

## A Appendix

### A.1 Detail of Feature-Transforming Methods

**$L_2$  normalization ( $L_2$ -norm)** From equation [6], normalizing the norm of the feature vectors can improve the performance of a prototype classifier. We denote a function normalizing the norm as  $\psi_{L_2}$  given by

$$\psi_{L_2}(\phi(\mathbf{x})) = \frac{\phi(\mathbf{x})}{\|\phi(\mathbf{x})\|}. \quad (10)$$

**LDA** From equation [8], decreasing the ratio of the within-class variance to the between-class variance can improve the performance of a prototype classifier. LDA [30] is widely used to search for a projection space that maximizes the between-class variance and minimizes the within-class variance. It computes the eigenvectors of a matrix  $\hat{\Sigma}_\tau^{-1} \hat{\Sigma}$ , where  $\hat{\Sigma}$  is the covariance matrix of prototypes and  $\hat{\Sigma}_\tau$  is the class-conditioned covariance matrix. Since the number of data is small in few-shot settings,  $\hat{\Sigma}_\tau^{-1}$  cannot be estimated stably and we add a regularizer term to  $\hat{\Sigma}_\tau$  and define it as  $\hat{\Sigma}_{\tau\text{reg}}$ .

$$\hat{\Sigma}_{\tau\text{reg}} = \hat{\Sigma}_\tau + \lambda I, \quad (11)$$

where  $\lambda \in R^1, I \in R^{D \times D}$  is identity matrix.

**EST** Since computation of LDA is unstable, we also analyze the effect of EST [2]. EST computes eigenvectors of a matrix  $\hat{\Sigma} - \rho \hat{\Sigma}_\tau$ : the difference between the covariance matrix of the class mean vectors and the class mean covariance matrix with weight parameter  $\rho$ . Similar to LDA, EST also searches for the projection space that maximizes the  $\Sigma$  and minimizes the  $\Sigma_\tau$ .

**EST+ $L_2$ -norm** We hypothesize that the combination of the transforming methods can improve the performance of a prototype classifier independently from each other. Specifically, we focus on reducing equation [6] and equation [8] by combining **EST** and  **$L_2$ -norm**. We first apply EST to reduce equation [8] and after that we apply  $L_2$ -norm to reduce equation [6]. At the end of the operation we want the variance of the norm to be 0, thus we apply EST and after that we apply  $L_2$ -norm.

**LDA+ $L_2$ -norm** We focus on reducing equation [8] and equation [7] by combining **LDA** and  **$L_2$ -norm**. Following the similar procedure of **EST+ $L_2$ -norm**, we first apply LDA and after that we apply  $L_2$ -norm.

### A.2 Existing Upper Bound on Expected Risk for Prototype Classifier

To analyze the behavior of a prototype classifier, we start from the current study [2]. The following theorem is the upper bound of the expected risk of prototypical networks with the next conditions.

- The probability distribution of an extracted feature  $\phi(\mathbf{x})$  given its class  $y = c$  is Gaussian i.e  $\mathcal{D}_y = \mathcal{N}(\mu_c, \Sigma_c)$ , where  $\mu_c = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_c}[\phi(\mathbf{x})]$  and  $\Sigma_c = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_c}[(\phi(\mathbf{x}) - \mu_c)(\phi(\mathbf{x}) - \mu_c)^\top]$ .
- All class-conditioned distributions have the same covariance matrix, i.e.,  $\forall(c, c'), \Sigma_c = \Sigma_{c'}$ .

**Theorem 2** (Cao et al. [2]). *Let  $\mathcal{M}$  be an operation of a prototype classifier on binary classification defined by equation [1]. Then for  $\mu = \mathbb{E}_{c \sim \tau}[\mu_c]$  and  $\Sigma = \mathbb{E}_{c \sim \tau}[(\mu_c - \mu)(\mu_c - \mu)^\top]$ , the misclassification risk of the prototype classifier on binary classification  $R_{\mathcal{M}}$  satisfies*

$$R_{\mathcal{M}}(\phi) \leq 1 - \frac{4 \text{Tr}(\Sigma)^2}{8(1 + 1/k)^2 \text{Tr}(\Sigma_c^2) + 16(1 + 1/k) \text{Tr}(\Sigma \Sigma_c) + \mathbb{E} \text{dist}_{L_2}^2(\mu_{c_1}, \mu_{c_2})}, \quad (12)$$

where  $\mathbb{E} \text{dist}_{L_2}^2(\mu_{c_1}, \mu_{c_2}) = \mathbb{E}_{c_1, c_2} \left[ ((\mu_{c_1} - \mu_{c_2})^\top (\mu_{c_1} - \mu_{c_2}))^2 \right]$ .

We show the detail of the derivation in Appendix [A.3]

### A.3 Reviewing Derivation Details of Theorem 2 (Cao et al. [2])

We briefly review the derivation of Theorem 2. In prototype classifier, from equation 1 and equation 3,  $R_{\mathcal{M}}$  is written with sigmoid function  $\sigma$  as follows:

$$\begin{aligned} R_{\mathcal{M}}(\phi) &= \Pr_{c_1, c_2 \sim \tau, \mathbf{x} \sim \mathcal{D}_{c_1}, S \sim \mathcal{D}^{\otimes 2K}} \left( \sigma(\|\phi(\mathbf{x}) - \overline{\phi(S_{c_2})}\| - \|\phi(\mathbf{x}) - \overline{\phi(S_{c_1})}\|) \leq \frac{1}{2} \right) \\ &= \Pr(\alpha < 0), \end{aligned} \quad (13)$$

where  $\alpha \triangleq \|\phi(\mathbf{x}) - \overline{\phi(S_{c_2})}\| - \|\phi(\mathbf{x}) - \overline{\phi(S_{c_1})}\|$ . we bound equation 13 with expectation and variance of  $\alpha$  by following proposition.

**Proposition 1.** *From the one-sided Chebyshev's inequality, it immediately follows that:*

$$R_{\mathcal{M}}(\phi) = \Pr(\alpha < 0) \leq 1 - \frac{\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha]^2}{\text{Var}_{S, c_1, c_2, \mathbf{x}} [\alpha] + \mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha]^2}. \quad (14)$$

equation 13 can be further write down as follows by Law of Total Expectation

$$\begin{aligned} \text{Var}_{S, c, \mathbf{x}}(\alpha) &= \mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_c} [\alpha^2] - (\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha])^2 \\ &= \mathbb{E}_{c_1, c_2} \mathbb{E}_{\mathbf{x}, S} [\alpha^2 | c_1, c_2] - (\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha])^2 \\ &= \mathbb{E}_{c_1, c_2} [\text{Var}_{\mathbf{x}, S}(\alpha | c_1, c_2) + \mathbb{E}_{\mathbf{x}, S} [\alpha | c_1, c_2]^2] - \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} \mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} [\alpha]^2. \end{aligned}$$

Therefore,

$$\Pr(\alpha < 0) \leq 1 - \frac{\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_c} [\alpha]^2}{\mathbb{E}_{c_1, c_2} [\text{Var}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S} [\alpha] + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S} [\alpha]^2]}. \quad (15)$$

We write down the expection and variance of  $\alpha$  with following Lemmas 3 and 4

**Lemma 3.** *Under the same notation and assumptions as Theorem 2 then,*

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha] &= (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}) \top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}) \\ \mathbb{E}_{c_1, c_2 \sim \tau} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} \mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} [\alpha] &= 2\text{Tr}(\Sigma). \end{aligned}$$

**Lemma 4.** *Under the same notation and assumptions as Theorem 2 then,*

$$\mathbb{E}_{c_1, c_2} [\text{Var}_{\mathbf{x}, S} [\alpha | c_1, c_2]] \leq 8 \left( 1 + \frac{1}{K} \right) \text{Tr} \left( \Sigma_{\tau} \left( \left( 1 + \frac{1}{K} \right) \Sigma_{\tau} + 2\Sigma \right) \right). \quad (16)$$

The proofs of the above lemmas are in the current study [2].

With Proposition 1 and Lemma 3, Lemma 4, we obtain Theorem 2.

### A.4 Derivation Details of Theorem 1

We will describe the detailed derivation of Theorem 1 in this section. Our derivation is different from Cao et al. [2]'s study in the following points.

1. We re-derived Lemma 3 because the term of the difference between the trace of the class covariance matrices is erased in the lemma. This term cannot omit in our derivation since we do not assume the class covariance matrix to be the same among classes.
2. We re-derived the bound on the variance of squared Euclidean distance of two vectors, e.g Lemma 4. The derivation of Cao et al. [2] uses the property of quadratic forms of normally distributed random variables and the fact that the sum of normally distributed random variables is also distributed in Gaussian distribution. The calculation of the variance of squared  $L_2$ -norm without depending on the property of some distributions is not straightforward [38]. We divide the variance of squared Euclidean distance of two vectors into the variance of the norm of the feature vectors and the variance of the inner-product of vectors. Then we apply Cauchy–Schwarz inequality to the inner-product.

We start the proof from equation [15](#). We first prove the following Lemma [5](#) related to the expectation statistics of  $\alpha$  in equation [15](#).

**Lemma 5.** *Under the same notations and assumptions as Theorem [1](#) then,*

$$\begin{aligned}\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha] &= \frac{1}{K} (\text{Tr}(\Sigma_{c_2}) - \text{Tr}(\Sigma_{c_1})) + (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}) \\ \mathbb{E}_{c_1, c_2 \sim \mathcal{T}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} \mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} [\alpha] &= 2\text{Tr}(\Sigma).\end{aligned}$$

*Proof.* First, from the definition of  $\alpha$ , we split  $\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha]$  in to two parts and examine them separately.

$$\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha] = \underbrace{\mathbb{E} \left[ \left\| \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right\|^2 \right]}_{(i)} - \underbrace{\mathbb{E} \left[ \left\| \phi(\mathbf{x}) - \overline{\phi(S_{c_1})} \right\|^2 \right]}_{(ii)}. \quad (17)$$

In regular conditions, for random vector  $X$ , the expectation of the norm is

$$\mathbb{E}[X^\top X] = \text{Tr}(\text{Var}(X)) + \mathbb{E}[X]^\top \mathbb{E}[X], \quad (18)$$

and the variance of the vector is

$$\text{Var}(X) = \mathbb{E}[X X^\top] - \mathbb{E}[X] \mathbb{E}[X]^\top \quad (19)$$

$$\Sigma_{c_i} \triangleq \text{Var}_{\mathbf{x} \sim \mathcal{D}_{c_i}} (\phi(\mathbf{x})). \quad (20)$$

Hence,

$$\begin{aligned}(i) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} \mathbb{E}_S \left[ \left\| \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right\|^2 \right] \\ &= \text{Tr} \left( \text{Var}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S} \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right] \right) + \mathbb{E}_{\mathbf{x}} \mathbb{E}_S \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right]^\top \mathbb{E}_{\mathbf{x}} \mathbb{E}_S \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right],\end{aligned} \quad (21)$$

where the first term inside the trace can be expanded as:

$$\begin{aligned}\text{Var} \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right] &= \mathbb{E} \left[ \left( \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right) \left( \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right)^\top \right] - (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})(\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top \\ &= \text{Var} [\phi(\mathbf{x})] + \mathbb{E} [\phi(\mathbf{x})] \mathbb{E} [\phi(\mathbf{x})]^\top + \text{Var} \left[ \overline{\phi(S_{c_2})} \right] + \mathbb{E} \left[ \overline{\phi(S_{c_2})} \right] \mathbb{E} \left[ \overline{\phi(S_{c_2})} \right]^\top \\ &\quad - \boldsymbol{\mu}_{c_2} \boldsymbol{\mu}_{c_1}^\top - \boldsymbol{\mu}_{c_1} \boldsymbol{\mu}_{c_2}^\top - (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})(\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top \\ &= \Sigma_{c_1} + \frac{1}{K} \Sigma_{c_2} \quad (\text{Last terms cancel out}).\end{aligned} \quad (22)$$

The second term in equation [21](#) is simply as follows.

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} \mathbb{E}_S \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right] = \boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}. \quad (23)$$

From equation [22](#) and equation [23](#) we obtain

$$(i) = \text{Tr}(\Sigma_{c_1}) + \frac{1}{K} \text{Tr}(\Sigma_{c_2}) + (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}). \quad (24)$$

Similarly for ii,

$$\begin{aligned}(ii) &= \text{Tr} \left( \text{Var}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S} \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_1})} \right] \right) + \mathbb{E}_{\mathbf{x}} \mathbb{E}_S \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_1})} \right]^\top \mathbb{E}_{\mathbf{x}} \mathbb{E}_S \left[ \phi(\mathbf{x}) - \overline{\phi(S_{c_1})} \right] \\ &= \text{Tr}(\Sigma_{c_1}) + \frac{1}{K} \text{Tr}(\Sigma_{c_1}).\end{aligned} \quad (25)$$

From (i) and (ii), and equation [17](#)

$$\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}} [\alpha] = \frac{1}{K} (\text{Tr}(\Sigma_{c_2}) - \text{Tr}(\Sigma_{c_1})) + (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}). \quad (26)$$

Since  $\mathbb{E}_{c_1, c_2, \mathbf{x}, S}[\alpha] = \mathbb{E}_{c_1, c_2 \sim \tau} [\mathbb{E}_{S \sim \mathcal{D}^{\otimes 2k}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}}[\alpha]]$ ,

$$\begin{aligned}
\mathbb{E}_{c_1, c_2, \mathbf{x}, S}[\alpha] &= \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \frac{1}{K} (\text{Tr}(\Sigma_{c_2}) - \text{Tr}(\Sigma_{c_1})) + (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}) \right] \\
&= \mathbb{E}_{c_1, c_2 \sim \tau} \left[ (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}) \right] \\
&= \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} + \boldsymbol{\mu}_{c_2}^\top \boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_2}^\top \boldsymbol{\mu}_{c_1} \right] \\
&= \text{Tr}(\Sigma) + \boldsymbol{\mu}^\top \boldsymbol{\mu} + \text{Tr}(\Sigma) + \boldsymbol{\mu}^\top \boldsymbol{\mu} - 2\boldsymbol{\mu}^\top \boldsymbol{\mu} \\
&= 2\text{Tr}(\Sigma). \tag{27}
\end{aligned}$$

□

Next we prove the following Lemma 6 related to the conditioned variance of  $\alpha$ .

**Lemma 6.** *Under the same notation and assumptions as Theorem 7*

$$\mathbb{E}_{c_1, c_2} \text{Var}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S \sim \mathcal{D}^{\otimes 2N}}[\alpha] \leq \frac{4}{K} \mathbb{E}_{c \sim \tau} \text{Var} \left[ \|\phi(\mathbf{x})\|^2 \right] + \frac{4}{K} \text{Var}_{c \sim \tau} [\text{Tr}(\Sigma_c)] + \text{V}_{\text{wit}}(\Sigma_\tau, \Sigma, \boldsymbol{\mu}), \tag{28}$$

where

$$\begin{aligned}
\text{V}_{\text{wit}}(\Sigma_\tau, \Sigma, \boldsymbol{\mu}) &= \frac{8}{K} (\text{Tr}(\Sigma_\tau)) (\text{Tr}(\Sigma) + \boldsymbol{\mu}^\top \boldsymbol{\mu}) + 4 (\text{Tr}(\Sigma) + \boldsymbol{\mu}^\top \boldsymbol{\mu})^2 \\
&\quad + 4 \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \text{Tr}(\Sigma_{c_1}) (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right]. \tag{29}
\end{aligned}$$

*Proof.* We start with the inequality between the variance of 2 random variables. We define  $\text{Cov}(A, B)$  as covariance of 2 random variables  $A, B$ .

$$\begin{aligned}
\text{Var}[A + B] &= \text{Var}[A] + \text{Var}[B] + 2\text{Cov}(A, B) \\
&\leq \text{Var}[A] + \text{Var}[B] + 2\sqrt{\text{Var}[A]\text{Var}[B]} \\
&\leq \text{Var}[A] + \text{Var}[B] + 2 \cdot \frac{\text{Var}[A] + \text{Var}[B]}{2} \\
&= 2\text{Var}[A] + 2\text{Var}[B]. \tag{30}
\end{aligned}$$

For  $\text{Var}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S}[\alpha]$ ,

$$\begin{aligned}
\text{Var}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S}[\alpha] &= \text{Var} \left[ \left\| \phi(\mathbf{x}) - \overline{\phi(S_{c_2})} \right\|^2 - \left\| \phi(\mathbf{x}) - \overline{\phi(S_{c_1})} \right\|^2 \right] \\
&= \text{Var} \left[ \left\| \overline{\phi(S_{c_1})} \right\|^2 - \left\| \overline{\phi(S_{c_2})} \right\|^2 - 2\phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})}) \right] \\
&\leq 2\text{Var} \left[ \left\| \overline{\phi(S_{c_1})} \right\|^2 - \left\| \overline{\phi(S_{c_2})} \right\|^2 \right] \\
&\quad + 4\text{Var} \left[ \phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})}) \right] \quad (\because \text{equation 30}) \\
&= 2\text{Var} \left[ \left\| \overline{\phi(S_{c_1})} \right\|^2 \right] + 2\text{Var} \left[ \left\| \overline{\phi(S_{c_2})} \right\|^2 \right] + 4\text{Var} \left[ \phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})}) \right]. \tag{31}
\end{aligned}$$

From 3rd line to 4th line we decompose the variance of  $\left\| \overline{\phi(S_{c_1})} \right\|^2 - \left\| \overline{\phi(S_{c_2})} \right\|^2$  use the independence of  $\overline{\phi(S_{c_1})}$  and  $\overline{\phi(S_{c_2})}$  with their class given.

Next we focus on  $\text{Var} \left[ \phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})}) \right]$ .

$$\begin{aligned}
\text{Var} \left[ \phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})}) \right] &= \mathbb{E} \left[ \left( \phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})}) \right)^2 \right] \\
&\quad - \left( \mathbb{E} [\phi(\mathbf{x})]^\top \mathbb{E} \left[ \overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})} \right] \right)^2 \\
&\leq \mathbb{E} \left[ \left( \|\phi(\mathbf{x})\|^2 \left\| \overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})} \right\|^2 \right) \right] - (\boldsymbol{\mu}_{c_1}^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}))^2 \\
&= \mathbb{E} \left[ \|\phi(\mathbf{x})\|^2 \right] \mathbb{E} \left[ \left\| \overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})} \right\|^2 \right] - (\boldsymbol{\mu}_{c_1}^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}))^2.
\end{aligned} \tag{32}$$

From the 2nd line to the 3rd line we use Cauchy–Schwarz inequality.

For  $\mathbb{E} \left[ \left\| \overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})} \right\|^2 \right]$ , with equation 18

$$\mathbb{E} \left[ \left\| \overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})} \right\|^2 \right] = \frac{1}{K} (\text{Tr}(\boldsymbol{\Sigma}_{c_1}) + \text{Tr}(\boldsymbol{\Sigma}_{c_2})) + (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}). \tag{33}$$

Thus  $\mathbb{E}_{c_1, c_2 \sim \tau} \left[ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1, S}} \left[ \|\phi(\mathbf{x})\|^2 \right] \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1, S}} \left[ \left\| \overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})} \right\|^2 \right] \right]$  is calculated as follows.

$$\begin{aligned}
&\mathbb{E}_{c_1, c_2 \sim \tau} \left[ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1, S}} \left[ \|\phi(\mathbf{x})\|^2 \right] \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1, S}} \left[ \left\| \overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})} \right\|^2 \right] \right] \\
&= \mathbb{E}_{c_1, c_2 \sim \tau} \left[ (\text{Tr}(\boldsymbol{\Sigma}_{c_1}) + \boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1}) \left( \frac{1}{K} (\text{Tr}(\boldsymbol{\Sigma}_{c_1}) + \text{Tr}(\boldsymbol{\Sigma}_{c_2})) + (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right) \right] \\
&= \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \frac{1}{K} \left( \text{Tr}(\boldsymbol{\Sigma}_{c_1})^2 + \text{Tr}(\boldsymbol{\Sigma}_{c_1}) \text{Tr}(\boldsymbol{\Sigma}_{c_2}) \right) \right] \\
&\quad + \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \frac{2}{K} (\text{Tr}(\boldsymbol{\Sigma}_\tau)) \boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} + \boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right] \\
&\quad + \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \text{Tr}(\boldsymbol{\Sigma}_{c_1}) (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right] \\
&= \frac{1}{K} \left( \mathbb{E}_{c_1, c_2 \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_{c_1})^2] + \mathbb{E}_{c_1, c_2 \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_{c_1})]^2 \right) \\
&\quad + \frac{2}{K} (\text{Tr}(\boldsymbol{\Sigma}_\tau)) (\text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu}) + \mathbb{E} [\boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})] \\
&\quad + \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \text{Tr}(\boldsymbol{\Sigma}_{c_1}) (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right] \\
&= \frac{1}{K} \text{Var}_{c \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_c)] \\
&\quad + \frac{2}{K} \text{Tr}(\boldsymbol{\Sigma}_\tau)^2 + \frac{2}{K} (\text{Tr}(\boldsymbol{\Sigma}_\tau)) (\text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu}) + \mathbb{E}_{c_1, c_2} \left[ \boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right] \\
&\quad + \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \text{Tr}(\boldsymbol{\Sigma}_{c_1}) (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right] \\
&= \frac{1}{K} \text{Var}_{c \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_c)] \\
&\quad + \frac{2}{K} \text{Tr}(\boldsymbol{\Sigma}_\tau) (\text{Tr}(\boldsymbol{\Sigma}_\tau) + \text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu}) + \mathbb{E}_{c_1, c_2} \left[ \boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right] \\
&\quad + \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \text{Tr}(\boldsymbol{\Sigma}_{c_1}) (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) \right].
\end{aligned} \tag{34}$$

Now we take into account the term  $-(\boldsymbol{\mu}_{c_1}^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}))^2$  in equation [32](#),

$$\begin{aligned}
& \mathbb{E}_{c_1, c_2} [\boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}) - (\boldsymbol{\mu}_{c_1}^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1}))^2] \\
&= \mathbb{E}_{c_1, c_2} [\boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} \boldsymbol{\mu}_{c_2}^\top \boldsymbol{\mu}_{c_2} - (\boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_2})^2] \\
&\leq \mathbb{E}_{c_1, c_2} [\boldsymbol{\mu}_{c_1}^\top \boldsymbol{\mu}_{c_1} \boldsymbol{\mu}_{c_2}^\top \boldsymbol{\mu}_{c_2}] \\
&= (\text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu})^2.
\end{aligned} \tag{35}$$

Thus  $\mathbb{E}_{c_1, c_2} [\text{Var}[\phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})})]]$  is calculated as follows.

$$\begin{aligned}
& \mathbb{E}_{c_1, c_2} [\text{Var}[\phi(\mathbf{x})^\top (\overline{\phi(S_{c_2})} - \overline{\phi(S_{c_1})})]] \\
&= \frac{1}{K} \text{Var}_{c \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_c)] + \frac{2}{K} \text{Tr}(\boldsymbol{\Sigma}_\tau) (\text{Tr}(\boldsymbol{\Sigma}_\tau) + \text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu}) + (\text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu})^2 \\
&\quad + \mathbb{E}_{c_1, c_2 \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_{c_1}) (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})].
\end{aligned} \tag{36}$$

Regarding  $\|\overline{\phi(S_c)}\|^2$ , since the function computing square norm is convex, next equation holds with  $D$ -dimensional Jensen's inequality [39](#):

$$\begin{aligned}
\|\overline{\phi(S_c)}\|^2 &= \left\| \frac{1}{K} \sum_{\substack{i=0 \\ \mathbf{x} \in S_c}} \phi(\mathbf{x}) \right\|^2 \\
&\leq \frac{1}{K} \left\| \sum_{\substack{i=0 \\ \mathbf{x} \in S_c}} \phi(\mathbf{x}) \right\|^2.
\end{aligned} \tag{37}$$

Combining equation [31](#), equation [36](#), and equation [37](#) we obtain

$$\begin{aligned}
& \mathbb{E}_{c_1, c_2} \text{Var}_{\mathbf{x} \sim D_{c_1}, S \sim \mathcal{D}^{\otimes 2N}} [\alpha] \\
&\leq \frac{4}{K} \mathbb{E}_{c \sim \tau} \text{Var}_{\mathbf{x} \sim \mathcal{D}_c} [\|\phi(\mathbf{x})\|^2] + \frac{4}{K} \text{Var}_{c \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_c)] \\
&\quad + \frac{8}{K} \text{Tr}(\boldsymbol{\Sigma}_\tau) (\text{Tr}(\boldsymbol{\Sigma}_\tau) + \text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu}) + 4 (\text{Tr}(\boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \boldsymbol{\mu})^2 \\
&\quad + 4 \mathbb{E}_{c_1, c_2 \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_{c_1}) (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})^\top (\boldsymbol{\mu}_{c_2} - \boldsymbol{\mu}_{c_1})].
\end{aligned} \tag{38}$$

□

To complete the proof of Theorem [1](#), we calculate  $\mathbb{E}_{c_1, c_2} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S \sim \mathcal{D}^{\otimes 2K}} [\alpha]^2$ .

$$\begin{aligned}
& \mathbb{E}_{c_1, c_2} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{c_1}, S \sim \mathcal{D}^{\otimes 2K}} [\alpha]^2 \\
&= \mathbb{E}_{c_1, c_2} \left[ \left( \frac{1}{K} (\text{Tr}(\boldsymbol{\Sigma}_{c_2}) - \text{Tr}(\boldsymbol{\Sigma}_{c_1})) + (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}) \right)^2 \right] \\
&= \mathbb{E}_{c_1, c_2} \left[ \left( \frac{1}{K} (\text{Tr}(\boldsymbol{\Sigma}_{c_2}) - \text{Tr}(\boldsymbol{\Sigma}_{c_1})) \right)^2 \right] + \mathbb{E}_{c_1, c_2} \left[ ((\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}))^2 \right] \\
&\quad + 2 \mathbb{E}_{c_1, c_2} \left[ \left( \frac{1}{K} (\text{Tr}(\boldsymbol{\Sigma}_{c_2}) - \text{Tr}(\boldsymbol{\Sigma}_{c_1})) \right) ((\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})) \right] \\
&= \text{Var}_{c_1, c_2 \sim \tau} \left[ \frac{1}{K} (\text{Tr}(\boldsymbol{\Sigma}_{c_2}) - \text{Tr}(\boldsymbol{\Sigma}_{c_1})) \right] + \left( \mathbb{E}_{c_1, c_2 \sim \tau} \left[ \frac{1}{K} (\text{Tr}(\boldsymbol{\Sigma}_{c_2}) - \text{Tr}(\boldsymbol{\Sigma}_{c_1})) \right] \right)^2 \\
&\quad + \mathbb{E}_{c_1, c_2} \left[ ((\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}))^2 \right] \\
&= \frac{2}{K^2} \text{Var}_{c \sim \tau} [\text{Tr}(\boldsymbol{\Sigma}_c)] + \mathbb{E}_{c_1, c_2} \left[ ((\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2})^\top (\boldsymbol{\mu}_{c_1} - \boldsymbol{\mu}_{c_2}))^2 \right].
\end{aligned} \tag{39}$$

From 2nd line to 3rd line, we use the symmetry of the last term with respect to  $c_1$  and  $c_2$  and erase the term.

Combining equation [15](#) Lemma [5](#) Lemma [6](#) and equation [39](#) we obtain the bound.

### A.5 Theorem [1](#) with $N$ -way Classification

The upper bound on the risk of  $N$ -way prototype classifier is as follows.

**Theorem 7.** *Let operation of binary class prototype classifier  $\mathcal{M}$  as defined in equation [1](#). Then for  $\overline{\phi(S_c)} = \frac{1}{K} \sum_{\mathbf{x} \in S_c} \phi(\mathbf{x})$ ,  $\mu_c = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_c}[\phi(\mathbf{x})]$ ,  $\Sigma_c = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_c}[(\phi(\mathbf{x}) - \mu_c)(\phi(\mathbf{x}) - \mu_c)^\top]$ ,  $\mu = \mathbb{E}_{c \sim \tau}[\mu_c]$ ,  $\Sigma = \mathbb{E}_{c \sim \tau}[(\mu_c - \mu)(\mu_c - \mu)^\top]$ ,  $\Sigma_\tau = \mathbb{E}_{c \sim \tau}[\Sigma_c]$ , if  $\phi(\mathbf{x})$  has its fourth moment, miss classification risk of binary class prototype classifier  $R_{\mathcal{M}}$  satisfy*

$$\begin{aligned} & R(\mathcal{M}(\phi, \mathbf{x}, \{S_i\}_{i=1}^N), y) \\ & \leq N - 1 - \sum_{\substack{c=1 \\ c \neq y}}^N \frac{4(\text{Tr}(\Sigma))^2}{\mathbb{E}V[h_{L_2}(\phi(\mathbf{x}))] + V_{\text{Tr}}(\Sigma_y) + V_{\text{wit}}(\Sigma_\tau, \Sigma, \boldsymbol{\mu}) + \mathbb{E} \text{dist}_{L_2}^2(\boldsymbol{\mu}_y, \boldsymbol{\mu}_c)}, \end{aligned} \quad (40)$$

where

$$\begin{aligned} \mathbb{E}V[h_{L_2}(\phi(\mathbf{x}))] &= \frac{4}{K} \mathbb{E}_{y \sim \tau} \left[ \text{Var}_{\mathbf{x}_c \sim \mathcal{D}_c} \left[ \|\phi(\mathbf{x})\|^2 \right] \right], \\ V_{\text{Tr}}(\Sigma_y) &= \left( \frac{4}{K} + \frac{2}{K^2} \right) \text{Var}_{c \sim \tau} [\text{Tr}(\Sigma_c)], \\ V_{\text{wit}}(\Sigma_\tau, \Sigma, \boldsymbol{\mu}) &= \frac{8}{K} (\text{Tr}(\Sigma_\tau)) (\text{Tr}(\Sigma) + \boldsymbol{\mu}^\top \boldsymbol{\mu}) + 4 (\text{Tr}(\Sigma) + \boldsymbol{\mu}^\top \boldsymbol{\mu})^2 \\ & \quad + 4 \mathbb{E}_{c \sim \tau} \left[ \text{Tr}(\Sigma_y) (\boldsymbol{\mu}_y - \boldsymbol{\mu}_c)^\top (\boldsymbol{\mu}_y - \boldsymbol{\mu}_c) \right], \\ \mathbb{E} \text{dist}_{L_2}^2(\boldsymbol{\mu}_y, \boldsymbol{\mu}_c) &= \mathbb{E}_{y,c} \left[ \left( (\boldsymbol{\mu}_y - \boldsymbol{\mu}_c)^\top (\boldsymbol{\mu}_y - \boldsymbol{\mu}_c) \right)^2 \right]. \end{aligned} \quad (41)$$

*Proof.* Let  $x, y$  be the input and its class of a query data. We define  $\alpha_c$  by  $\alpha_c = \left\| \phi(\mathbf{x}) - \overline{\phi(S_c)} \right\|^2 - \left\| \phi(\mathbf{x}) - \overline{\phi(S_y)} \right\|^2$ . Then a prototype classifier miss-classify a class of input  $x, \hat{y}$ , if  $\exists c \in [1, N], c \neq y, \alpha_c < 0$ . Hence:  $R_{\mathcal{M}}(\phi) = \Pr(\bigcup_{\substack{c=1 \\ c \neq y}}^N \alpha_c < 0)$

By Frechet's inequality, next equation holds:

$$R_{\mathcal{M}}(\phi) \leq \sum_{\substack{c=1 \\ c \neq y}}^N \Pr(\alpha_i < 0).$$

Noting that Theorem [1](#) can be applied to each term in the summation and then we obtain Theorem [7](#)  $\square$

### A.6 Detailed performance results

We show in Table [5](#) the detailed performance results of standard object recognition in this section with 95% confidence margin. The table shows the similar result with Table [1](#). We can observe that the prototype classifier with  $L_2$ -norm, EST+ $L_2$ -norm, LDA+ $L_2$ -norm performs comparably with ProtoNet and the linear-evaluation-based approach in most of the settings.

Table 5: Classification accuracies with ResNet18 on *mini*ImageNet, *tiered*ImageNet, CIFAR100, linear-evaluation-based methods [3], centering with  $L_2$ -norm [24], and ours. The Baseline without linear-evaluation methods with accuracy greater than the lower 95% confidence margin of the accuracy of ProtoNet and Baseline are in bold. Regarding to Baseline++, Baseline+++ without linear-evaluation methods with accuracy greater than the lower 95% confidence margin of the accuracy of ProtoNet and Baseline+++ are in bold.

	<i>mini</i> ImageNet			<i>tiered</i> ImageNet			CIFAR-100			FC100		
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
<i>methods with meta-learning or linear-evaluation</i>												
ProtoNet [1]	56.74 ± 0.84	75.64 ± 0.62	61.75 ± 0.94	81.56 ± 0.68	65.46 ± 0.96	79.52 ± 0.66	35.92 ± 0.71	50.86 ± 0.71				
Baseline [3]	55.41 ± 0.82	76.95 ± 0.61	63.38 ± 0.91	83.18 ± 0.64	65.12 ± 0.88	79.68 ± 0.64	40.06 ± 0.68	57.04 ± 0.71				
<i>A prototype classifier with feature-transformation methods</i>												
Baseline-w/o-linear, centering+ $L_2$ -norm [24]	<b>57.67 ± 0.83</b>	75.50 ± 0.67	<b>65.26 ± 0.88</b>	81.63 ± 0.64	<b>65.26 ± 0.86</b>	78.58 ± 0.66	<b>41.51 ± 0.72</b>	<b>56.44 ± 0.74</b>				
Baseline-w/o-linear	43.86 ± 0.80	72.36 ± 0.92	56.16 ± 0.89	80.33 ± 0.66	54.06 ± 0.85	77.76 ± 0.64	35.90 ± 0.61	55.20 ± 0.77				
Baseline-w/o-linear, $L_2$ -norm	<b>56.57 ± 0.80</b>	<b>76.44 ± 0.61</b>	<b>65.19 ± 0.87</b>	<b>82.93 ± 0.66</b>	<b>64.76 ± 0.87</b>	<b>79.90 ± 0.66</b>	<b>40.96 ± 0.71</b>	<b>57.71 ± 0.76</b>				
Baseline-w/o-linear, EST	44.19 ± 0.82	70.85 ± 0.92	57.05 ± 0.93	73.16 ± 0.78	57.16 ± 0.88	78.98 ± 0.69	<b>48.43 ± 0.91</b>	<b>63.75 ± 0.73</b>				
Baseline-w/o-linear, EST+ $L_2$ -norm	<b>56.39 ± 0.79</b>	<b>76.24 ± 0.64</b>	<b>64.71 ± 0.91</b>	<b>83.24 ± 0.67</b>	<b>65.54 ± 0.71</b>	<b>79.80 ± 0.71</b>	<b>50.50 ± 0.97</b>	<b>65.31 ± 0.77</b>				
Baseline-w/o-linear, LDA	47.14 ± 0.81	69.87 ± 0.65	57.05 ± 0.90	81.19 ± 0.65	56.49 ± 0.86	78.88 ± 0.66	<b>50.29 ± 0.85</b>	<b>63.56 ± 0.77</b>				
Baseline-w/o-linear, LDA+ $L_2$ -norm	<b>56.37 ± 0.81</b>	<b>76.39 ± 0.66</b>	<b>65.44 ± 0.94</b>	<b>83.16 ± 0.62</b>	<b>65.64 ± 0.84</b>	<b>80.72 ± 0.67</b>	<b>50.40 ± 0.97</b>	<b>65.31 ± 0.77</b>				
<i>methods with meta-learning or linear-evaluation</i>												
Baseline++ [3]	55.07 ± 0.81	74.71 ± 0.61	64.02 ± 0.92	83.18 ± 0.64	65.64 ± 0.93	79.80 ± 0.66	36.93 ± 0.70	50.41 ± 0.73				
<i>A prototype classifier with feature-transformation methods</i>												
Baseline+++w/o-linear, centering+ $L_2$ -norm [24]	<b>57.00 ± 0.64</b>	74.06 ± 0.61	<b>64.92 ± 0.91</b>	81.41 ± 0.66	<b>66.14 ± 0.93</b>	<b>80.03 ± 0.71</b>	<b>38.30 ± 0.74</b>	<b>51.06 ± 0.70</b>				
Baseline+++w/o-linear	36.80 ± 0.76	63.76 ± 0.71	48.27 ± 0.91	75.87 ± 0.71	50.18 ± 0.96	74.23 ± 0.70	30.76 ± 0.62	47.62 ± 0.69				
Baseline+++w/o-linear, $L_2$ -norm	<b>57.21 ± 0.83</b>	<b>74.89 ± 0.65</b>	<b>66.67 ± 0.94</b>	<b>82.49 ± 0.68</b>	<b>66.84 ± 0.92</b>	<b>80.49 ± 0.71</b>	<b>38.55 ± 0.72</b>	<b>51.15 ± 0.71</b>				
Baseline+++w/o-linear, EST	47.21 ± 0.77	69.11 ± 0.64	53.49 ± 0.90	71.81 ± 0.75	53.36 ± 0.93	73.70 ± 0.74	<b>39.92 ± 0.88</b>	<b>53.61 ± 0.70</b>				
Baseline+++w/o-linear, EST+ $L_2$ -norm	<b>57.00 ± 0.75</b>	<b>76.13 ± 0.64</b>	<b>65.52 ± 0.97</b>	<b>79.96 ± 0.74</b>	<b>66.57 ± 0.95</b>	<b>79.42 ± 0.70</b>	<b>42.89 ± 0.82</b>	<b>56.82 ± 0.70</b>				
Baseline+++w/o-linear, LDA	45.95 ± 0.78	68.34 ± 0.68	49.09 ± 0.90	76.17 ± 0.68	52.98 ± 0.91	73.36 ± 0.75	<b>40.27 ± 0.79</b>	<b>55.80 ± 0.74</b>				
Baseline+++w/o-linear, LDA+ $L_2$ -norm	<b>56.78 ± 0.87</b>	<b>75.83 ± 0.65</b>	<b>66.68 ± 0.90</b>	<b>82.53 ± 0.67</b>	<b>65.89 ± 0.93</b>	<b>79.47 ± 0.70</b>	<b>44.21 ± 0.86</b>	<b>58.10 ± 0.70</b>				



### A.7 Visualization of the feature distribution

We show in Figure 3 the distribution of features on testset of FC100 before and after applying the data transformation. From figure 3, we can observe that the feature-transformation-methods slightly change the distributions even the transformations improve the performance of a prototype classifier. The projection into a lower dimensional space for visualization does not accurately represent the relationship of a higher dimension.

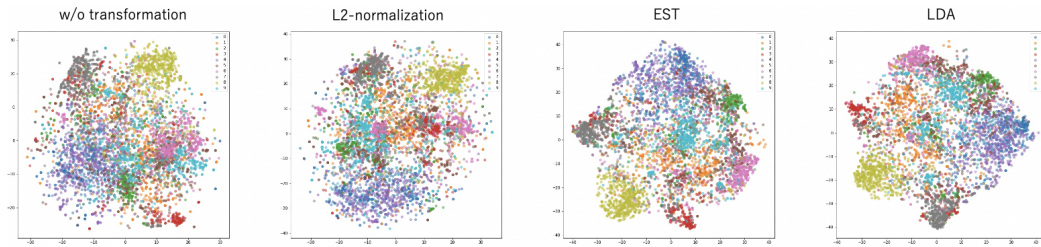


Figure 3