
 178 of unknown parameters of pdf of data can improve the accuracy of the classifier. On the other side,  
 179 having a good classifier means having a good estimator of unknowns of the pdf of input data. In  
 180 many practical problems, the optimal maximum likelihood classifier may not be achievable, but the  
 181 likelihood function of the classifier provides an optimal bound of the probability of error. In such  
 182 cases, we can improve the accuracy of sub-optimal classifiers and that is the main focus of this paper.  
 183 Fig. 1 illustrates the proposed theorems visually.

#### 184 4 PROPOSED METHOD: DETECTION BOOSTER TRAINING (DBT)

185 In this section, we propose the *detection booster training (DBT)* method based on the achieved  
 186 theorems in the previous section to improve the accuracy of deep networks. Specifically, we divide  
 187 a deep model into two parts - early and later layers. We apply a detector (detection here means  
 188 detecting a target pattern from noise/background) on the early layers of the neural network in order

Loss	Ver. Acc. (%)
ResNet-50-DBT (CE)	98.96
ResNet-50-DBT	<b>99.12</b>

Table 1: Verification accuracy on LFW dataset for two different  $\mathcal{L}_{ABC}$  trained using CASIA Yi et al. (2014) dataset.

Loss	Acc.	Acc. on H-set
ResNet-100-AF	78.85	00.04
ResNet-100-DBT	<b>81.11</b>	<b>21.00</b>

Table 2: Comparison of Rank-1 identification accuracy on the IJB-B, with animal distractors.

189 to improve the estimation of unknown parameters of the family of pdf (based on Theorem 2). A  
 190 better estimation of unknown parameters corresponds to better feature representations in the early  
 191 layers and these features are input to the rest of the layers to construct the deterministic discriminative  
 192 functions (DDF) useful for the in-domain data classification (based on Theorem 3).

193 A general schema for dividing a deep model into two sub-models namely PEF (parameter estimator  
 194 functions) and DDF is depicted in Figure 2. The early layers of the model estimate the unknown  
 195 parameters of pdf of data while the later layers construct the discriminative functions essential for  
 196 classification. Based on this scheme, we formally define the three main components of DBT as  
 197 follows:

- 198 • *parameter estimator functions* (PEF): The sub-network from input layer to the  $k$ th layer, where  $k$  is  
 199 a hyperparameter in the DBT approach.
- 200 • *auxiliary binary classification* (ABC): Some additional layers are attached to the end of PEF,  
 201 mapping the output of the  $k$ th layer to a one-dimensional vector.
- 202 • *deterministic discriminative functions* (DDF): The sub-network from  $k$ th layer to the output of the  
 203 model. The output of model is a vector equal to the length of the number of classes  $n$ .

204 From Theorem 2, we showed that unknown parameter estimation can be improved using a detection  
 205 approach. During training, we apply a binary classification on the early layers (PEF) of the model to  
 206 improve the estimation of unknown parameters of pdf and subsequently provide rich feature vectors  
 207 for DDF. We define the *auxiliary binary classification problem* (ABC problem) as follows:

- 208 • Class 1 (alternative hypothesis) of ABC problem denoted by  $\mathcal{H}_1$  is set of all data points of classes  
 209 of  $\mathcal{C}_1$  to  $\mathcal{C}_n$ , i.e.  $\theta \in \cup_{i=1}^n \Theta_i$ .
- 210 • Class 0 (null hypothesis) of ABC problem denoted by  $\mathcal{H}_0$  is a dataset of data points from same  
 211 distribution  $p(\mathbf{x}, \theta)$  but  $\theta \notin \cup_{i=1}^n \Theta_i$ . We also define the dataset of Class 0 of ABC as *ABC-noise*  
 212 dataset, i.e., the ABC is given by the following hypothesis testing problem:  $\mathcal{H}_1 : \theta \in \cup_{i=1}^n \Theta_i$  versus  
 213  $\mathcal{H}_0 : \theta \notin \cup_{i=1}^n \Theta_i$ . In many practical problems, the noise, background or interference data related to  
 214 the in-domain dataset have same type of probability distribution but different pdf parameters. Hence,  
 215 using that dataset is a cheap and adept choice for the null hypothesis of ABC.

216 The Auxiliary Binary Classification problem influences only the PEF and ABC units while the main  
 217 classification problem with  $n$  classes updates the parameters of both PEF and DDF using in-domain  
 218 data. Since the auxiliary classifier is only used during training, the *inference model* (IM) consists of  
 219 only PEF and DDF and hence, there is no additional computation cost during inference. We formulate  
 220 the aforementioned method in the following notations and loss functions. Assume that  $\mathbf{x}$  is a data  
 221 point that belongs to Class  $\mathcal{C}_i$ ,  $i \in \{1, \dots, n\}$  or Class  $\mathcal{H}_0$  of ABC. Here, we define two type of  
 222 labels denoted by  $l_{ABC}$  and  $l_{MC}$ , where the subscription "MC" stands for multi-classes. So, if  $\mathbf{x}$   
 223 belongs to class  $\mathcal{C}_i$ , then  $l_{ABC} = 1$  and  $l_{MC} = i - 1$ , else if  $\mathbf{x}$  is a ABC-noise data point,  $l_{ABC} = 0$   
 224 and  $l_{MC}$  is NONE. Therefore, the loss function is defined as:

$$\mathcal{L}_{tot} = \mathcal{L}_{ABC}(Q_{ABC}(Q_{PEF}(\mathbf{x})), l_{ABC}) + \lambda l_{ABC} \mathcal{L}_{MC}(Q_{DDF}(Q_{PEF}(\mathbf{x})), l_{MC}), \quad (2)$$

225 where  $Q_{PEF}$ ,  $Q_{ABC}$  and  $Q_{DDF}$  are the functions of PEF, ABC and DDF blocks, respectively. We  
 226 set the hyperparameter  $\lambda = 1$  to balance the two loss terms. It is seen that, the second term of the  
 227 total loss is zero if  $l_{ABC} = 0$ .  $\mathcal{L}_{ABC}$  and  $\mathcal{L}_{MC}$  are selected based on the problem definition and  
 228 datasets. For classification, a simple selection for them can be binary cross-entropy and cross-entropy,  
 229 respectively. For a given task and deep neural network, the choice of  $k$  and  $\mathcal{L}_{ABC}$  influences the  
 230 feature representation of early layers differently and consequently the accuracy of the model. We  
 231 provide empirical studies in the next section to verify the same.



Figure 3: Maximally activated receptive fields of layer 15 of Inception-ResNet-v1 with (top row) and without (bottom row) DBT.



Figure 4: Examples of mis-identified faces along with their corresponding animal distractors on the IJB-B for ArcFace.

232 5 EXPERIMENTAL STUDY OF DBT

233 FACE RECOGNITION

234 We conduct experiments on face recognition benchmarks and show that the DBT method learns rich  
 235 features essential for face recognition. We also discover an important observation that current state-  
 236 of-the-art (SOTA) face recognition models are very sensitive to non-face data, in particular, animal  
 237 faces. Fig. 4 shows a few examples of misidentified faces and their corresponding animal distractors  
 238 from the IJB-B dataset using the ArcFace (Deng et al. (2019)) model. We show that our DBT method  
 239 not only improves the verification accuracy but also implicitly tackles this robustness issue of current  
 240 models against non-face data. Implementation details are provided in the appendix.

241 We consider the PEF discussed in Section 4 to be the first three layers of the model and DDF to be  
 242 the rest of layers. Ablation studies on the choice of PEF and DDF are provided in the supplementary  
 243 material. We define  $\mathcal{L}_{MC}$  in (2) as the SOTA ArcFace loss function proposed in (Deng et al. (2019)).  
 244 The ABC-noise is a non-face dataset containing 500K images that we collected from background  
 245 patches of MS1MV2 (Guo et al. (2016)) (More details in Appendix). We experimented with two  
 246 different loss functions for  $\mathcal{L}_{ABC}$ . For the first one, since popular face recognition models (Deng et al.  
 247 (2019); Wang et al. (2018)) use normalized output features and compute the losses on a hypersphere,  
 248 we select  $\mathcal{L}_{ABC}$  as follows. Let  $p_f \in \mathbb{R}^d$  and  $p_{n_f} \in \mathbb{R}^d$  denote the prototypes for faces and non-  
 249 faces, respectively. Following (Mettes et al. (2019)), we constrain the face/non-face prototypes on  
 250 diametrically opposite directions i.e  $\cos(\theta_{p_f p_{n_f}}) = -1$  and normalize the output feature vectors for  
 251 faces and non-faces such that  $\|p_{f_i}\| = \|p_{n_{f_i}}\| = 1$ . We then define the  $\mathcal{L}_{ABC}$  as,

$$\mathcal{L}_{ABC} = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{e^{s(\cos(m_1 \theta_{y_i} + m_2) - m_3)}}{e^{s(\cos(m_1 \theta_{y_i} + m_2) - m_3)} + e^{s \cos \theta_2}} \right) + \frac{1}{N} \sum_{i=1}^N (-1 - |p_{f_i} \cdot p_{n_{f_i}}|)^2, \quad (3)$$

252 where  $\theta_{y_i}$  and  $\theta_2$  correspond to the angles between the weights and the features for face and non-face  
 253 labels, respectively;  $m_1, m_2, m_3$  are the angular margins;  $s$  denotes the radius of the hypersphere. For  
 254 the second choice, we use simple binary cross entropy for  $\mathcal{L}_{ABC}$ . Table 1 shows that the verification  
 255 accuracy on LFW (Huang et al. (2007)) using (3) is 0.16% higher than simple cross entropy loss. This  
 256 also shows that choosing a task-specific  $\mathcal{L}_{ABC}$  is essential in obtaining more accurate results. We use  
 257 Eqn.1 as the default for  $\mathcal{L}_{ABC}$  in all our face recognition experiments, unless otherwise stated.

258 Table 3 compares the verification accuracy of our method versus the current SOTA method ArcFace  
 259 on five different test sets, LFW, CPLFW (Zheng & Deng (2018)), CALFW (Zheng et al. (2017)),  
 260 CFP-FP (Sengupta et al. (2016)) and AgeDb-30 (Moschoglou et al. (2017)). For the LFW test set,  
 261 we follow the unrestricted with labeled outside data protocol to report the performance. We trained  
 262 ResNet-50 and ResNet-100 using ArcFace and DBT approaches on CASIA (small) and MS1MV2  
 263 (large) datasets, respectively. The results show that DBT method outperforms ArcFace on all datasets.  
 264 Table 7 shows the angle statistics of the trained ArcFace and DBT models on the LFW dataset. Min.  
 265 Inter and Inter refer to the mean of minimum angles and mean of all angles between the template  
 266 embedding features of different classes (mean of the embedding features of all images for each class),  
 267 respectively. Intra refers to the mean of angles between  $x_i$  and template embedding feature for each  
 268 class. From Table 7, we infer that DBT extracts better face features and hence reduces the intra-class  
 269 variations. Directly from Tables 3 and 7, we infer that first, DBT consistently improves the accuracy

Method	LFW	CALFW	CPLFW	CFPFP	AgeDb-30
ResNet-50-AF (ArcFace)	98.46	89.48	80.88	86.74	88.98
ResNet-50-DBT	<b>99.12</b>	<b>91.38</b>	<b>87.10</b>	<b>94.95</b>	<b>91.23</b>
ResNet-100-AF (ArcFace)	99.61	94.50	89.35	96.14	95.33
ResNet-100-DBT	<b>99.75</b>	<b>95.13</b>	<b>90.70</b>	<b>96.90</b>	<b>96.16</b>

Table 3: ArcFace vs. DBT-ArcFace: verification(%) accuracy on LFW, CALFW, CPLFW, CFP-FP and AgeDb-30 of models ResNet-100 and ResNet-50.

270 on all test sets. Second, learning better features in the early layers is crucial to obtain rich face feature  
 271 embeddings. Third, the achieved gain using DBT is more pronounced on models trained using a  
 272 smaller (CASIA) dataset (it has fewer identities and images). This shows that DBT can address the  
 273 issue of the lack of in-domain data using cheap ABC-noise data.

274 We also provide the results of training Inception-ResNet-V1 and ResNet-64 models using DBT on  
 275 MS1MV2 to show the generalization capacity of the DBT method. For the Inception-ResNet-V1 and  
 276 ResNet-64, the PEF is set to be the first six layers and the DDF is the rest of the model. We use large  
 277 margin cosine loss (LMCL) Wang et al. (2018) for  $\mathcal{L}_{MC}$  and Cross entropy (CE) for  $\mathcal{L}_{ABC}$ . Table 4  
 278 shows the verification accuracy on LFW for Inception-ResNet-V1 and ResNet-64 models trained  
 279 on MS1MV2 with and without DBT. The results show that DBT method is independent of model  
 280 depth or architectures or loss functions and thereby consistently improves the accuracy compared  
 281 to baseline results. Table 4 also compares the DBT method with state-of-the-art methods on LFW  
 282 and YTF datasets. DBT method notably improves the baselines that are comparable to ArcFace and  
 283 superior to all the other methods. We were not able to reproduce the results of the ArcFace paper  
 284 using our Tensorflow implementation and dataset. We believe that using the original implementation  
 285 and dataset from ArcFace will achieve superior results over the baselines on the benchmark datasets  
 286 as evident from the results of our implementation. Finally, we compare the result ArcFace and DBT  
 287 on IJB-B and IJB-C, in Table 5. It is seen that DBT provides a notable boost on both IJB-B and  
 288 IJB-C by a considerable margin. DBT improves the verification accuracy as high as **1.94** % on IJB-B  
 289 and **2.57** % on IJB-C dataset at  $10^{-4}$  false alarm rate (FAR). We plot the receptive fields of the top  
 290 ten maximally activated neurons of an intermediate layer of the face recognition model to visualize  
 291 the features learned using the DBT method. Fig. 3 shows that the receptive fields of layer 15 of  
 292 the inception-resnet-v1 model trained using DBT attends to the regions of eyes, nose and mouth as  
 293 compared to insignificant regions in the normal training method. This shows that DBT learns more  
 294 discriminative features essential to face recognition, corroborating our theoretical claims.

295 To show that current SOTA models are not robust to animal faces, we performed a 1:N identification  
 296 experiment with approximately 3000 animal distractors on the IJB-B (Whitelam et al. (2017)) dataset.  
 297 We trained the face recognition model with about 500K non-face data which contains 200 animal  
 298 faces. This is disjoint from the 3000 distractors used in the identification experiment. We collected the  
 299 animal faces from web images using MTCNN (Zhang et al. (2016a)) face detector which are the false  
 300 positives from the face detector. Table 2 shows the Rank-1 identification accuracy of ResNet-100  
 301 on IJB-B dataset, trained on MS1MV2 using the ArcFace loss (ResNet-100-AF) versus our DBT  
 302 approach (ResNet-100-DBT). The third column of Table 2 denotes the accuracy on a hard subset  
 303 of images (false positives from ArcFace model) on the IJB-B dataset denoted by H-set. Results  
 304 of Table 2 show that current face recognition models are unable to discriminate out-of-distribution  
 305 (non-face) images from face images. Our ResNet-100-DBT significantly (as high as **21**%) reduces the  
 306 misidentification rate as compared to the ArcFace model which shows that DBT method inherently  
 307 overcomes this issue while also improving face recognition accuracy.

### 308 IMAGE CLASSIFICATION

309 In this section, we evaluate ResNet-110 and ResNext-101 models trained with and without DBT on  
 310 image classification problem using CIFAR-10, CIFAR-100, and ImageNet. We also show the power  
 311 of DBT to compensate for the smaller in-domain training set. For all implementations, PEF is defined  
 312 to be the first three layers and DDF is the rest of the model.  $\mathcal{L}_{ABC}$  and  $\mathcal{L}_{MC}$  are set to cross-entropy  
 313 loss. ABC-noise is the same data used in face recognition experiments. We follow the same training  
 314 configurations from (He et al. (2016); Xie et al. (2017)).

315 To study the efficacy of the DBT method in augmenting smaller in-domain training datasets, we  
 316 also trained ResNet-100 and ResNext-101 using partial training data on CIFAR-10 and CIFAR-100.

Model	Loss	LFW	Method	LFW	YTF
Inception Resnet	CE	99.45	Center Loss	99.28	94.9
Inception Resnet-DBT	CE	<b>99.50</b>	Range Loss	99.52	93.7
Inception Resnet	LMCL	99.55	SphereFace	99.42	95.0
Inception Resnet-DBT	LMCL	<b>99.60</b>	SphereFace+	99.47	-
Resnet 64	CE	99.55	CosFace	99.73	97.6
Resnet 64-DBT	CE	<b>99.63</b>	ArcFace	<b>99.82</b>	<b>98.02</b>
Resnet 64	LMCL	99.65	ArcFace**	99.61	97.31
Resnet 64-DBT	LMCL	<b>99.68</b>	ResNet-100-DBT	<b>99.75</b>	<b>97.67</b>

Table 4: Comparison of DBT models with SOTA methods on LFW and YTF. ArcFace \*\* refers to our arcface implementation.

Method	IJB-B						IJB-C					
	10 <sup>-6</sup>	10 <sup>-5</sup>	10 <sup>-4</sup>	10 <sup>-3</sup>	10 <sup>-2</sup>	10 <sup>-1</sup>	10 <sup>-6</sup>	10 <sup>-5</sup>	10 <sup>-4</sup>	10 <sup>-3</sup>	10 <sup>-2</sup>	10 <sup>-1</sup>
ArcFace	38.47	65.60	82.97	91.11	96.01	98.91	61.96	73.22	83.84	91.85	96.51	<b>99.08</b>
DBT	<b>47.01</b>	<b>72.70</b>	<b>84.91</b>	<b>91.92</b>	<b>96.37</b>	<b>99.03</b>	<b>67.42</b>	<b>77.33</b>	<b>86.41</b>	<b>92.75</b>	<b>96.66</b>	99.06

Table 5: 1:1 verification: ResNet-100: DBT vs. ArcFace on the IJB-B and IJB-C datasets

Method	Top-1	Top-5	Method	Min. Inter	Intra	Inter
ResNet	22.10	6.15	ArcFace	53.23	7.2	<b>88.73</b>
ResNet-DBT	<b>21.82</b>	<b>6.02</b>	ResNet-DBT	<b>52.96</b>	<b>7.16</b>	88.52

Table 6: Top-1 and Top-5 error rates (%) on ILSVRC15 benchmark for ResNet w/o DBT.

Table 7: Comparison of inter and intra angles (degrees) for different methods on LFW.

ResNet Models	CIFAR-10	CIFAR-100	ResNext Models	CIFAR-10	CIFAR-100
He et al. (2016)*	5.84	22.15	Xie et al. (2017)*	5.03	21.24
DBT (5/5)	<b>5.25</b>	<b>21.53</b>	DBT (5/5)	<b>4.68</b>	<b>19.79</b>
ResNet (4/5)	5.89	24.23	ResNext (4/5)	4.93	23.52
DBT (4/5)	<b>5.36</b>	<b>23.98</b>	DBT (4/5)	<b>4.76</b>	<b>22.56</b>
ResNet (3/5)	6.61	27.99	ResNext (3/5)	5.38	27.25
DBT (3/5)	<b>5.44</b>	<b>26.81</b>	DBT (3/5)	<b>4.77</b>	<b>26.04</b>
ResNet (2/5)	7.06	33.81	ResNext (2/5)	5.85	33.62
DBT (2/5)	<b>5.94</b>	<b>31.95</b>	DBT (2/5)	<b>5.05</b>	<b>30.73</b>
ResNet (1/5)	8.20	47.43	ResNext (1/5)	7.24	48.05
DBT (1/5)	<b>6.86</b>	<b>43.65</b>	DBT (1/5)	<b>6.05</b>	<b>42.56</b>

Table 8: Comparison of Top-1 error rates (%) for CIFAR-10 and CIFAR-100 datasets w/o DBT.\* denotes our implementation. (x/5) denotes the fraction of training data used for training that model.

Method	VoxC (top 1)	VoxC (top 5)	Librispeech	VCTK	ELSDSR
VGG-M CNN	80.5	92.1	93.12	82.52	79.98
VGG-M CNN-DBT	<b>82.3</b>	<b>95.8</b>	<b>95.62</b>	<b>88.14</b>	<b>81.56</b>

Table 9: Accuracy of speaker identification (%) for different datasets.

Method	CIFAR-10	CIFAR-100
ResNet-Back	5.65	21.84
ResNet-DBT	<b>5.25</b>	<b>21.53</b>
ResNext-Back	4.97	21.65
ResNext-DBT	<b>4.68</b>	<b>19.79</b>

Table 10: Comparison of top-1 error rates on CIFAR-10 and CIFAR-100 using an additional background class vs DBT.



Method	LFW	CALFW	CPLFW	CFP-FP	AgeDb-30
ResNet+mod	99.16	91.46	86.11	93.81	92.71
ResNet-DBT+mod	<b>99.65</b>	<b>95.05</b>	<b>90.08</b>	<b>96.20</b>	<b>95.87</b>

Table 11: Ablation study on the verification performance of adding background class to the model on MS1MV2 dataset.

317 We randomly selected a fraction of the training data to be our training set, e.g.,  $k/5$  of dataset  
318 means that we only used  $k$  fifth of total samples for training. From first row of Table 8, we find that  
319 models trained with DBT show **0.59%** and **0.35%** improvement on CIFAR-10, **0.62%** and **1.45%**  
320 improvement on CIFAR-100 over baseline models for ResNet-110 and ResNext-101 architectures,  
321 respectively. Furthermore, using partial training data with our DBT method achieves superior results  
322 (as high as 5.49 % on ResNext (1/5) CIFAR-100) as compared to normal training. Table 6 shows  
323 the results on Imagenet. We see that DBT improves the accuracy by **0.28%** on Top-1 accuracy. This  
324 shows that the DBT method consistently improves the results on both small and large datasets.

### 325 SPEAKER IDENTIFICATION

326 We consider the problem of speaker identification using the VGG-M (Chatfield et al. (2014)) model.  
327 We set PEF as the first two CNN layers and DDF as the remaining CNN layers.  $\mathcal{L}_{ABC}$  and  $\mathcal{L}_{MC}$   
328 are defined to be the cross-entropy loss. The ABC-noise is generated from the silence intervals of  
329 VoxCeleb (Nagrani et al. (2017)) augmented with Gaussian noise with variance one. The input to the  
330 model is the short-time Fourier transformation of speech signals with a hamming sliding window  
331 of width 25 ms and step 10 ms. Table 9 provides the accuracies of VGG-M model trained with and  
332 without DBT on VoxCeleb, Librispeech (Panayotov et al. (2015)), VCTK (Veaux et al. (2016)) and  
333 ELSDR (L. (2004)) datasets. Table 9 shows that the trained models using DBT significantly improves  
334 the accuracy (as high as **5.62%**) for all datasets. Implementation details are provided in the appendix.  
335

### 336 MISCELLANEOUS EXPERIMENTS

337 In this section, we experiment with the naive way of using background data by considering non-faces  
338 as a separate class in the final classification layer. For face recognition, Table 11 shows the results  
339 of training with an additional background class on MS1MV2 dataset with and without using DBT.  
340 ResNet+mod refers to a model trained with ArcFace loss and  $n + 1$  classes where the additional class  
341 corresponds to the non-faces. ResNet-DBT+mod refers to a model trained with both DBT and the  
342 additional non-face class. We find that adding the additional non-face class hurts the performance  
343 of the model whereas ResNet-DBT+mod improves the results significantly relative to ResNet+mod  
344 model. Since the non-face dataset is sampled from a wide range of a family of distributions compared  
345 with faces, it has a larger range of unknown parameters, then the sufficient statistic of them should be  
346 larger than the sufficient statistics of face data. Thus, when we restrict faces and non-faces on the  
347 surface of a hypersphere, the non-face data is more spread on the surface compared with each of the  
348 other face classes. We demonstrate this effect with the help of a toy example in Fig. 6 in the appendix.  
349 We also conduct this experiment on CIFAR-10/CIFAR-100 and report it in Table 10. We see that  
350 naively incorporating the background class is inferior to DBT showing that DBT is an effective  
351 technique to utilize background data to boost the performance of classification models.

## 352 6 CONCLUSION

353 In this paper, we presented a detailed theoretical analysis of the dual relationship between estimating  
354 the unknown pdf parameters and classification accuracy. Based on the theoretical study, we presented  
355 a new method called DBT using ABC-noise data for improving in-distribution classification accuracy.  
356 We showed that using ABC-noise data helps in better estimation of unknown parameters of pdf of  
357 input data and thereby improves the feature representations and consequently the accuracy in image  
358 classification, speaker classification, and face recognition benchmarks. It also augments the training  
359 data when only limited labeled data is available by improving accuracy. We showed that the concept  
360 of DBT is generic and generalizes well across domains through extensive experiments using different  
361 model architectures and datasets. Our framework is complementary to existing training methods and  
362 hence, it can be easily integrated with current and possibly future classification methods to enhance  
363 accuracy. In summary, the proposed DBT method is a powerful technique that can augment limited  
364 training data and improve classification accuracy in deep neural networks.

## 365 REFERENCES

- 366 M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean,  
367 M. Devin, and S. Ghemawat and. TensorFlow: Large-scale machine learning on heterogeneous  
368 systems, 2015.
- 369 Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing  
370 mutual information across views. In *Advances in Neural Information Processing Systems*, pp.  
371 15535–15545, 2019.
- 372 T. L. Berg, A. C. Berg, J. Edwards, and D. A. Forsyth. Who’s in the picture. *NeurIPS*, 2004.
- 373 David Berthelot, Colin Raffel, Aurko Roy, and Ian Goodfellow. Understanding and improving  
374 interpolation in autoencoders via an adversarial regularizer. *ICLR*, 2019.
- 375 Ben-Zion Bobrovsky, E Mayer-Wolf, and M Zakai. Some classes of global cramér-rao bounds. *The*  
376 *Annals of Statistics*, pp. 1421–1438, 1987.
- 377 Ken Chatfield, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Return of the devil in the  
378 details: Delving deep into convolutional nets. *BMVC*, 2014.
- 379 Ernesto Conte, Antonio De Maio, and Giuseppe Ricci. Grlt-based adaptive detection algorithms for  
380 range-spread targets. *IEEE transactions on signal processing*, 49(7):1336–1348, 2001.
- 381 J. Deng, J. Guo, X. Niannan, and S. Zafeiriou. Arcface: Additive angular margin loss for deep face  
382 recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 383 Richard O Duda, Peter E Hart, and David G Stork. *Pattern classification*. John Wiley & Sons, 2012.
- 384 Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew  
385 Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of*  
386 *computer vision*, 111(1):98–136, 2015.
- 387 F. F. Li, R. Fergus, and P. Perona. Learning generative visual models from few training examples: An  
388 incremental bayesian approach tested on 101 object categories. In *CVPR Workshop*, pp. 178–178,  
389 2004.
- 390 Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by  
391 predicting image rotations. *ICLR*, 2018.
- 392 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
393 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, pp. 2672–2680, 2014.
- 394 Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. Ms-celeb-1m: A dataset and benchmark for large-scale  
395 face recognition. *ECCV*, 9907:87–102, 2016.
- 396 Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle  
397 for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on*  
398 *Artificial Intelligence and Statistics*, pp. 297–304, 2010.
- 399 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *Computer Vision*  
400 *and Pattern Recognition*, pp. 770–778, 2016.
- 401 G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for  
402 studying face recognition in unconstrained environments. In *Technical Report*, 2007.
- 403 Aapo Hyvärinen. Survey on independent component analysis. 1999.
- 404 S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing  
405 internal covariate shift. *International Conference on Machine Learning*, 37:448–456, 2015.
- 406 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *ICLR*, 2014.
- 407 Feng L. Speaker recognition, informatics and mathematical modelling. *Technical University of*  
408 *Denmark, DTU*, 2004.

- 409 C-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu. Deeply-Supervised Nets. *Proceedings of*  
410 *Machine Learning Research (PMLR)*, 38:562–570, 2015.
- 411 Youngjo Lee, John A Nelder, and Yudi Pawitan. *Generalized linear models with random effects:*  
412 *unified analysis via H-likelihood*, volume 153. CRC Press, 2018.
- 413 E. L. Lehmann and G. Casella. *Theory of point estimation*, 1998. 2ndn ed.
- 414 Erich L Lehmann and George Casella. *Theory of point estimation*. Springer Science & Business  
415 Media, 2006.
- 416 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
417 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*  
418 *conference on computer vision*, pp. 740–755. Springer, 2014.
- 419 Bernard Lindgren. *Statistical theory*. Routledge, 2017.
- 420 B. Maze, J. Adams, J. A. Duncan, N. Kalka, T. Miller, C. Otto, A. K. Jain, W. T. Niggel, J. Anderson,  
421 J. Cheney, and P. Grother. Iarpa janus benchmark - c: Face dataset and protocol. *International*  
422 *Conference on Biometrics*, pp. 158–165, 2018.
- 423 P. Mettes, E. van der Pol, and C. Snoek. Hyperspherical prototype networks. *NeuRIPS*, 01 2019.
- 424 S. Moschoglou, A. Papaioannou, C. Sagonas, J. Deng, I. Kotsia, and S. Zafeiriou. Agedb: the first  
425 manually collected, in-the-wild age database. *CVPR Workshop*, 2(3):5, 2017.
- 426 Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. Voxceleb: a large-scale speaker identifica-  
427 tion dataset. *arXiv preprint arXiv:1706.08612*, 2017.
- 428 M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles.  
429 *ECCV*, 2016.
- 430 M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representa-  
431 tions using convolutional neural networks. *CVPR*, pp. 1717–1724, 2014.
- 432 V. Panayotov, G. Chen, D. Povey, and S. Khudanpur. Librispeech: an asr corpus based on public  
433 domain audio books. *International Conference on Acoustics, Speech and Signal Processing*  
434 *(ICASSP)*, pp. 5206–5210, 2015.
- 435 BLS Prakasa Rao. Cramer-rao type integral inequalities for estimators of functions of multidimen-  
436 sional parameter. *Sankhyā: The Indian Journal of Statistics, Series A*, pp. 53–73, 1992.
- 437 Ali Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf:  
438 an astounding baseline for recognition. *CVPR Workshops*, 2014.
- 439 O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla,  
440 M. Bernstein, A. C. Berg, and F.F Li. ImageNet Large Scale Visual Recognition Challenge.  
441 *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- 442 Christian S., Vincent V., Sergey I., Jon S., and ZB W. Rethinking the inception architecture for  
443 computer vision. *CVPR*, 2016.
- 444 S. Sengupta, J. Chen, C. Castillo, V. M. Patel, R. Chellappa, and D. W. Jacobs. Frontal to profile face  
445 verification in the wild. In *Winter Conference on Applications of Computer Vision (WACV)*, pp.  
446 1–9, 2016.
- 447 N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way  
448 to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15(1):1929–1958, 2014.
- 449 C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and  
450 A. Rabinovich. Going deeper with convolutions. *CVPR*, pp. 1–9, 2015.
- 451 S. Teerapittayanon, B. McDanel, and H. T. Kung. Branchynet: Fast inference via early exiting from  
452 deep neural networks. *ICPR*, 2016.

- 453 Christophe Veaux, Junichi Yamagishi, Kirsten MacDonald, et al. Superseded-cstr vctk corpus:  
454 English multi-speaker corpus for cstr voice cloning toolkit. *University of Edinburgh. The Centre  
455 for Speech Technology Research (CSTR)*, 2016.
- 456 H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu. Cosface: Large margin  
457 cosine loss for deep face recognition. *CVPR*, pp. 5265–5274, 2018.
- 458 C. Whitelam, E. Taborsky, A. Blanton, B. Maze, J. Adams, T. Miller, N. Kalka, A. K. Jain, J. A.  
459 Duncan, K. Allen, J. Cheney, and P. Grother. Iarpa janus benchmark-b face dataset. *CVPR  
460 Workshops*, pp. 592–600, 2017.
- 461 L. Wolf, T. Hassner, and I. Maoz. Face recognition in unconstrained videos with matched background  
462 similarity. *CVPR*, pp. 529–534, 2011.
- 463 S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural  
464 networks. *CVPR*, pp. 5987–5995, 2017.
- 465 D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv*, abs/1411.7923,  
466 2014.
- 467 J. Yosinski, J. Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural  
468 networks? *NIPS*, 2014.
- 469 Ofer Zeitouni, Jacob Ziv, and Neri Merhav. When is the generalized likelihood ratio test optimal?  
470 *IEEE Transactions on Information Theory*, 38(5):1597–1602, 1992.
- 471 K. Zhang, Z. Zhang, Z. Li, and Y. Qiao. Joint face detection and alignment using multi-task cascaded  
472 convolutional networks. *Signal Processing Letters*, 23(10):1499–1503, 2016a.
- 473 L. Zhang, J. Song, A. Gao, J. Chen, C. Bao, and K. Ma. Be your own teacher: Improve the  
474 performance of convolutional neural networks via self distillation. *ICCV*, 2019a.
- 475 Linfeng Zhang, Zhanhong Tan, Jiebo Song, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Scan:  
476 A scalable neural networks framework towards compact and efficient models. *NeurIPS*, 2019b.
- 477 Richard Zhang, Phillip Isola, and Alexei Efros. Colorful image colorization. *ECCV*, 2016b.
- 478 T. Zheng and W. Deng. Cross-pose lfw: A database for studying cross-pose face recognition in un-  
479 constrained environments. *Technical Report, Beijing University of Posts and Telecommunications*,  
480 2018.
- 481 T. Zheng, W. Deng, and J. Hu. Cross-age lfw: A database for studying cross-age face recognition in  
482 unconstrained environments. *arXiv*, abs/1708.08197, 2017.

## 483 APPENDIX

### 484 IN-DOMAIN FAMILY OF PDFS AND THE EXTENDED FAMILY OF DISTRIBUTIONS

485 In this section, we discuss about background/noise and in-domain data points and their corresponding  
486 distributions to clarify the definition of those concepts in this paper. Consider a random vector denoted  
487 by  $\mathbf{s}$ . Assume that the corresponding distribution is Gaussian with mean and variance given by  $\alpha \neq 0$   
488 and  $\sigma = 1$ , respectively. Now, assume that we observed  $\mathbf{x} = \mathbf{s} + \mathbf{n}$ , where the pdf of  $\mathbf{n}$  is assumed to  
489 be Gaussian with zero mean and variance  $\sigma_n^2$ , hence the pdf of  $\mathbf{x}$  is Gaussian with mean  $\alpha$  and variance  
490  $1 + \sigma_n^2$ . Here,  $\mathbf{n}$  is the background or noise data and the vector of unknowns is given by,  $\boldsymbol{\theta} = [\alpha, \sigma_n^2]$ .  
491 The *in-domain* family of pdfs for  $\mathbf{x}$  is then given by  $\mathcal{P}_{\mathbf{x}} = \{\mathcal{N}(\alpha, 1 + \sigma_n^2) | \alpha \neq 0, \sigma_n^2 > 0\}$ . If we  
492 include the family of pdf of  $\mathbf{n}$  to  $\mathcal{P}_{\mathbf{x}}$ , then we can extend  $\mathcal{P}_{\mathbf{x}}$  as  $\mathcal{P} = \{\mathcal{N}(\alpha, 1 + \sigma_n^2) | \alpha \in \mathbb{R}, \sigma_n^2 > 0\}$ .  
493 So  $\mathcal{P}$  is the union of family of pdfs of in-domain data points and noise/background data. From  
494 estimation theory, we know that the sufficient statistics and the unknown parameters of  $\mathcal{P}$  can also  
495 represent the sufficient statistics and the unknown parameters of  $\mathcal{P}_{\mathbf{x}}$ . In other words, an estimation of  
496  $\alpha$  can help us detect if the observed data point is from  $\mathbf{s} + \mathbf{n}$  or  $\mathbf{n}$  by comparing it with a threshold.  
497 Thus, estimating the unknown parameters of the family of pdfs using  $\mathcal{P}$  can provide more information  
498 about the observed data useful for tasks such as classification.

In general, we can assume that a generalized family of pdfs is given by the family of pdf of noise or background along with the family of pdfs of in-domain data. Hence, estimating from the extended

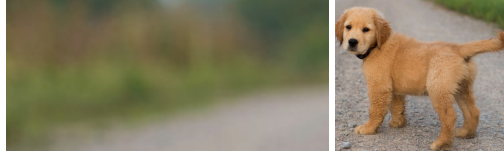


Figure 5: In-domain data point versus background data point. The background is cropped from the in-domain image and provides complementary information to the main data, thereby we can provide a better estimation of the pdf parameters of in-domain data.

family of distribution can provide more information about the in-domain distribution. Let us consider that the pdf of in-domain data points is given by  $p_{\mathbf{x}}(\mathbf{x}, [\boldsymbol{\theta}_s, \boldsymbol{\theta}_n])$  and the pdf of noise/background is given by  $p_{\mathbf{n}}(\mathbf{x}, \boldsymbol{\theta}_n)$ , so the extended pdf can be represented by

$$h(p_{\mathbf{n}}(\mathbf{x}, \boldsymbol{\theta}_n), p_{\mathbf{x}}(\mathbf{x}, [\boldsymbol{\theta}_s, \boldsymbol{\theta}_n])),$$

where  $h$  is a function that combines two pdfs in a general structure. So a general family of distribution can be denoted as follows:

$$\mathcal{P} = \{h(p_{\mathbf{n}}(\mathbf{x}, \boldsymbol{\theta}_n), p_{\mathbf{x}}(\mathbf{x}, [\boldsymbol{\theta}_s, \boldsymbol{\theta}_n])) | \boldsymbol{\theta} := [\boldsymbol{\theta}_s, \boldsymbol{\theta}_n] \in \Theta_{s,n}\},$$

499 where  $\boldsymbol{\theta}$  is defined as a new set of parameters in a higher dimension and  $\Theta_{s,n}$  are set of all possible  
 500  $[\boldsymbol{\theta}_s, \boldsymbol{\theta}_n]$  that belongs to  $p_{\mathbf{n}}$  and  $p_{\mathbf{x}}$ . The extended family of pdf provides more information about  
 501 the nuisance parameters of pdf of in-domain datapoints. Inspired by this observation, we develop  
 502 our detection booster training method using background/noise data. Figure 5 shows an example of  
 503 background and in-domain data point.

#### 504 PROOF OF THEOREM 1

505 Let  $t_i(\cdot)$  denote deterministic discriminative function of  $\Theta_i$ . Since the efficient minimum variance  
 506 and unbiased estimation of  $t_i(\boldsymbol{\theta})$  exists, we have

$$\frac{\partial \ln(p(\mathbf{x}, \boldsymbol{\theta}))}{\partial t_i(\boldsymbol{\theta})} = I_{t_i}(\boldsymbol{\theta})(\hat{t}_i(\mathbf{x}) - t_i(\boldsymbol{\theta})), \quad (4)$$

where  $\hat{t}_i(\mathbf{x})$  is the minimum variance and unbiased estimation of  $t_i(\boldsymbol{\theta})$  using the data point  $\mathbf{x}$  and  $I_{t_i}(\mathbf{x})$  is the Fisher information function of  $t_i(\boldsymbol{\theta})$ , which is given by

$$I_{t_i}(\boldsymbol{\theta}) = \frac{\partial t_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}^T \mathbf{I}(\boldsymbol{\theta}) \frac{\partial t_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \geq 0,$$

507 where  $T$  denotes the transpose and  $\mathbf{I}(\boldsymbol{\theta})$  is the Fisher information matrix of  $\boldsymbol{\theta}$ . Now we show that  
 508 the log-likelihood ratio is an increasing function in  $\hat{t}_i(\mathbf{x})$ . Note that  $I_{t_i}(\boldsymbol{\theta}) \geq 0$  (Lehmann & Casella  
 509 (2006)).

510 On the other hand, we have  $d \ln(p(\mathbf{x}, \boldsymbol{\theta})) = \sum_j \frac{\partial \ln(p(\mathbf{x}, \boldsymbol{\theta}))}{\partial \theta_j} d\theta_j$ , therefore,

$$\begin{aligned} \ln(p(\mathbf{x}, \boldsymbol{\theta})) + k(\mathbf{x}) &= \sum_j \int \frac{\partial \ln(p(\mathbf{x}, \boldsymbol{\theta}))}{\partial \theta_j} d\theta_j = \sum_j \int \frac{\partial \ln(p(\mathbf{x}, \boldsymbol{\theta}))}{\partial t_i(\boldsymbol{\theta})} \frac{\partial t_i(\boldsymbol{\theta})}{\partial \theta_j} d\theta_j = \\ &= \int \frac{\partial \ln(p(\mathbf{x}, \boldsymbol{\theta}))}{\partial t_i(\boldsymbol{\theta})} \sum_j \frac{\partial t_i(\boldsymbol{\theta})}{\partial \theta_j} d\theta_j = \int (I_{t_i}(\boldsymbol{\theta})(\hat{t}_i(\mathbf{x}) - t_i(\boldsymbol{\theta}))) \sum_j \frac{\partial t_i(\boldsymbol{\theta})}{\partial \theta_j} d\theta_j = \alpha(\boldsymbol{\theta})\hat{t}_i(\mathbf{x}) - \beta(\boldsymbol{\theta}) \end{aligned} \quad (5)$$

where the third equality is archived based on the third property of  $t_i(\cdot)$  in its definition and the forth equality is given by replacing (4);  $k(\mathbf{x})$  is the constant of integration. Finally, the last equality is given by defining the following terms

$$\alpha(\boldsymbol{\theta}) := \int I_{t_i}(\boldsymbol{\theta}) \sum_j \frac{\partial t_i(\boldsymbol{\theta})}{\partial \theta_j} d\theta_j, \quad \beta(\boldsymbol{\theta}) := \int I_{t_i}(\boldsymbol{\theta}) t_i(\boldsymbol{\theta}) \sum_j \frac{\partial t_i(\boldsymbol{\theta})}{\partial \theta_j} d\theta_j, \quad (6)$$

511 thus  $\frac{d\alpha(\boldsymbol{\theta})}{dt_i(\boldsymbol{\theta})} = I_{t_i}(\boldsymbol{\theta}) \geq 0$ , i.e.,  $\alpha(\boldsymbol{\theta})$  is increasing in  $t_i(\boldsymbol{\theta})$ . Since,  $t_i$  is a deterministic discriminative  
 512 function of  $\Theta_i$ , so for each  $j \neq i$  and  $\boldsymbol{\theta}_i \in \Theta_i$  and  $\boldsymbol{\theta}_j \in \Theta_j$ , we have  $t_i(\boldsymbol{\theta}_i) > t_i(\boldsymbol{\theta}_j)$ , therefore

513  $\alpha(\theta_i) \geq \alpha(\theta_j)$ . The later inequality is achieved based on the increasing property of  $\alpha(\theta)$  with  
514 respect to  $t_i(\theta)$ .

515 Using (5), the log likelihood ratio of class  $i$  against the rest of classes is given by  $\text{LLR} :=$   
516  $\ln(p(\mathbf{x}, \theta_i)) - \ln(p(\mathbf{x}, \theta_j))$ , so we have  $\text{LLR} = (\alpha(\theta_i) - \alpha(\theta_j))\hat{t}_i(\mathbf{x}) - (\beta(\theta_i) - \beta(\theta_j))$ .  $\text{LLR}$   
517 depends on  $\mathbf{x}$  only via  $\hat{t}_i(\mathbf{x})$  and since for each  $j \neq i$  and  $\theta_i \in \Theta_i$  and  $\theta_j \notin \Theta_i$ ,  $\alpha(\theta_i) - \alpha(\theta_j) > 0$ ,  
518 then  $\text{LLR}$  is increasing in  $\hat{t}_i(\mathbf{x})$ .  $\square$

519 **PROOF OF THEOREM 2**

The probability of true positive of class  $i$  of  $r_j$  is given by

$$P_{tp,i,j} = \Pr_{\theta}(h_j(\mathbf{x}) > \tau) = 1 - F_{j\theta}(\tau),$$

where  $F_{i\theta}(\cdot)$  denotes the Cumulative distribution function (CDF) of  $h_j$ . Since the probability of true positive of class  $i$  of  $r_1$  is greater than  $r_2$  for all  $\tau$ ,  $F_{1\theta}(\tau) < F_{2\theta}(\tau)$ , for all  $\tau$ . Now we define a function as follows

$$u(\tau, \theta) := F_{2\theta}(\tau) - F_{1\theta}(\tau).$$

Since the CDFs are increasing in  $\tau$  and tend to 1 and the number of inflection points of these CDFs are bounded, there is an  $h_{\min}$  such that, for  $\tau > h_{\min}$ , such that  $u(\tau, \theta)$  is a monotonically decreasing function in  $\tau$ . Thus for any  $\theta$  that satisfies  $d(\theta) > h_{\min}$  we have

$$u(d(\theta) + \epsilon, \theta) < u(d(\theta) - \epsilon, \theta).$$

520 Replacing  $u(h, \theta) = F_{2\theta}(h) - F_{1\theta}(h)$  in the last inequality, we have

$$F_{2\theta}(d(\theta) + \epsilon) - F_{1\theta}(d(\theta) + \epsilon) < F_{2\theta}(d(\theta) - \epsilon) - F_{1\theta}(d(\theta) - \epsilon) \Rightarrow \quad (7)$$

$$F_{2\theta}(d(\theta) + \epsilon) - F_{2\theta}(d(\theta) - \epsilon) < F_{1\theta}(d(\theta) + \epsilon) - F_{1\theta}(d(\theta) - \epsilon). \quad (8)$$

521 Based on the definition of CDF, we have

$$\begin{aligned} \Pr_{\theta}(|h_2(\mathbf{x}) - d(\theta)| < \epsilon) &= \Pr_{\theta}(d(\theta) - \epsilon < h_2(\mathbf{x}) < d(\theta) + \epsilon) < \\ \Pr_{\theta}(d(\theta) - \epsilon < h_1(\mathbf{x}) < d(\theta) + \epsilon) &= \Pr_{\theta}(|h_1(\mathbf{x}) - d(\theta)| < \epsilon). \end{aligned} \quad (9)$$

522  $\square$

523 **PROOF OF THEOREM 3**

524 First, we prove the following claim,

525 Claim: For any open set, there exists a set of disjoint countable open balls such that their union equals  
526 the origin open set.

527 Proof of claim: Consider an open set  $\mathcal{O}$ , and also consider  $x_0 \in \mathcal{O}$ , such that  $B(x_0, r_0) \subseteq \mathcal{O}$   
528 and  $r_0$  is the greatest possible radius between all possible open balls in  $\mathcal{O}$ , where  $B(x_0, r_0)$  is the  
529 open ball with radius  $r_0$  at point  $x_0$ . Now, we define  $x_1 \in \mathcal{O} - \overline{B(x_0, r_0)}$ , where  $\overline{B(x_0, r_0)}$  is  
530 the closure of  $B(x_0, r_0)$ , as the point with greatest radius in  $\mathcal{O} - \overline{B(x_0, r_0)}$  and similarly  $x_i \in$   
531  $\mathcal{O} - \bigcup_{k=0}^{i-1} \overline{B(x_k, r_k)}$  such that  $B(x_i, r_i)$  provides the greatest radius in  $\mathcal{O} - \bigcup_{k=0}^{i-1} \overline{B(x_k, r_k)}$ . So  
532 we have  $\mathcal{O} = \bigcup_{k=0}^{\infty} B(x_k, r_k)$ . This is because, if the latest equality is not valid, then there exists  
533 an open ball in  $\mathcal{O} - \bigcup_{k=0}^{\infty} \overline{B(x_k, r_k)}$  hence another open ball with greatest radius will be added to  
534  $\bigcup_{k=0}^{\infty} B(x_k, r_k)$ , which has a contradiction with the definition of  $\bigcup_{k=0}^{\infty} B(x_k, r_k)$ . The claim is proven  
535 at this point.

Now, we show the true positive probability of  $r_1$  is greater than  $r_2$ . Let  $\Theta'_m$  be the set of interior points of  $\Theta_m$ , then, there exists a union of disjoint open balls such that  $\Theta'_m = \bigcup_{k=0}^{\infty} B(x_k, r_k)$ . From assumptions in the theorem, we have  $\Pr(\|\hat{\theta}_1 - \theta\| \leq \epsilon) \geq \Pr(\|\hat{\theta}_2 - \theta\| \leq \epsilon)$ , then

$$\Pr_{\theta}(\hat{\theta}_1 \in B(x_k, r_k)) \geq \Pr_{\theta}(\hat{\theta}_2 \in B(x_k, r_k)),$$

536 where  $\theta \in \Theta_m$ . Based on the claim we have

$$\Pr_{\theta}(\hat{\theta}_1 \in \Theta'_m) \geq \Pr_{\theta}(\hat{\theta}_2 \in \Theta'_m). \quad (10)$$

Moreover, based on definition of  $r_i$ , the true positive probability of class  $m$  is given by

$$p_{tp,i} = \Pr_{\theta}(\hat{\theta}_i \in \Theta_m) = \Pr_{\theta}(\hat{\theta}_i \in \Theta'_m) + \Pr_{\theta}(\hat{\theta}_i \in \Theta_m - \Theta'_m),$$

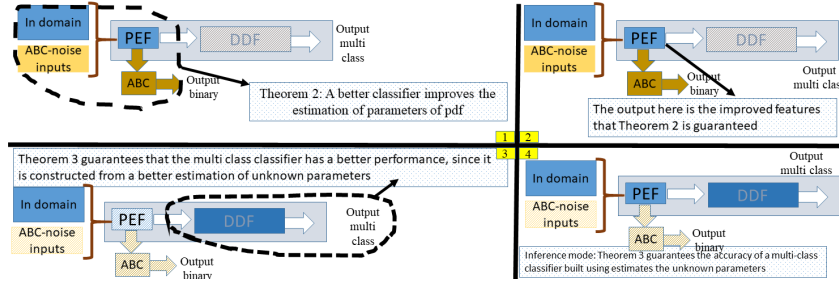


Figure 6: Relationship between the theorems in Section 3 and the proposed method in Section 4.

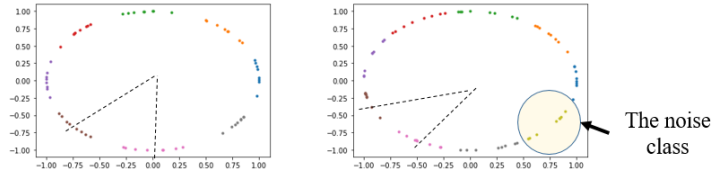


Figure 7: Feature distance between different classes with and without additional background class for a toy example. Left: Contains 8 classes and the feature separation is visibly larger; Right: Contains an additional noise class that decreases the feature distance for all the other classes.

for  $i = 1, 2$ . Additionally, from the Cauchy–Schwarz inequality, we have

$$\Pr_{\theta}(\hat{\theta}_i \in \Theta_m - \Theta'_m) \leq \mu_i(\Theta_m - \Theta'_m) = 0,$$

537 So,  $p_{tp,i} = \Pr_{\theta}(\hat{\theta}_i \in \Theta'_m)$  and from (10) the true positive probability of class  $i$  of  $r_1$  is greater than  
 538  $r_1$ .

539 The error probability of  $r_j$  is given by  $p_{er,j} = 1 - \sum_{i=1}^n P_i P_{tp,i,j}$ , where  $P_i$  is the prior probability  
 540 of class  $i$ . Therefore,  $p_{er,1} \leq p_{er,2}$ .  $\square$

541

#### 542 CONNECTING THE THEOREMS WITH THE PROPOSED METHOD

543 Fig. 6 shows the connection between the proposed theorems and the approach. In part 1, Theorem  
 544 2 connects the estimation of unknown parameters to the auxiliary classifier. In part 2, the learned  
 545 features are passed to a decision making network (result of Theorem 2). In part 3, Theorem 3  
 546 guarantees that the multi-class classifier outperforms other classifiers, because it is using the features  
 547 from a better estimation of unknown parameters of pdf.

548

#### 549 TOY EXAMPLE:

550 We demonstrate the effect of adding background class to the original classifier with a toy example  
 551 and visualize it in Fig. 7. In this example, the input is a sequence of binary bits (+1 and -1) with  
 552 length 3 in white Gaussian noise. the classifier is constructed using two fully connected layers with  
 553 sigmoid and the last layer is normalized on unit circle. As seen from Fig. 7, adding an additional  
 554 noise class visibly reduces the feature separation between all the other classes.

#### 555 IMPLEMENTATION DETAILS

##### 556 FACE RECOGNITION

557 We use Tensorflow (Abadi et al. (2015)) to conduct all our experiments. We train with a batch  
 558 size of 256 on two NVIDIA TeslaV100 (32G) GPUs. We train our models following small (less  
 559 than 1M training images) and large (more than 1M training images) protocol conventions. We use  
 560 CASIA-Webface (Yi et al. (2014)) dataset for small protocol and MS1MV2 dataset for the large  
 561 protocol. We use ResNet-50 (He et al. (2016)) and ResNet-100 models for small and large protocols,  
 562 respectively. The PEF is selected as the first three layers. Following (Deng et al. (2019)), we apply

563 BN (Ioffe & Szegedy (2015)), dropout (Srivastava et al. (2014)) to the last feature map layer followed  
564 by a fully connected layer and batch normalization to obtain the 512-D embedding vector. We set  
565 the feature scale  $s$  parameter to 64 following (Wang et al. (2018); Deng et al. (2019)) and set the  
566 margin parameters  $(m_1, m_2, m_3)$  to (1, 0.5, 0), respectively. For small scale protocol, we start the  
567 learning rate at 0.01 and divide the learning rate by 10 at 40K, 80K, and 100K iterations. We train for  
568 120K iterations. For large scale protocol, we start the learning rate at 0.01 and divide the learning  
569 rate by 10 at 80K, 100K, and 200K iterations. We train for 240K iterations. We use Momentum  
570 optimizer and set the momentum to 0.9 and weight decay to  $5e-4$ . We use the feature centre of all  
571 images from a template or all frames from a video in order to report the results on IJB-B, IJB-C and  
572 YTF datasets. For ABC-noise data, we cropped background images patches from MS1MV2 (Guo  
573 et al. (2016)) dataset and cropped hard examples from the Caltech-101 (F. F. Li et al. (2004)) dataset  
574 plus a few open sourced images (animal faces) using MTCNN (Zhang et al. (2016a)) face detector.  
575 We generated roughly 500K non-face images for training the ABCs.

#### 576 SPEAKER IDENTIFICATION

577 L2 loss and dropout with a rate of 0.2 are applied during training for generalization. The ABC-noise  
578 is collected from silence intervals of the VoxCeleb dataset, where an energy-based voice activity  
579 detection (VAD) is applied to detect the silence intervals. To augment the ABC-noise, Gaussian  
580 noise is added to the silence intervals. Each batch size is set to 64 and the optimizer is ADAM with  
581 a learning rate of 0.001. The VoxCeleb dataset is trained for 11 epochs and the other datasets are  
582 trained for 6 epochs.

#### 583 LFW AND YTF DATASETS

584 LFW database contains the annotations for 5171 faces in a set of 2845 images taken from the Faces  
585 in the Wild data set (Berg et al. (2004)). YouTubeFaces (Wolf et al. (2011)) contains 3,425 videos of  
586 1,595 people. Following the standard convention, we report the results on 5000 video pairs using  
587 unrestricted with labeled outside data protocol.

#### 588 IJB-B AND IJB-C DATASETS

589 The IJB-B contains 1,845 subjects with 21.8K still images and 55K frames from 7,011 videos. In  
590 total, there are 12,115 templates with 10,270 genuine matches and 8M impostor matches. The IJB-C  
591 dataset (Maze et al. (2018)) is a further extension of IJB-B, having 3,531 subjects with 31.3K still  
592 images and 117.5K frames from 11,779 videos. In total, there are 23, 124 templates with 19,557  
593 genuine matches and 15, 639K impostor matches.