

A Statistical Model for Predicting Generalization in Few-Shot Classification

Yassir Bendou¹, Vincent Gripon¹, Bastien Pasdeloup¹, Giulia Lioi¹, Lukas Mauch², Stefan Uhlich², Fabien Cardinaux², Ghouthi Boukli Hacene^{2,3} and Javier Alonso Garcia²

¹IMT Atlantique, Lab-STICC, UMR CNRS 6285, F-29238 Brest, France

²Sony Europe, R&D Center, Stuttgart Laboratory 1, Germany

³Mila, Montréal, Canada

¹name.surname@imt-atlantique.fr, ²name.surname@sony.com

Abstract—The estimation of the generalization error of classifiers often relies on a validation set. Such a set is hardly available in few-shot learning scenarios, a highly disregarded shortcoming in the field. In these scenarios, it is common to rely on features extracted from pre-trained neural networks combined with distance-based classifiers such as nearest class mean. In this work, we introduce a Gaussian model of the feature distribution. By estimating the parameters of this model, we are able to predict the generalization error on new classification tasks with few samples. We observe that accurate distance estimates between class-conditional densities are the key to accurate estimates of the generalization performance. Therefore, we propose an unbiased estimator for these distances and integrate it in our numerical analysis. We empirically show that our approach outperforms alternatives such as the leave-one-out cross-validation strategy.

Index Terms—few-shot learning, classification, deep learning, generalization

I. INTRODUCTION

During the last decade, the problem of few-shot classification, that is to say classification with very few training samples (typically less than ten per class [1]), has known a large number of contributions [2]–[5]. Many current state-of-the-art solutions consist in using a pre-trained deep feature extractor to embed samples in a feature space where classes are expected to be easier to discriminate. Then, distance-based classifiers such as nearest-class mean (NCM) are applied to the obtained feature vectors [6]–[8].

In the few-shot scenario, performance prediction is not straightforward. The classical approach is to perform cross-validation which is hardly feasible with few training samples. Indeed, removing training samples and using them as a holdout test set severely impacts the generalization ability and at the same time, using too few such samples to measure accuracy on previously unseen data leads to poor performance estimates [9]. The usual strategy consists in “leaving-one-out cross-validation” [10] where one averages the generalization estimated by removing a single training sample of each class over multiple such independent random removals.

Finding an alternative to cross-validation to measure generalization ability has been extensively studied in the literature for standard classification settings where a large number of training samples is available [11]. In this work, we are mainly

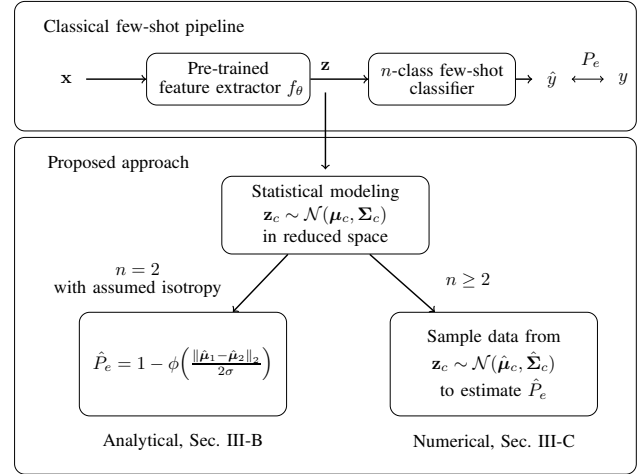


Fig. 1: The classical few-shot pipeline consists in using a pre-trained feature extractor to embed data x into features $z = f_\theta(x)$ and combine it with a classifier. Our approach consists of modeling the class-densities of the features as Gaussian distributions in a lower dimensional subspace and predicting the probability error of the classifier \hat{P}_e . Depending on the number of classes, we either use an analytical form of the probability error or estimate it by sampling a large number of datapoints.

interested in proposing an alternative to cross-validation and to existing generalization prediction methods in the context of few-shot classification. The proposed method estimates the parameters of a statistical model of the class-conditional densities at the output of the feature extractor. The model is then used to predict the probability of error (c.f. Figure 1).

There are two key difficulties here: 1) The first one is that we need to identify a statistical model that is both expressive enough to accurately capture the data distribution and at the same time depends on as few parameters as possible, such that we can accurately estimate them even in the very low data regime. 2) The second one is that our derived expression for the probability of error depends on the distances between class centers. We observe that the naive estimate for the distances is

biased, which leads to underestimating the probability of error especially when working in high-dimensional spaces with very few samples.

The main contributions¹ of our work are as follows:

- We introduce a statistical model of class-conditional densities in the feature space;
- Using this model, we obtain statistical bounds on the generalization error;
- We show that the naive estimator for the distances between class centers is biased and, to alleviate this problem, we propose an unbiased estimator instead;
- We provide experiments and show that our method outperforms other model-free generalization error predictors, such as the leaving-one-out cross-validation, in the context of standardized few-shot classification benchmarks.

II. RELATED WORK

A. Classification in few-shot learning

Since there are few training samples in few-shot classification tasks, training a deep neural network architecture that is typically constructed from a large number of parameters is out of the question. The main approach is to use a pre-trained deep feature extractor and combine it with a simple classifier trained for the few-shot task. In this paper, we refer to P for probabilities and p for probability density functions. We formalize few-shot classification as follows:

Definition 1 (Few-shot classification). Let $\mathcal{D}_B = \{(\mathbf{x}'_i, y'_i)\}_{i=1}^m$ be a large base dataset where $\forall i, (\mathbf{x}'_i, y'_i)$ are i.i.d samples drawn from the true joint probability distribution $p_{\mathcal{X}_B, \mathcal{Y}_B}$ and $\forall i, y'_i \in \mathcal{Y}_B$ is the class associated with \mathbf{x}'_i . We are also given a small few-shot training dataset $\mathcal{D} = \{(\mathbf{x}_j, y_j)\}_{j=1}^\ell$, where $\forall j, y_j \in \mathcal{Y} \neq \mathcal{Y}_B$ and $\forall j, (\mathbf{x}_j, y_j) \sim p_{\mathcal{X}, \mathcal{Y}}$.

The goal of few-shot classification is to train a classifier C using \mathcal{D} , with the potential help of \mathcal{D}_B .

In the literature, most methods consist in training a deep feature extractor f_θ with parameter set θ using \mathcal{D}_B . This feature extractor is then either adapted or used as is on \mathcal{D} to produce feature vectors $\mathbf{z} = f_\theta(\mathbf{x})$ in a Euclidean space. The few-shot classifier then works with the modified training dataset $f_\theta(\mathcal{D}) = \{(\mathbf{z}_j, y_j)\}_{j=1}^\ell$.

There are multiple strategies on how to train f_θ that can roughly be classified as optimization-based approaches and distance metric approaches.

Optimization-based methods aim to effectively adapt f_θ to new few-shot tasks. Meta-learning has been a popular method especially with the introduction of MAML and its variants [12], [13], [13]–[15] which aim to learn a good initialization of the neural network such that it can quickly adapt to new tasks with few gradient steps. On the other hand, distance-based approaches aim to learn a good feature

extractor [16]. The main idea is to find an appropriate feature space where each image is well represented for transfer tasks.

In the distance-based category, the feature extractor can be trained in two different ways. The first one relies on episodic training where the idea is to reproduce the same conditions of the few-shot adaptation phase during the pre-training of the feature extractor [17], [18]. Episodic training is not specific to distance-based approaches as it is also used in meta-learning. The second way to train a feature extractor is to use a standard cross-entropy loss. This is usually referred to as transfer learning, which has been successful in recent years and largely adopted due to its competitive performance compared to episodic training, while being relatively simple to implement [6], [8], [19], [20].

A common approach in transfer learning is to rely on data augmentation, self-supervision and manifold mixup [21] to learn a robust feature extractor for better generalization ability. In this paper, we adopt the transfer learning strategy with the above-mentioned training techniques.

Various few-shot classifiers have been proposed in the literature such as fine tuning a multi-layer perceptron with a cross-entropy loss [16], which has been criticized for being biased in few-shot regimes [22]. Other distance-based approaches such as using a nearest class mean classifier [8] or an earth distance metric using optimal transport [23] have also been studied. We adhere to the nearest class mean (NCM) approach due to its well established performance and its simplicity.

Definition 2 (Nearest class mean classifier). A nearest class mean classifier C_{NCM} is the optimal classifier when class-conditional densities follow a Gaussian distribution with equal isotropic covariance and uniform prior across classes [24]:

$$p(\mathbf{z} | y = c) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_c, \sigma^2 \mathbf{I}), \quad (1)$$

where $\boldsymbol{\mu}_c$ is the center of class $c \in \mathcal{Y}$, σ is the standard deviation and \mathbf{I} the identity matrix. The classification of a new sample \mathbf{z} is performed according to:

$$C_{NCM}(\mathbf{z}) = \arg \min_{c \in \mathcal{Y}} \|\mathbf{z} - \boldsymbol{\mu}_c\|_2. \quad (2)$$

In practice we estimate the class centers from the training data \mathcal{D} using the empirical average of each class.

Once the class centers are estimated, there are a maximum of $(n - 1)$ dimensions of interest in the considered Euclidean space, which correspond to directions between class centers. Remaining ones can be disregarded, as they produce contributions that are orthogonal to the axes between class centers, and thus contribute equally to any distance computation in Equation (2). Projecting the data onto this lower dimensional subspace is preferable in few-shot classification as it allows to work with lower dimensions [25]. Such a projection can be performed using for example a QR decomposition [26] to reduce the dimension of our data. Indeed, as shown in [25], a subspace of dimension $(n - 1)$ does not impact the boundary decisions of a NCM classifier.

¹The code to reproduce the results of our experiments is available in the following link: <https://github.com/ybendou/fs-generalization>.

B. Predicting Generalization

Predicting generalization is one of the most important topics of ML. It was brought into focus by [27], who asked the question of how one can measure the generalization from training data, showing that neural networks can easily fit randomly labeled data with high accuracy but with low generalization capabilities. We can define the problem of predicting generalization as follows [28]:

Definition 3 (Predicting generalization). *Given an underlying joint probability distribution $p_{\mathcal{X},\mathcal{Y}}$ and a classifier C trained on \mathcal{D} , the goal of predicting generalization is to estimate the error defined as:*

$$\mathcal{R}(C) = \mathbb{E}_{(\mathbf{x},y) \sim p_{\mathcal{X},\mathcal{Y}}} [\mathbb{1}_{C(\mathbf{x}) \neq y}]. \quad (3)$$

Many works have been proposed on generalization of a neural network trained on standard and large training datasets such as ImageNet [29] and on finding a measure which highly correlates with the accuracy instead of predicting the accuracy itself. In that regard, the Kendall’s rank-correlation coefficient is commonly used [11], since usually the ordering between different models is what matters for applications such as finding the best hyper-parameters or the best architecture in the field of neural architecture search.

The proposed methods in the literature can be summarized into few different families. The first one is PAC-Bayes methods where the generalization behavior of a model is described by probably approximately correct (PAC) bounds [30]; these methods often provide an upper-bound on the generalization error and the results are often restricted to a small set of models (e.g., no depth variations). The second family is norm-based methods, which analyze the neural network weights. These methods have shown to perform poorly [11]. The last family of methods aims to analyze the intermediate representation of the training data in the feature space. Many methods have been proposed in this line of work such as using the Davies-Bouldin Index [31] which is a clustering measure of the training data. Another approach [32] measures the distance of the training data to the decision boundaries in multiple intermediate features. Note that the main focus of these methods is to predict the generalization of a model trained for a certain task where large training data is available. Predicting generalization when working with few labeled samples has mainly been addressed when using meta-learning methods [33], [34]. The closest work to ours is [35], where some of the strategies for predicting generalization mentioned before have been tested for few-shot classification using transfer learning.

Differently to previously mentioned works, in this paper we aim at deriving a statistical model of the class-conditional densities in the feature space and to use this model to estimate the generalization error. As we will demonstrate in the experiments, the proposed methodology can outperform the previously mentioned ones in few-shot settings.

III. METHODOLOGY

A. Statistical model

As previously mentioned, the first step of our proposed methodology consists in proposing a statistical model for class-conditional densities in the feature space.

Let us assume that each class follows a Gaussian distribution with a uniform prior across the classes, i.e, $p(y_{\mathbf{z}} = c) = \frac{1}{n}, \forall c \in \mathcal{Y}$, where n is the cardinal of \mathcal{Y} , which is reasonable to assume in a few-shot setting [25]. The conditional densities are defined as:

$$p(\mathbf{z} | y = c) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c), \quad (4)$$

where $\boldsymbol{\Sigma}_c$ is the covariance matrix of class c .

Our hypothesis stands from the fact that given a well pre-trained feature extractor, each class in the feature space should follow a multivariate Gaussian distribution centered around a class center. This assumption has been largely adopted in the few-shot literature [36]–[39]. Furthermore, the performance we obtain through our experimental results in Section IV show that this model fits our data.

Predicting generalization can be defined as predicting the probability of error from the classifier C . Let $R_c = \{\mathbf{z} | C(\mathbf{z}) = c\}$ be the decision region for class c using the classifier C and $R = \cup_{c \in \mathcal{Y}} R_c$. The theoretical error of our problem is defined by the sum of integrals:

$$P_e = \sum_{c \in \mathcal{Y}} \int_{R \setminus R_c} p(\mathbf{z} | y = c) p(y = c) d\mathbf{z}. \quad (5)$$

B. Analytical insight

A closed form solution of Equation 5 when $n > 2$ is often intractable. Therefore, in this section we focus on the case of binary classification. We derive an analytical expression for P_e and propose a statistical bound for its estimate.

For the case of a binary classifier of isotropic Gaussian data with equal covariance $\boldsymbol{\Sigma}_c = \sigma^2 \mathbf{I}_d, \forall c$, where d is the dimension of \mathbf{z} , and class centers $\boldsymbol{\mu}_a$ and $\boldsymbol{\mu}_b$, the probability error has a closed form which only depends on the distance between the class centers $r = \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\|_2$ and the shared standard deviation:

$$P_e = 1 - \phi\left(\frac{r}{2\sigma}\right), \quad (6)$$

where ϕ is the cumulative distribution function of $\mathcal{N}(0, 1)$.

To estimate P_e , we typically estimate r and σ using i.i.d samples such that $\hat{r} = \|\hat{\boldsymbol{\mu}}_a - \hat{\boldsymbol{\mu}}_b\|_2$, where $\hat{\boldsymbol{\mu}}$ is the empirical mean estimate.

a) *Statistical bound of the probability error:* Under this analytical form we derive in the univariate case a statistical bound for $\hat{P}_e = 1 - \phi\left(\frac{\hat{r}}{2\hat{\sigma}}\right)$.

Lemma 1. *Given two classes represented by two univariate isotropic distributions $\mathcal{N}_a(\mu_a, \sigma^2)$ and $\mathcal{N}_b(\mu_b, \sigma^2)$ with shared and known standard deviations, the true probability of error P_e is bounded by:*

$$P\left(\left|P_e - \hat{P}_e\right| \leq \frac{\phi'(0)}{\sqrt{2k}} \left| \phi^{-1}\left(1 - \frac{\alpha}{2}\right) \right|\right) \geq 1 - \alpha, \quad (7)$$

with a probability $1 - \alpha$, where \hat{P}_e is the probability error estimate from a sequence of k i.i.d random variables from the two distributions with empirical mean estimates $\hat{\mu}_a$ and $\hat{\mu}_b$, ϕ is the cumulative distribution function of $\mathcal{N}(0, 1)$, ϕ^{-1} its inverse and ϕ' its derivative, with $\phi'(0) = \frac{1}{\sqrt{2\pi}}$.

The proof of Lemma 1 can be found in the Appendix. The bound for the true probability error is thus $\mathcal{O}(\frac{1}{\sqrt{k}})$. We compare this behavior on real data in the Appendix. Furthermore, this result holds for multivariate isotropic distributions as we can simply project on the axis between the two class centers to work with univariate distributions. In the remaining sections, we consider multivariate distributions.

b) *Bias of the naive distance estimator:* Estimating the probability error depends on estimating the distance between class centers. The naive approach for estimating distances is usually performed by estimating the means of each distribution using the empirical mean estimate and computing $\hat{r} = \|\hat{\mu}_a - \hat{\mu}_b\|_2$ to which we refer to as the naive estimator. However, this estimation of the distance between class centers is biased.

Lemma 2. Let $(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k)$ and $(\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k)$ be two sequences of i.i.d random variables drawn from their respective multivariate probability distributions p_a and p_b assumed independent with finite expected values μ_a and μ_b and finite second order moment with covariance matrices Σ_a and Σ_b . Let \hat{r} be the naive estimator for the distances using $\hat{\mu}_a = \frac{1}{k} \sum_{i=1}^k \mathbf{a}_i$ and $\hat{\mu}_b = \frac{1}{k} \sum_{i=1}^k \mathbf{b}_i$ the mean estimator of each of the two sequences, then:

$$\mathbb{E}_{\substack{\mathbf{a} \sim p_a \\ \mathbf{b} \sim p_b}}(\hat{r}^2) - r^2 = \frac{\text{Tr}(\Sigma_a + \Sigma_b)}{k}. \quad (8)$$

The proof of Lemma 2 is included in the Appendix. This Lemma holds for any two independent distributions. The bias is a function of the noise and the number of samples k . In the case of two isotropic multivariate distributions $\mathcal{N}_a(\mu_a, \sigma_a^2 \mathbf{I}_d)$ and $\mathcal{N}_b(\mu_b, \sigma_b^2 \mathbf{I}_d)$, the bias becomes

$$\mathbb{E}_{\substack{\mathbf{a} \sim \mathcal{N}_a \\ \mathbf{b} \sim \mathcal{N}_b}}(\hat{r}^2) - r^2 = \frac{d(\sigma_a^2 + \sigma_b^2)}{k}. \quad (9)$$

Interestingly, the bias scales with the ratio $\frac{d}{k}$. We include this bias reduction step as part of our numerical approach. The correction is performed in the original high dimensional space. We show the extent of this bias in the experimental results in section IV. Furthermore, to show the importance of this step, in our experiments we provide results with and without correcting the bias.

C. Numerical insight

Computing P_e analytically in equation 5 for $n > 2$ is hard. However, in practice we can approximate P_e through a Monte Carlo method. For each class, we draw a large number of data points from the Gaussian distributions that we fitted to the few-shot training dataset. This way, we artificially enrich our dataset and compute the classifier's decisions on a *virtual validation set*. We sample in the reduced subspace

with $n - 1$ dimensions. Note that the decision regions in the lower dimensional and in the original subspace are the same, as already stated in Section II.

Similar to the binary case, we need to take into account that the distances between the estimated class centers are positively biased. This leads to underestimated P_e (see Appendix). Correcting this bias is key to accurately estimating P_e . In fact, for a distance-based classifier, the absolute positioning of the class centers do not affect its decisions. In order to perform the sampling, we generate a set of points which respect the new estimated distances using the Nonmetric Multidimensional Scaling algorithm [40].

Estimating P_e by sampling means that there are no restrictions for choosing the covariance of the data. In section IV we compare the performance for different covariance matrices, i.e.: 1) using the identity matrix, 2) using a shared isotropic covariance matrix across classes, 3) using isotropic covariance matrix per class, 4) using full covariance matrix per class. In section IV we run an experiment to validate our choice.

IV. EXPERIMENTS

A. Datasets and implementation details

We run our experiments on three standardized few-shot vision classification benchmarks:

- Mini-ImageNet [1]: A dataset for in-domain few-shot learning extracted from ImageNet. The dataset is divided into three splits with disjoint classes. The first split with 64 classes (base classes) is used to train the feature extractor, the second split contains 16 classes for validation and the third split contains 20 classes (novel classes) for test. Each class contains 600 images.
- Tiered-ImageNet [41]: Another subset of ImageNet for in-domain classification with 351 base classes, 97 classes for validation and 160 novel classes. Each class contains around 1300 samples. This dataset has classes with a hierarchy structure.
- Meta-dataset [42]: The classical Meta-dataset benchmark consists of 10 different datasets from which we use 4 for our experiments. The first dataset is ImageNet for in-domain classification which is also split into three disjoint subsets for training the feature extractor, validating and testing. The remaining 3 datasets are for cross-domain classification (VGG Flower, CUB-200-2011 and Describable Textures). Note that in cross-domain the training split of ImageNet is used to train the feature extractor and the 3 remaining datasets are used for testing.

Each of these benchmarks contains a large number of samples (≥ 500 samples per class) from which we artificially and randomly sample, for each dataset, 10^3 few-shot classification problems to run our experiments.

We use the standard training procedure proposed in [21] to pre-train a ResNet-18 architecture [43]. This architecture along the training procedure has been vastly used in the few-shot classification literature with feature vectors of 512 dimensions [6], [16].

B. Estimating the first and second order moments

First, we conduct an experiment to validate the choice of the statistical model: particularly, which covariance matrix should be selected for the Gaussian model. We generate 10^3 few-shot problems. For each of these, we estimate the mean and covariance matrix of each class under different models and compare it with the true distribution of each class from all the available samples in the original dataset (typically ≥ 500 per class). We use the Kullback-Leibler (KL) divergence as a metric and vary the number of samples per class. The evolution of the error per model is given in Figure 2.

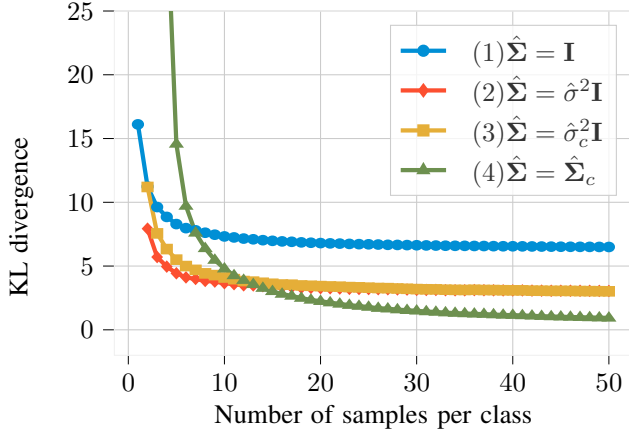


Fig. 2: Average KL divergence between a Gaussian distribution fitted with a limited number of samples and the closest Gaussian approximation using a larger number of samples. We average over 10^3 5-class few-shot classification problems from the test set of Meta-dataset ImageNet.

Figure 2 shows a trade-off between the expressive power of each model and its capacity to overfit. As the number of samples per class increases (≥ 14), the error of the model with more parameters (model 4) becomes lower than simpler models and hence it becomes a better choice. For very few samples, it is desirable to choose a simpler model with fewer parameters (model 2). Therefore, in our experiments we vary the model used depending on the number of samples. We use a shared isotropic covariance matrix across classes (model 2) for $k \leq (n-1)^2$ which is 16 for 5 class classification tasks and is roughly where the intersection between the model 2 and 4 is in Figure 2. For more samples, we use a free covariance matrix (model 4).

C. Unbiased estimator for the distances

In order to quantify the bias of the naive estimator for the distances from Lemma 2, we generate synthetic data from two isotropic distributions $\mathcal{N}(\mu_1, \sigma^2 \mathbf{I}_d)$ and $\mathcal{N}(\mu_2, \sigma^2 \mathbf{I}_d)$ under different settings of Signal-to-Noise ratios $\text{SNR}_{\text{DB}} = 10 \log_{10} \left(\frac{r}{\sqrt{2}\sigma} \right)$, where $r = \|\mu_1 - \mu_2\|_2$. Figure 3 shows the bias of the naive estimator for the distance and the proposed unbiased estimator $\hat{r}_{\text{unbiased}}^2 = \hat{r}_{\text{naive}}^2 - \frac{2d\hat{\sigma}^2}{k}$, where k is the number of available samples and $d = 512$, the

number of dimensions of the original features in our data. For a low SNR, the bias of the naive estimator becomes very important. On the other hand, the proposed estimator is unbiased. When estimating means from few data points, the estimated distance can be seen as the sum of the distance between the true means and the added distance in the orthogonal directions to the axis between the true means due to noise.

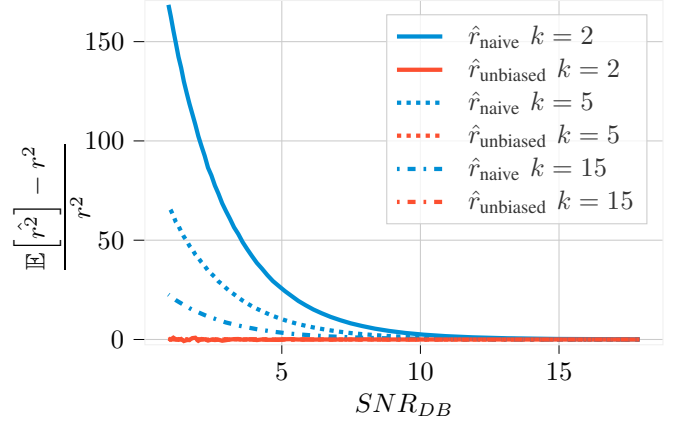


Fig. 3: Bias of the naive estimator. Here, the task is to estimate the distance between the center of two normal distributions. The true generative parameters are denoted r and σ . We normalize the bias by the true distance and plot it against the SNR for different sample sizes k . The proposed estimator is unbiased while the naive estimator has a large bias especially in low SNR regimes with low samples.

D. Performance in predicting generalization

In this section, we report the performance of our method under different settings. For each dataset, we run 10^3 few-shot classification problems between 5 classes which is the standard number of classes in the standardized few-shot benchmarks [8]. For each run, we predict the accuracy $(1 - \hat{P}_e)$ and compare it to the real accuracy obtained using the large number of samples available from the original datasets. We compare our method to two other approaches. The first one is the leaving-one-out cross-validation from the few available samples. The second method is the one used by [35] which computes the Davies-Bouldin (DB) Index, a clustering measure of inner-class and outer-class variance. For a fair comparison, we use the validation split of the original datasets to train a linear regression for the DB Index method and apply it on the few-shot tasks. For the cross-domain datasets, we use the validation split of ImageNet.

Our main metric is the Mean Absolute Percentage Error defined as: $\text{MAPE}(\hat{P}_e, P_e) = \frac{|P_e - \hat{P}_e|}{1 - P_e}$ averaged over 10^3 problems. The results of our experiments are reported in Figure 4. For each dataset², we plot the MAPE of the different

²Some datasets have a limited number of samples per class (VGG-Flower and CUB-200-2011). In order to ensure that we have enough samples to measure the ground-truth accuracy, we plot less than 50 samples for these datasets.

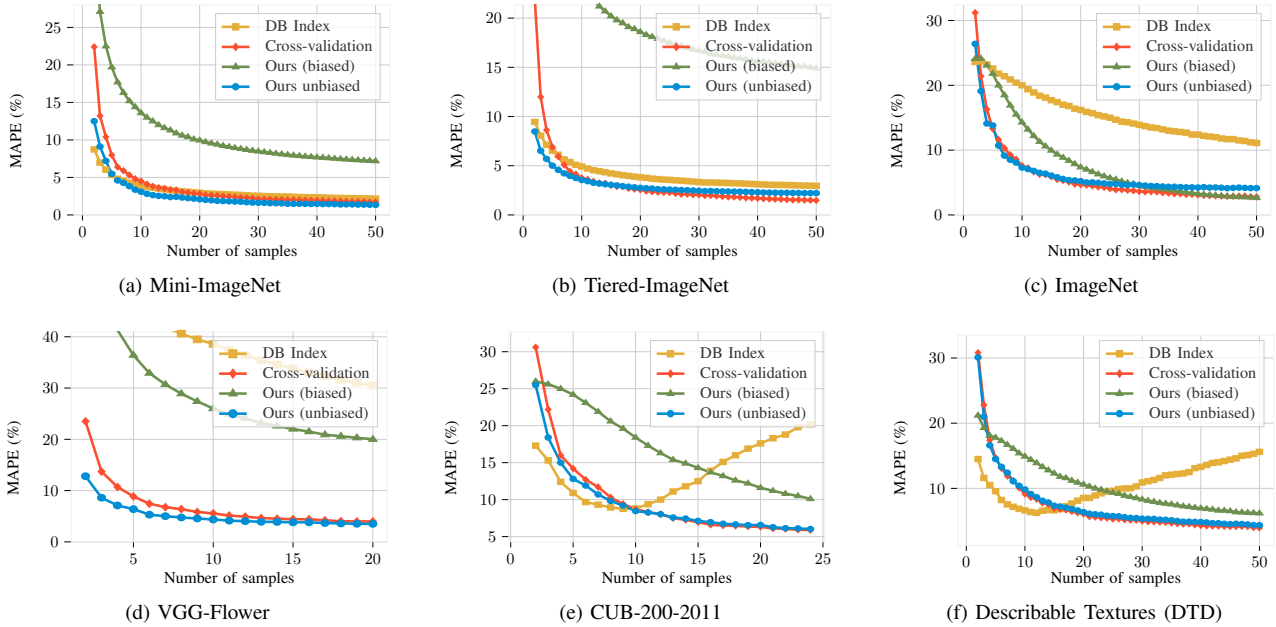


Fig. 4: Prediction error of different generalization predictors over 10^3 few-shot classification problems with 5 classes. Figures (a,b,c) are in-domain datasets and Figures (d,e,f) are cross-domain datasets. We plot the Mean-Absolute-Percentage Error against the number of samples and compare our method (unbiased) to cross-validation or Davies-Bouldin Index.

approaches against the number of samples per class. We observe that our proposed method (Ours unbiased) outperforms cross-validation when very few labeled samples are available ($k \leq 10$). The performance of cross-validation becomes more competitive when more samples are available, which is expected given the robustness of cross-validation with large data. Furthermore, the DB Index method is only competitive on Mini-ImageNet and tends to collapse in cross-domain or harder datasets. Moreover, our method without the bias correction (Ours biased) does not yield good results, showing the importance of the unbiased estimator for the distances.

Furthermore, our method is more efficient at predicting generalization. For example on Tiered-ImageNet, when the cross-validation method would require 17 samples to estimate generalization at a certain performance, our method would only require 10 samples to match it. On other datasets such as ImageNet and DTD, the gap is smaller, however on average our method does not perform worse than the cross-validation on the tested datasets.

Figure 5 shows a scatter plot of the true accuracy and the predicted accuracy of each method for few-shot problems sampled from ImageNet. The predictions from our method are aligned with the ground truth accuracies and are less scattered than the cross-validation approach.

E. Predicting generalization as a binary classification problem

Let us consider now that we want to classify a problem into two classes: hard or easy. To artificially create such classes, we use a threshold of 85% accuracy on unseen samples. We

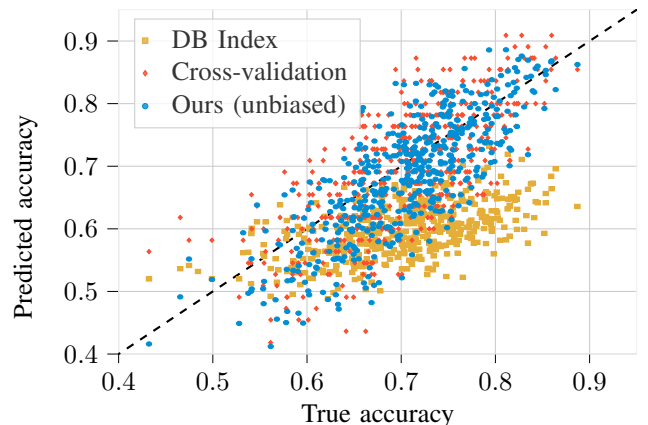


Fig. 5: Scatter plot of 5-class few-shot problems with 10 samples per class from ImageNet. Each point represents a different problem with a true ground-truth accuracy plotted against the predicted accuracy from the different methods.

want to investigate the ability of the proposed methodology to better classify between the two classes compared to DB Index or cross-validation. We thus compute the ROC curve when generating 10^3 few-shot problems from Mini-ImageNet. The results are depicted in Figure 6, where we see that the gap between the proposed methodology and available alternatives is even higher than in Figure 4. This is because we focus here on high accuracy, where our model typically reaches the best estimation of error.

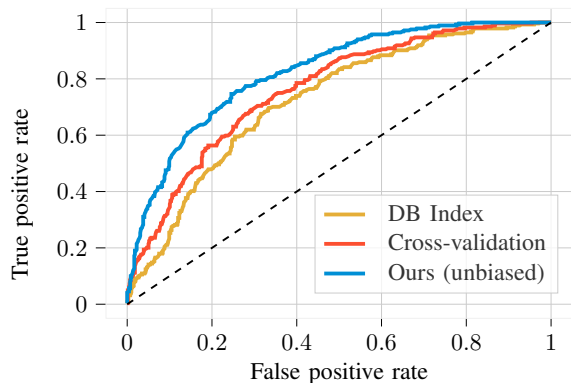


Fig. 6: Receiver operating characteristic (ROC) curve after binarizing the prediction of generalization for few-shot problems from Mini-ImageNet.

V. DISCUSSION AND LIMITATIONS

A first observation is that the DB Index performs poorly in our experiments. We believe that this is due to the large domain gap, causing a mismatch in the learned parameters of the linear regression. We can observe this behaviour in Figure 5 where the DB Index predictions are misaligned with the ground-truth accuracies. We have experimented with different learned functions using polynomial functions with different degrees, the linear regression yielded the best results. On the other hand, our method and the cross-validation do not suffer from this behaviour.

Furthermore, when predicting generalization from few annotated samples, all the existing techniques, except for cross-validation, are relying on a measure of clustering which depends on computing the distances and the variances of the classes, these methods could benefit from the bias correction step especially when working with few samples. The unbiased estimator is not restricted to our method. Other methods relying on the neural network gradients during training [44] or the analysis of the function space defined by the network [28] require training the feature extractor on the few-shot task which needs a large number of samples. In case of binary classification, our method can be seen as a clustering score similar to the DB Index on which we apply a Gaussian kernel. However, the analogy only holds for binary classification. While clustering scores are either a global measure for the data distribution (DB Index) or could be used to compute pairwise probability errors between classes, our method has the advantage of estimating the class-conditional densities to compute the overlap between all the classes.

Predicting error depends on the accuracy of the few-shot problems. Hard few-shot problems have a low SNR which makes the estimation of accuracy more difficult. On the other hand, good separation between the classes in the latent space leads to a better estimation of generalization. This explains the better performance on in-domain datasets compared to cross-domain, as well as the results of Figure 6. We run an experiment in Figure 7 to show the effect of SNR on the

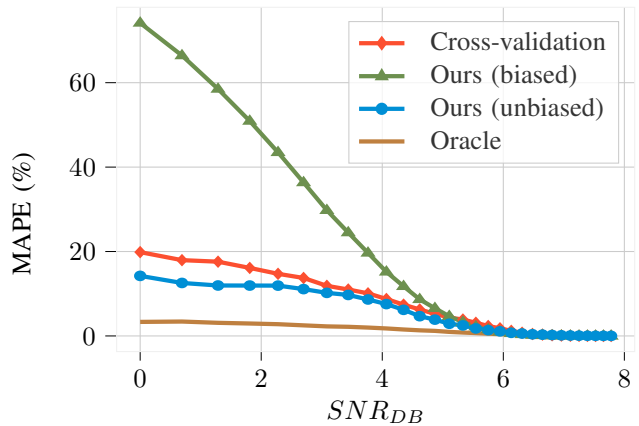


Fig. 7: Prediction error on 10^4 binary 10-shot classification problems generated from isotropic Gaussian distributions with different SNR values. The oracle model is using our approach with the true parameter distributions.

performance of both our method and cross-validation with 10 samples per class. The performance drastically drops in low SNR regimes especially for our method without correcting for the bias. This is expected since the bias is higher when the noise is high. We also include the performance of an oracle model, which is using our method with the true parameters of underlying distribution of the data. As the SNR value gets lower, the separation between the Gaussian distributions becomes more difficult which also impacts the estimation of the accuracy.

Finally, although the proposed estimator is unbiased, it still suffers from a large variance. We show this in Figure 8 where the variance of the estimator is slightly higher than the naive estimator. The variance of both estimators are particularly high in low SNR regimes.

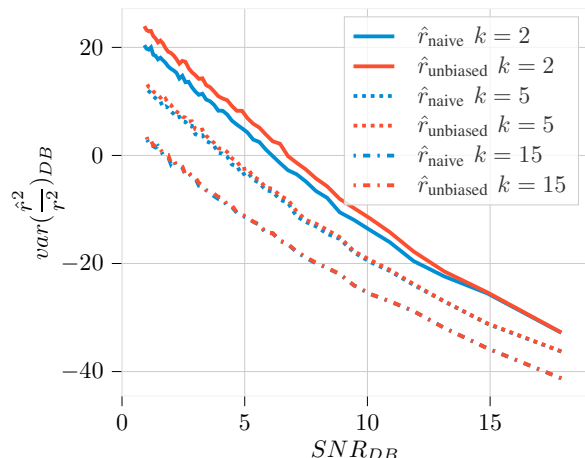


Fig. 8: Variance of both estimators against different SNR values in log-log scale. We normalize the estimates by the real distances. The unbiased estimator has a slightly higher variance than the naive estimator. The variance of the two estimators are asymptotically of $\mathcal{O}(\frac{1}{\text{SNR}})$.

VI. CONCLUSION

We propose a model based approach to estimate the generalization capability of few-shot classifiers and provide statistical performance bounds. Our method outperforms the leave-one-out cross-validation and the Davis-Bouldin score-based estimator for different Signal-to-Noise regimes and a small number of labeled samples. This is especially important in the few-shot context. Our method strongly relies on unbiased estimates of the inter-class distances, which is a key contribution of this paper. Note that our method can be generalized to transfer-based few-shot learners with any distance-based classifier. Although we improve upon existing methods, we think that it opens up interesting new directions for further research.

REFERENCES

- [1] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, “Matching networks for one shot learning,” *Advances in neural information processing systems*, vol. 29, 2016.
- [2] I. Ziko, J. Dolz, E. Granger, and I. B. Ayed, “Laplacian regularized few-shot learning,” in *International conference on machine learning*. PMLR, 2020, pp. 11 660–11 670.
- [3] M. N. Rizve, S. Khan, F. S. Khan, and M. Shah, “Exploring complementary strengths of invariant and equivariant representations for few-shot learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 836–10 846.
- [4] Y. Hu, S. Pateux, and V. Gripon, “Squeezing backbone feature distributions to the max for efficient few-shot learning,” *Algorithms*, vol. 15, no. 5, p. 147, 2022.
- [5] P. Bateni, J. Barber, J.-W. van de Meent, and F. Wood, “Enhancing few-shot image classification with unlabelled examples,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 2796–2805.
- [6] Y. Bendou, Y. Hu, R. Lafargue, G. Lioi, B. Pasdeloup, S. Pateux, and V. Gripon, “Easy: Ensemble augmented-shot y-shaped learning: State-of-the-art few-shot classification with simple ingredients,” 2022. [Online]. Available: <http://arxiv.org/abs/2201.09699>
- [7] Y. Tian, Y. Wang, D. Krishnan, J. B. Tenenbaum, and P. Isola, “Rethinking few-shot image classification: a good embedding is all you need?” in *European Conference on Computer Vision*. Springer, 2020, pp. 266–282.
- [8] Y. Wang, W.-L. Chao, K. Q. Weinberger, and L. van der Maaten, “Simpleshot: Revisiting nearest-neighbor classification for few-shot learning,” *arXiv preprint arXiv:1911.04623*, 2019.
- [9] Y. Bengio and Y. Grandvalet, “No unbiased estimator of the variance of k-fold cross-validation,” *Advances in Neural Information Processing Systems*, vol. 16, 2003.
- [10] Q. F. Gronau and E.-J. Wagenmakers, “Limitations of bayesian leave-one-out cross-validation for model selection,” *Computational brain & behavior*, vol. 2, no. 1, pp. 1–11, 2019.
- [11] Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio, “Fantastic generalization measures and where to find them,” *arXiv preprint arXiv:1912.02178*, 2019.
- [12] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” *International Conference on Machine Learning*, pp. 1126–1135, 2017.
- [13] J. Liu, F. Chao, and C.-M. Lin, “Task augmentation by rotating for meta-learning,” *arXiv preprint arXiv:2003.00804*, 2020.
- [14] K. Lee, S. Maji, A. Ravichandran, and S. Soatto, “Meta-learning with differentiable convex optimization,” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10 657–10 665, 2019.
- [15] N. Fei, Z. Lu, T. Xiang, and S. Huang, “Melr: Meta-learning via modeling episode-level relationships for few-shot learning,” *International Conference on Learning Representations*, 2020.
- [16] G. S. Dhillon, P. Chaudhari, A. Ravichandran, and S. Soatto, “A baseline for few-shot image classification,” *arXiv preprint arXiv:1909.02729*, 2019.
- [17] J. Snell, K. Swersky, and R. Zemel, “Prototypical networks for few-shot learning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [18] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, “Matching networks for one shot learning,” *Advances in neural information processing systems*, vol. 29, 2016.
- [19] X. Luo, L. Wei, L. Wen, J. Yang, L. Xie, Z. Xu, and Q. Tian, “Rectifying the shortcut learning of background for few-shot learning,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [20] J. Ma, H. Xie, G. Han, S.-F. Chang, A. Galstyan, and W. Abd-Almageed, “Partner-assisted learning for few-shot image classification,” *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10 573–10 582, 2021.
- [21] P. Mangla, N. Kumari, A. Sinha, M. Singh, B. Krishnamurthy, and V. N. Balasubramanian, “Charting the right manifold: Manifold mixup for few-shot learning,” *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2218–2227, 2020.
- [22] S. Ghaffari, E. Saleh, D. Forsyth, and Y.-X. Wang, “On the importance of firth bias reduction in few-shot classification,” in *International Conference on Learning Representations*, 2022.
- [23] C. Zhang, Y. Cai, G. Lin, and C. Shen, “Deepemd: Few-shot image classification with differentiable earth mover’s distance and structured classifiers,” *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12 203–12 213, 2020.
- [24] K. P. Murphy, *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [25] Y. Hu, S. Pateux, and V. Gripon, “Adaptive dimension reduction and variational inference for transductive few-shot classification,” *arXiv preprint arXiv:2209.08527*, 2022.
- [26] W. Gander, “Algorithms for the qr decomposition,” *Res. Rep.*, vol. 80, no. 02, pp. 1251–1268, 1980.
- [27] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding deep learning requires rethinking generalization,” *arXiv preprint arXiv:1611.03530*, 2016.
- [28] G. Ortiz-Jiménez, S.-M. Moosavi-Dezfooli, and P. Frossard, “What can linearized neural networks actually say about generalization?” *Advances in Neural Information Processing Systems*, vol. 34, pp. 8998–9010, 2021.
- [29] Y. Jiang, P. Natekar, M. Sharma, S. K. Aithal, D. Kashyap, N. Subramanyam, C. Lassance, D. M. Roy, G. K. Dziugaite, S. Gunasekar *et al.*, “Methods and analysis of the first competition in predicting generalization of deep learning,” in *NeurIPS 2020 Competition and Demonstration Track*. PMLR, 2021, pp. 170–190.
- [30] N. Ding, X. Chen, T. Levinboim, S. Goodman, and R. Soricut, “Bridging the gap between practice and pac-bayes theory in few-shot meta-learning,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 29 506–29 516, 2021.
- [31] P. Natekar and M. Sharma, “Representation based complexity measures for predicting generalization in deep learning,” *arXiv preprint arXiv:2012.02775*, 2020.
- [32] Y. Jiang, D. Krishnan, H. Mobahi, and S. Bengio, “Predicting the generalization gap in deep networks with margin distributions,” *arXiv preprint arXiv:1810.00113*, 2018.
- [33] A. Farid and A. Majumdar, “Generalization bounds for meta-learning via pac-bayes and uniform stability,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 2173–2186, 2021.
- [34] Q. Chen, C. Shui, and M. Marchand, “Generalization bounds for meta-learning: An information-theoretic analysis,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 25 878–25 890, 2021.
- [35] M. Bontonou, L. Béthune, and V. Gripon, “Predicting the accuracy of a few-shot classifier,” *arXiv preprint arXiv:2007.04238*, 2020.
- [36] T. Chobola, D. Vašata, and P. Kordík, “Transfer learning based few-shot classification using optimal transport mapping from preprocessed latent space of backbone neural network,” in *AAAI Workshop on Meta-Learning and MetaDL Challenge*. PMLR, 2021, pp. 29–37.
- [37] Y. Hu, V. Gripon, and S. Pateux, “Leveraging the feature distribution in transfer-based few-shot learning,” in *International Conference on Artificial Neural Networks*. Springer, 2021, pp. 487–499.
- [38] J. Xu and H. Le, “Generating representative samples for few-shot classification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 9003–9013.
- [39] T. Cao, M. Law, and S. Fidler, “A theoretical analysis of the number of shots in few-shot learning,” *arXiv preprint arXiv:1909.11722*, 2019.

- [40] J. B. Kruskal, "Nonmetric multidimensional scaling: a numerical method," *Psychometrika*, vol. 29, no. 2, pp. 115–129, 1964.
- [41] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. B. Tenenbaum, H. Larochelle, and R. S. Zemel, "Meta-learning for semi-supervised few-shot classification," *arXiv preprint arXiv:1803.00676*, 2018.
- [42] E. Triantafillou, T. Zhu, V. Dumoulin, P. Lamblin, U. Evci, K. Xu, R. Goroshin, C. Gelada, K. Swersky, P.-A. Manzagol *et al.*, "Meta-dataset: A dataset of datasets for learning to learn from few examples," *arXiv preprint arXiv:1903.03096*, 2019.
- [43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [44] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang, "On large-batch training for deep learning: Generalization gap and sharp minima," *arXiv preprint arXiv:1609.04836*, 2016.

A Statistical Model for Predicting Generalization in Few-Shot Classification

Supplementary Materials

Yassir Bendou¹, Vincent Gripon¹, Bastien Pasdeloup¹, Giulia Lioi¹, Lukas Mauch², Stefan Uhlich², Fabien Cardinaux², Ghouthi Boukli Hacene^{2,3} and Javier Alonso Garcia²

¹IMT Atlantique, Lab-STICC, UMR CNRS 6285, F-29238 Brest, France

²Sony Europe, R&D Center, Stuttgart Laboratory 1, Germany

³Mila, Montréal, Canada

¹name.surname@imt-atlantique.fr, ²name.surname@sony.com

I. STATISTICAL BOUND FOR THE PROBABILITY ESTIMATE

Lemma 1. *Given two classes represented by two univariate isotropic distributions $\mathcal{N}_a(\mu_a, \sigma^2)$ and $\mathcal{N}_b(\mu_b, \sigma^2)$ with shared and known standard deviations, the true probability of error P_e is bounded by:*

$$P\left(\left|P_e - \hat{P}_e\right| \leq \frac{\phi'(0)}{\sqrt{2k}} \left| \phi^{-1}\left(1 - \frac{\alpha}{2}\right) \right|\right) \geq 1 - \alpha, \quad (1)$$

with a probability $1 - \alpha$, where \hat{P}_e is the probability error estimate from a sequence of k i.i.d random variables from the two distributions with empirical mean estimates $\hat{\mu}_a$ and $\hat{\mu}_b$, ϕ is the cumulative distribution function of $\mathcal{N}(0, 1)$, ϕ^{-1} its inverse and ϕ' its derivative, with $\phi'(0) = \frac{1}{\sqrt{2\pi}}$.

a) *Proof of Lemma 1:* Let (a_1, a_2, \dots, a_k) and (b_1, b_2, \dots, b_k) be two sequences of i.i.d random variables drawn from two respective univariate normal distributions $\mathcal{N}(\mu_a, \sigma^2)$ and $\mathcal{N}(\mu_b, \sigma^2)$. Let $\hat{\mu}_a = \frac{1}{k} \sum_{i=1}^k a_i$ and $\hat{\mu}_b = \frac{1}{k} \sum_{i=1}^k b_i$ be the empirical mean estimates of each of the two sequences. Then $\hat{\mu}_a - \hat{\mu}_b \sim \mathcal{N}(\mu_a - \mu_b, \frac{2\sigma^2}{k})$. Therefore, $\hat{t} = \frac{(\hat{\mu}_a - \hat{\mu}_b) - (\mu_a - \mu_b)}{\sigma\sqrt{2/k}} \sim \mathcal{N}(0, 1)$. With a probability $1 - \alpha$, the confidence bound for the standard normal is:

$$P\left(\left|\frac{(\hat{\mu}_a - \hat{\mu}_b) - (\mu_a - \mu_b)}{\sigma\sqrt{2/k}}\right| \leq \left|\phi^{-1}(1 - \alpha/2)\right|\right) = 1 - \alpha.$$

We also know that :

$$\left|P_e - \hat{P}_e\right| = \left|\phi\left(\frac{|\mu_a - \mu_b|}{2\sigma}\right) - \phi\left(\frac{|\hat{\mu}_a - \hat{\mu}_b|}{2\sigma}\right)\right|.$$

Given that ϕ is a differentiable function then it is Lipschitz continuous and since $\forall x \in \mathbb{R} : \phi'(x) \leq \phi'(0)$, then:

$$\forall (x, y) \in \mathbb{R}^2 : |\phi(x) - \phi(y)| \leq \phi'(0)|x - y|.$$

Using this property and the triangular inequality, we get:

$$\begin{aligned} \left|P_e - \hat{P}_e\right| &\leq \phi'(0) \left| \frac{|\mu_a - \mu_b|}{2\sigma} - \frac{|\hat{\mu}_a - \hat{\mu}_b|}{2\sigma} \right| \\ \frac{\sqrt{2k} |P_e - \hat{P}_e|}{\phi'(0)} &\leq \left| \frac{(\mu_a - \mu_b) - (\hat{\mu}_a - \hat{\mu}_b)}{\sigma\sqrt{2/k}} \right|. \end{aligned}$$

Hence with a probability $1 - \alpha$:

$$\frac{\sqrt{2k} |P_e - \hat{P}_e|}{\phi'(0)} \leq \left| \frac{(\mu_a - \mu_b) - (\hat{\mu}_a - \hat{\mu}_b)}{\sigma\sqrt{2/k}} \right| \leq \left| \phi^{-1}(1 - \alpha/2) \right|.$$

Finally, since density functions are positive, we can rewrite with a probability $1 - \alpha$:

$$P\left(\left|P_e - \hat{P}_e\right| \leq \frac{\phi'(0)}{\sqrt{2k}} \left|\phi^{-1}(1 - \alpha/2)\right|\right) \geq 1 - \alpha.$$

□

Figure 1 shows the $\frac{1}{\sqrt{k}}$ behaviour of the confidence bound on real data. We generate 10^3 few-shot problems and plot the average difference between the true probability error and the estimated probability error using $\frac{1}{|P_e - \hat{P}_e|^2}$ against the number of samples k . We run our experiment with both binary classification and 5-class problems and we fit a linear regression for each of the configurations with respective coefficients of determination of $R^2 = 0.9979$ and $R^2 = 0.9981$.

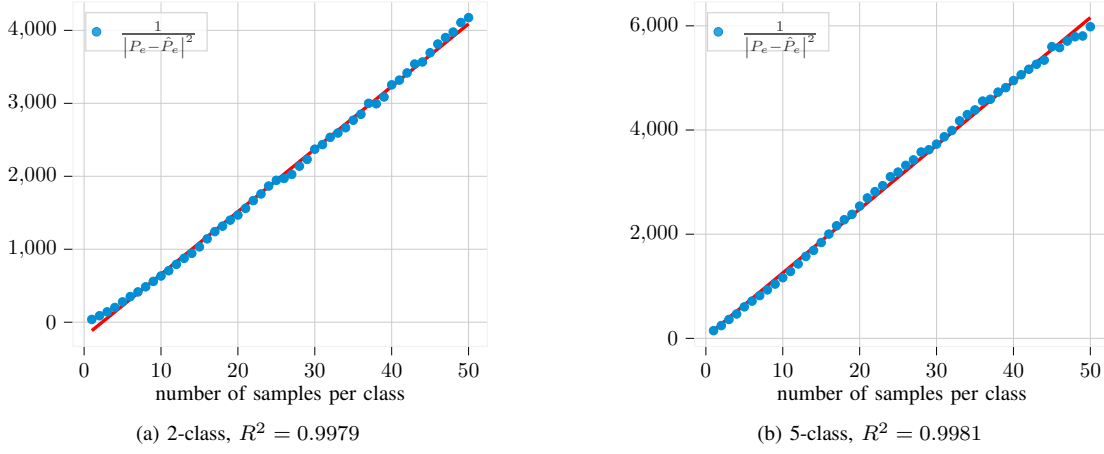


Fig. 1: Inverse quadratic difference between the estimate of the probability of error and the true probability of error. For each point, we average 10^3 few-shot problems from Mini-ImageNet for both binary and 5-class problems.

II. BIAS OF THE NAIVE DISTANCE ESTIMATE

Lemma 2. Let $(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k)$ and $(\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_k)$ be two sequences of i.i.d random variables drawn from their respective multivariate probability distributions p_a and p_b assumed independent with finite expected values $\boldsymbol{\mu}_a$ and $\boldsymbol{\mu}_b$ and finite second order moment with covariance matrices $\boldsymbol{\Sigma}_a$ and $\boldsymbol{\Sigma}_b$. Let \hat{r} be the naive estimator for the distances using $\hat{\boldsymbol{\mu}}_a = \frac{1}{k} \sum_{i=1}^k \mathbf{a}_i$ and $\hat{\boldsymbol{\mu}}_b = \frac{1}{k} \sum_{i=1}^k \mathbf{b}_i$ the mean estimator of each of the two sequences, then:

$$\mathbb{E}_{\substack{\mathbf{a} \sim p_a \\ \mathbf{b} \sim p_b}}(\hat{r}^2) - r^2 = \frac{\text{Tr}(\boldsymbol{\Sigma}_a + \boldsymbol{\Sigma}_b)}{k}. \quad (2)$$

a) *Proof of Lemma 2::* For this proof, we exploit the independence of the variables and the second order definition: $\mathbb{E}_{\mathbf{a} \sim p_a}[\mathbf{a}\mathbf{a}^T] = \boldsymbol{\Sigma}_a + \mathbb{E}_{\mathbf{a} \sim p_a}[\mathbf{a}]\mathbb{E}_{\mathbf{a} \sim p_a}[\mathbf{a}]^T$.

$$\begin{aligned} \mathbb{E}_{\substack{\mathbf{a} \sim p_a \\ \mathbf{b} \sim p_b}}\left[\|\hat{\boldsymbol{\mu}}_a - \hat{\boldsymbol{\mu}}_b\|_2^2\right] &= \frac{1}{k^2} \mathbb{E}_{\substack{\mathbf{a} \sim p_a \\ \mathbf{b} \sim p_b}}\left[\left\|\sum_{i=1}^k (\mathbf{a}_i - \mathbf{b}_i)\right\|_2^2\right] \\ &= \frac{1}{k^2} \mathbb{E}_{\substack{\mathbf{a} \sim p_a \\ \mathbf{b} \sim p_b}}\left[\text{Tr}\left(\sum_{i=1}^k \sum_{j=1}^k (\mathbf{a}_i - \mathbf{b}_i)(\mathbf{a}_j - \mathbf{b}_j)^T\right)\right] \\ &= \frac{1}{k^2} \sum_{\substack{i,j \\ i \neq j}}^k \text{Tr}\left(\mathbb{E}_{\substack{\mathbf{a}_i \sim p_a \\ \mathbf{b}_i \sim p_b}}(\mathbf{a}_i - \mathbf{b}_i)\mathbb{E}_{\substack{\mathbf{a}_j \sim p_a \\ \mathbf{b}_j \sim p_b}}(\mathbf{a}_j - \mathbf{b}_j)^T\right) + \frac{1}{k^2} \sum_{i=1}^k \text{Tr}\left(\mathbb{E}_{\substack{\mathbf{a}_i \sim p_a \\ \mathbf{b}_i \sim p_b}}[(\mathbf{a}_i - \mathbf{b}_i)(\mathbf{a}_i - \mathbf{b}_i)^T]\right) \\ &= \frac{k^2 - k}{k^2} \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\|_2^2 + \frac{1}{k} \text{Tr}(\mathbb{E}_{\mathbf{a} \sim p_a}[\mathbf{a}\mathbf{a}^T]) + \frac{1}{k} \text{Tr}(\mathbb{E}_{\mathbf{b} \sim p_b}[\mathbf{b}\mathbf{b}^T]) - \frac{2}{k} \text{Tr}\left(\mathbb{E}_{\substack{\mathbf{a} \sim p_a \\ \mathbf{b} \sim p_b}}[\mathbf{a}\mathbf{b}^T]\right) \\ &= \left(1 - \frac{1}{k}\right) \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\|_2^2 + \frac{1}{k} \text{Tr}\left(\boldsymbol{\Sigma}_a + \boldsymbol{\Sigma}_b + \mathbb{E}_{\mathbf{a} \sim p_a}[\mathbf{a}]\mathbb{E}_{\mathbf{a} \sim p_a}[\mathbf{a}]^T + \mathbb{E}_{\mathbf{b} \sim p_b}[\mathbf{b}]\mathbb{E}_{\mathbf{b} \sim p_b}[\mathbf{b}]^T - \mathbb{E}_{\mathbf{a} \sim p_a}[\mathbf{a}]\mathbb{E}_{\mathbf{b} \sim p_b}[\mathbf{b}]^T\right) \\ &= \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\|_2^2 + \frac{1}{k} \text{Tr}(\boldsymbol{\Sigma}_a + \boldsymbol{\Sigma}_b). \end{aligned}$$

III. ADDITIONAL EXPERIMENTS FOR PREDICTING GENERALIZATION

We experiment with different number of classes and show the performance of our method compared to DB-Index and cross-validation. In Figure 2 we show the results on 10-class few-shot problems. The trend is similar for different configurations of classes. We did not include DTD in the 10-class problems since the original dataset only has 7 classes. Regarding the performance, our method performs better on average than the alternatives. We only use a matrix with shared isotropic covariance given that the dimensionality of the projected space is 9 for 10-class problems.

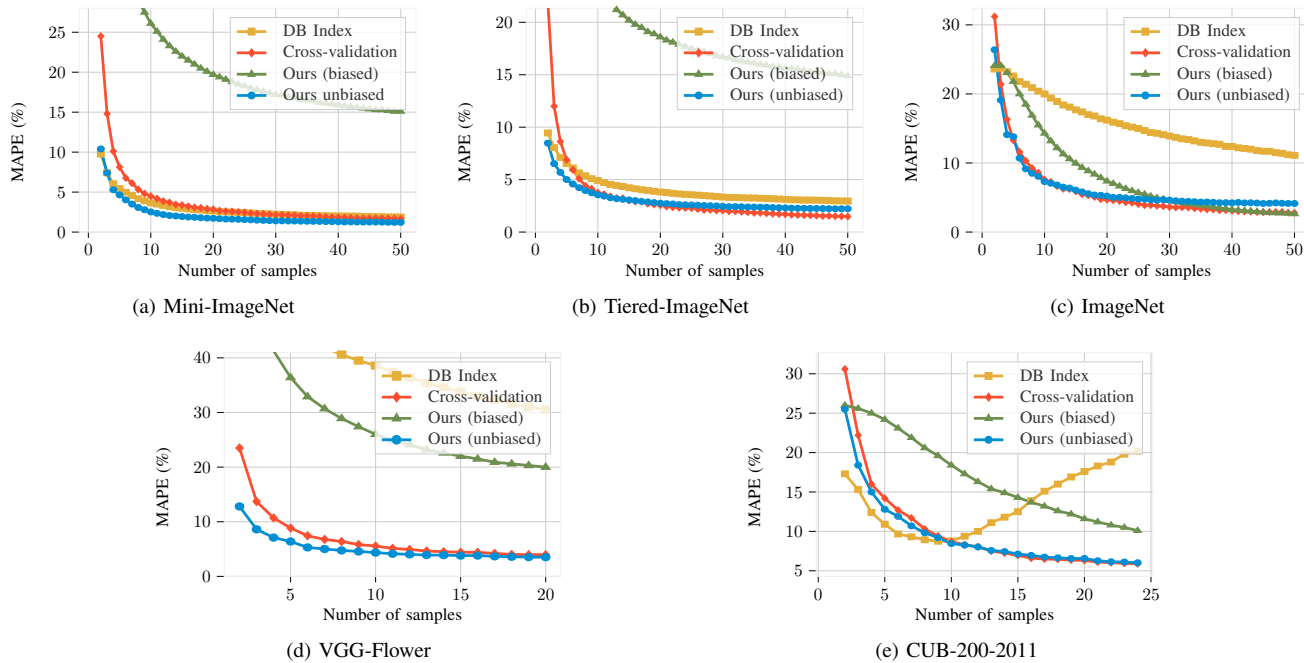


Fig. 2: Prediction error of different generalization predictors over 10^3 few-shot classification problems with 10 classes. Figures (a,b,c) are in-domain datasets and Figures (d,e) are cross-domain datasets. We plot the Mean-Absolute-Percentage Error against the number of samples and compare our method to cross-validation and Davies–Bouldin Index.

In Figure 3 we include an additional scatter plot for Mini-ImageNet. The predictions from our method are aligned with the ground truth accuracies and are less scattered than the cross-validation approach.

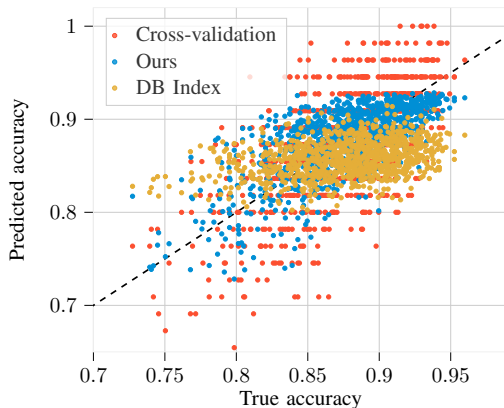


Fig. 3: Scatter plot of 5-class few-shot problems with 10 samples per class from Mini-ImageNet. Each point represents a different problem with a true ground-truth accuracy plotted against the predicted accuracy from the different methods.