

# Few-Shot Image Recognition with Manifolds

Debasmit  $Das^{(\boxtimes)}$ , J. H. Moon, and C. S. George Lee

School of Electrical and Computer Engineering, Purdue University, West Lafayette, IN, USA {das35,moon92,csglee}@purdue.edu

Abstract. In this paper, we extend the traditional few-shot learning (FSL) problem to the situation when the source-domain data is not accessible but only high-level information in the form of class prototypes is available. This limited information setup for the FSL problem deserves much attention due to its implication of privacy-preserving inaccessibility to the source-domain data but it has rarely been addressed before. Because of limited training data, we propose a non-parametric approach to this FSL problem by assuming that all the class prototypes are structurally arranged on a manifold. Accordingly, we estimate the novel-class prototype locations by projecting the few-shot samples onto the average of the subspaces on which the surrounding classes lie. During classification, we again exploit the structural arrangement of the categories by inducing a Markov chain on the graph constructed with the class prototypes. This manifold distance obtained using the Markov chain is expected to produce better results compared to a traditional nearestneighbor-based Euclidean distance. To evaluate our proposed framework, we have tested it on two image datasets - the large-scale ImageNet and the small-scale but fine-grained CUB-200. We have also studied parameter sensitivity to better understand our framework.

# 1 Introduction

Deep learning has produced breakthrough in many areas like computer vision [7, 9], speech recognition [1], natural language processing [4] etc., mainly due to the availability of lots of labeled data, complex neural network architectures and efficient training procedures. Even though these deep learning models are trained on large labeled datasets, they still fail to generalize to new classes or environments. Humans, on the other hand, can quickly recognize new objects from very few samples. They do that by using their previously obtained *knowledge* and apply it to new situations. This difference in the way machines and humans learn provides motivation to carry out research on few-shot learning (FSL). Accordingly, most few-shot learning methods strive for a transfer-learning approach where they extract useful transferable knowledge from data-abundant base classes and use it to recognize data-starved novel classes.

Most previous methods in FSL assumed that abundant labeled data is available across the base categories from which a robust and generalizable knowledge

representation can be learned. However, in certain situations, it is difficult to have access to all the labeled data of these base categories due to privacy restrictions and/or inefficiency in maintaining such a large database. Hence, alternatively class exemplars or prototypes can be retained. These prototypes summarize the class information by averaging over the data samples without revealing sensitive information about the data. For example, the prototype of the dog category can be the arithmetic mean of all the dog sample features. Previous work on *hypothe*sis transfer learning (HTL) [8, 10, 15] assumed access to base class models, where recognition performance would depend on the choice of these models. As a result, this HTL-based setting does not allow for fair comparison. On the other hand, our proposed setting involving base-class prototypes allow for fair comparison where performance depends only on the data and the proposed transfer learning approach. This proposed FSL setting is depicted in Fig. 1.

Previous FSL methods cannot be applied to this proposed restrictive setting. This is because they assume that lots of labeled data are available from the source domain. Consequently, they use neural-network-based parametric approaches. The neural-network-based approaches can be categorized depending on the type of transferable knowledge extracted from the base categories and encoded in the neural network architecture: (a) *Metric-learning* methods [13, 14, 16] learn a metric space; (b) *Metalearning* methods [2, 5, 12] learn the learning procedure; (c) *Generative* methods [6, 11, 17] generate data for the novel classes. However, the



Fig. 1. In this FSL setting, the base-class prototypes are known but not the novel-class proto-types. The spread of the classes (dashed boundaries) are also unknown.

neural-network-based parametric models might severely overfit if we only have limited data in the form of class prototypes available from the source domain. Hence, in order to solve this restrictive FSL problem, it is natural to seek a non-parametric approach.

In this paper, we address this restricted FSL setting by formulating it as a case of *ill-sampling*. As depicted in Fig. 1, the correct locations of the base category prototypes are known but that of the novel categories are unknown. This is because the few-shot data from a novel category might be sampled from the periphery of the class distribution. These non-representative samples when used for classification will result in poor recognition performance. Therefore, a non-parametric-based prior is used to produce a biased estimate of the novelclass prototype location.

For the non-parametric-based prior, we find inspiration from the idea [3] that data samples from one class lie on a low-dimensional subspace. Therefore, we can consider all the classes as a collection of piece-wise linear subspaces. This set of subspaces can be considered as an approximation of a non-linear manifold close to which the class-prototypes lie. This manifold serves as a prior to estimate the location of the novel-class prototype. The subspace near the

novel-class prototype is found by calculating the mean of the subspaces on which the nearby base classes lie. The subspace on which the nearby base classes lie is again found using their nearest neighbors as shown in Fig. 2. Finally, the novelclass sample can be projected onto the mean subspace to obtain the novel-class prototype.

Once the novel-class prototypes are estimated, one can use the nearestneighbor (NN) approach to assign a test sample to a class based on the Euclidean distance to all the prototypes. However, the estimation procedure for the novelclass prototype might still be error prone. Hence, there is a need to exploit the structural arrangement of the manifold containing all the classes to assign a class to a test sample. The neighboring class locations can provide better estimate of class-prototype distances. This can be achieved by constructing a graph using all the class prototypes and then using equilibrium probability of an induced absorbing Markov-chain process to output the most probable class. Finally, to validate the proposed approach, we perform experiments and analyses on this framework to set a benchmark for future research.



**Fig. 2.** The surrounding subspaces  $S_1$ ,  $S_2$  and  $S_3$  for a novel-class sample  $\mathbf{x}_n$  are found by using its Nearest Neighbors (NN) and the NNs of its NNs. Their mean can be calculated to obtain the subspace  $\overline{S}$  on which the novel-class sample  $\mathbf{x}_n$  is projected to obtain  $\mathbf{c}_p$ . The weighted average  $\mathbf{c}_d$  of the nearby prototypes are also used to obtain the novel-class prototype.

### 2 Proposed Approach

In this section, we describe the proposed framework, which consists of two steps estimating novel-class prototypes using the manifold approach and classification using the Markov chain method.

#### 2.1 Estimating Novel-Class Prototypes

Consider that we have access to the base-category prototypes collected in the form of a matrix  $\mathbf{C} \in \mathbb{R}^{n_b \times d}$ , where  $n_b$  is the number of base prototypes and d

is the dimensionality of the feature space. Let the one-shot sample from a novel class be  $\mathbf{x}_n \in \mathbb{R}^d$ . Our goal is to estimate the prototype location  $\mathbf{c}_n \in \mathbb{R}^d$  of the novel class. Our assumption is that all the base and novel-class prototypes, i.e., all rows of  $\mathbf{C}$  and  $\mathbf{c}_n$  lie close to a non-linear manifold. Since the mathematical expression of the manifold is unknown, we express it as a collection of piecewise linear subspaces. First, we find the surrounding classes of the novel class by finding r nearest-neighboring prototypes of the novel-class sample  $\mathbf{x}_n$ . Let these nearest-neighboring prototypes, we find q neighboring prototypes. These new prototypes can be expressed as  $\mathbf{c}_{nij} \in \mathbb{R}^d$  for  $j \in \{1, 2, ..., r\}$ . For each of the r neighboring prototypes, we find q neighboring prototypes. These new prototypes can be expressed as  $\mathbf{c}_{nij} \in \mathbb{R}^d$  for  $j \in \{1, 2, ..., q\}$  and  $i \in \{1, 2, ..., r\}$ . Hence,  $\mathbf{c}_{nij}$  represents the  $j^{th}$  nearest neighbor of the  $i^{th}$  nearest neighbor of the novel class sample  $\mathbf{x}_n$ . To represent the non-linear manifold, we form r linear subspaces using the r nearest neighbors of the novel sample as well as the q nearest neighbors of each of the r prototypes. The linear subspace  $S_i$  corresponding to the  $i^{th}$  nearest neighbor of  $\mathbf{x}_n$  is represented as a column

space such that  $S_i \equiv [\mathbf{c}_{ni} \vdots \mathbf{c}_{ni1} \vdots \mathbf{c}_{ni2} \vdots ... \vdots \mathbf{c}_{niq}]$ . The linear subspace can be orthonormalized to obtain  $S_i^{\perp}$  and the operation can be repeated for all the rnearest neighbors. The net result is r linear subspaces with dimensionality (q+1)surrounding the novel-class sample  $\mathbf{x}_n$ . In the example in Fig. 2, we chose r = 3and q = 3. These linear subspaces represent linearized localized versions of the non-linear manifold on which the class prototypes lie. The subspace on which the novel-class prototype lies close to can be found by averaging these surrounding r subspaces  $S_i^{\perp}$  for  $i \in \{1, 2, ..., r\}$ . For finding the average of these orthonormal subspaces, we use the concept of Grassmann manifold.

A Grassmann manifold  $\mathcal{G}(n,l)$  for n,l > 0 is the topological space composed of all *l*-dimensional linear subspaces embedded in an *n*-dimensional Euclidean space. A point on the Grassmann manifold is represented as an  $n \times l$  orthonormal matrix **S** whose columns span the corresponding linear subspace  $\mathcal{S}$ . It is represented as:  $\mathcal{G}(n,l) = \{\text{span}(\mathbf{S}): \mathbf{S} \in \mathbb{R}^{n \times l}, \mathbf{S}^T \mathbf{S} = \mathbf{I}_l\}$ , where  $\mathbf{I}_l$  is a  $l \times l$ dimensional identity matrix and superscript T indicates matrix transpose.

Following the definition, the r orthonormal subspaces  $S_i^{\perp}$  for  $i \in \{1, 2, ..., r\}$  are points lying on a  $\mathcal{G}(d, q + 1)$  Grassmann manifold. The average of these points on the Grassmann manifold will represent the linear subspace to which the novel-class prototype lies close to. The average of these points is found using the *extrinsic mean*. For a set of points on the Grassmann manifold  $\mathcal{G}(d, q+1)$ , the extrinsic mean is the point that minimizes the Frobenius-norm-squared difference of the projections of the points onto the space of (q + 1) ranked  $d \times d$  matrices. Therefore, the optimization problem for finding the extrinsic mean  $\mathbf{\bar{S}}$  is

$$\underset{\mathbf{S}^T \mathbf{S} = \mathbf{I}}{\operatorname{argmin}} \sum_{i=1}^{r} d(\mathbf{S}_i, \mathbf{S})^2, \quad \text{where} \quad d(\mathbf{S}_i, \mathbf{S}) = \frac{||\mathbf{S}\mathbf{S}^T - \mathbf{S}_i\mathbf{S}_i^T||_{\mathcal{F}}}{\sqrt{2}}.$$
 (1)

Here  $|| \cdot ||_{\mathcal{F}}$  is the Frobenius norm. Let  $\mathbf{S}^*$  be the solution to the optimization problem (1), which can be found using eigenvalue decomposition.  $\mathbf{S}^*$  is the spanning matrix of the extrinsic mean of the surrounding subspaces  $\mathcal{S}_i^{\perp}$ 's. Setting the extrinsic mean  $\mathbf{\bar{S}} = \mathbf{S}^*$ , we project the novel-class sample  $\mathbf{x}_n$  onto the

subspace spanned by the matrix  $\bar{\mathbf{S}}$ . The projected point  $\mathbf{c}_p$  can be obtained as  $\mathbf{c}_p = \bar{\mathbf{S}}\bar{\mathbf{S}}^+\mathbf{x}_n$ , where  $\bar{\mathbf{S}}^+ = (\bar{\mathbf{S}}^T\bar{\mathbf{S}})^{-1}\bar{\mathbf{S}}^T$ . Superscripts -1 and + indicate matrix inverse and matrix pseudo-inverse, respectively.

We also consider the direct contribution of the surrounding class prototypes into estimating the novel-class prototype. If  $\mathbf{C}_r \in \mathbb{R}^{r \times d}$  consists of the *r* nearest neighbors of the novel-class sample  $\mathbf{x}_n$ , then their contribution  $\mathbf{c}_d$  to the novelclass prototype location can be found using the equation  $\mathbf{c}_d = \mathbf{C}_r^T \mathbf{p}_d$ , where  $\mathbf{p}_d \in \mathbb{R}^r$  is the probability vector formed by carrying out the exponential mapping of the Euclidean distances of the class prototypes to  $\mathbf{x}_n$ , followed by normalization. Hence, the contributions  $\mathbf{x}_n$ ,  $\mathbf{c}_p$  and  $\mathbf{c}_d$  can be used to estimate the novel-class prototype location  $\mathbf{c}_n$  as

$$\mathbf{c}_n = \alpha_2 [\alpha_1 \mathbf{x}_n + (1 - \alpha_1) \mathbf{c}_p] + (1 - \alpha_2) \mathbf{c}_d, \tag{2}$$

where  $\alpha_1, \alpha_2 \in [0, 1]$  are scalar weights that are manually set. In case  $\mathbf{x}_n$  is very close to the novel-class prototype,  $\alpha_1 = \alpha_2 \approx 1$  will produce optimal classification performance.

#### 2.2 Classification Using Absorbing Markov Chain

Once the class prototype locations of the novel classes are known, the structural arrangement of the prototypes of both the base and novel classes are again used to recognize a test sample to obtain a more informed decision about the classification. The structural arrangement of the classes is represented using a k'-nearestneighbor (k'-NN) graph, where each node represents a class prototype. The k'-NN graph formulation allows nodes to only be connected to its k'-NN nodes. The weights between the nodes are defined using the exponential of the negative Euclidean distances. Upon defining this graph, an absorbing Markov-chain process Markov chain corresponds to a node in the transition probability from a state i to a st where  $w_{il}$  is the weight connecting nodes i a process, there are two kinds of states - transtate and the absorbing state have self-tranrespectively. This suggests that a random w transient node for the next step but for the



**Fig. 3.** Possible transitions are shown with directed arrows. The transition probability from a state *i* to a state *j* is  $p_{ij}$ . The transient state and the absorbing state have self-transition probabilities  $p_{ii}$  as 0 and 1, respectively.

graph, an absorbing Markov-chain process is induced on it. Each state of the Markov chain corresponds to a node in the graph and therefore a category. The transition probability from a state *i* to a state *j* is found as  $p_{ij} = w_{ij} / \sum_l w_{il}$ , where  $w_{il}$  is the weight connecting nodes *i* and *l*. In the absorbing Markov chain process, there are two kinds of states - *transient* and *absorbing*. The transient state and the absorbing state have self-transition probabilities  $p_{ii}$  as 0 and 1, respectively. This suggests that a random walker on a graph cannot stay on the transient node for the next step but for the absorbing node it will stay there forever. An example of an absorbing Markov chain is given in Fig. 3, where the arrows represent the possible transitions from one node to another. Overall, the Markov chain is represented using the transition matrix **P**, which models the

dynamics of the process. Using  $\mathbf{P}$ , the Markov-chain equations are described as follows:

$$\mathbf{u}^{t+1} = \mathbf{u}^{t} \mathbf{P}, \quad \text{where} \quad \mathbf{P} = \begin{bmatrix} \mathbf{T}_{n_{t} \times n_{t}} \ \mathbf{A}_{n_{t} \times n_{a}} \\ \mathbf{0}_{n_{a} \times n_{t}} \ \mathbf{I}_{n_{a} \times n_{a}} \end{bmatrix}, \mathbf{P}^{\infty} = \begin{bmatrix} \mathbf{0}_{n_{t} \times n_{t}} \ (\mathbf{I} - \mathbf{T})^{-1} \mathbf{A} \\ \mathbf{0}_{n_{a} \times n_{t}} \ \mathbf{I}_{n_{a} \times n_{a}} \end{bmatrix}.$$
(3)

 $\mathbf{u}^t$  and  $\mathbf{u}^{t+1}$  are the states of the process at instants t and t + 1, respectively, and they are represented as a probability vector over all the states.  $\mathbf{T}$  describes the transition probabilities from one transient state to another. A describes the transition probabilities from transient states to absorbing states.  $n_t$  and  $n_a$  are the number of transient and absorbing states, respectively, and the zero and identity matrices  $\mathbf{0}_{n_a \times n_t}$  and  $\mathbf{I}_{n_a \times n_a}$  imply that the process cannot leave the absorbing state. Our goal is to find the equilibrium state  $\mathbf{u}^t$  as  $t \to \infty$  for a given initial state  $\mathbf{u}^0$ . Accordingly,  $\mathbf{u}^{\infty} = \mathbf{u}^0 \mathbf{P}^m$  as  $m \to \infty$ . The closed-form solution of  $\mathbf{P}^m$  as  $m \to \infty$  is treated as  $\mathbf{P}^{\infty}$ . Using this formulation, the equilibrium state probabilities can only be distributed among the absorbing states with zero probabilities on the transient states.

The initial state  $\mathbf{u}^0$  is calculated using the Euclidean distances of the test sample to all the base- and novel-class prototypes and normalizing it to obtain a probability vector. Using the absorbing Markov-chain formulation in Eq. (3), we choose the novel categories and base categories as transient and absorbing states, respectively, to obtain the most probable base category from  $\mathbf{u}^{\infty}$ . Then, we interchange the order of absorbing and transient states to obtain the most probable novel category. In the final step, we apply one nearest neighbor approach on the test sample to choose the most probable class among the most probable base and novel categories obtained in the previous steps. The overall procedure from the novel-class prototype estimation to the Markov-chain-based prediction for a test sample is given in Algorithm 1. In case we have multiple samples for a novel class,  $\mathbf{x}_n$  is set as the mean of these samples.

### 3 Experiments and Discussions

#### 3.1 Dataset Description

To evaluate our proposed approach, we used two image recognition datasets – ImageNet and CUB-200. Originally, the ImageNet dataset consists of 21K categories of which we used 1000 for our experiments. These 1000 categories are accordingly split into base and novel classes. The CUB-200 is a fine-grained dataset that consists of 200 categories of different bird species. Of these 200 classes, we used a total of 150 of which 100 are base and 50 are novel classes. For both datasets, the image features used were the 2048-dimensional ResNet-101 [7].

#### 3.2 Effects of Varying the Number of Classes and Samples

In this section, we study how recognition performance is affected by the number of categories and the number of samples per category in the base and novel Algorithm 1: Proposed two-step FSL procedure using manifolds.

**Given:** Base category prototypes  $\mathbf{C} \in \mathbb{R}^{n_b \times d}$ . Novel class one-shot samples  $\mathbf{x}_n, n \in \{1, 2, \dots, n_{nov}\}$ , Test sample  $\mathbf{x}_{te}$ . **Parameters:**  $r, q, k', \alpha_1, \alpha_2$ **Goal:** Classify  $\mathbf{x}_{te}$  into one of the  $n_b + n_{nov}$  categories **Step 1** Estimate class prototype for each novel class for each novel class  $n \in \{1, 2, ..., n_{nov}\}$ Obtain r nearest base prototypes for  $\mathbf{x}_n$  to form  $\mathbf{C}_r \in \mathbb{R}^{r \times d}$ Obtain q nearest base prototypes for each of the r base prototypes Obtain orthonormal subspaces  $\mathcal{S}_i^{\perp}$  for  $i \in \{1, 2, ..., r\}$  using the q neighbors  $\bar{\mathcal{S}} \leftarrow \texttt{ExtrinsicManifoldMean}(\mathcal{S}_1^{\perp}, \mathcal{S}_2^{\perp}, ..., \mathcal{S}_r^{\perp})$ Project  $\mathbf{x}_n$  onto  $\bar{\mathcal{S}}$  to obtain  $\mathbf{c}_p$  followed by  $\mathbf{c}_d$ Obtain novel-class prototypes as  $\mathbf{c}_n \leftarrow \alpha_2(\alpha_1 \mathbf{x}_n + (1 - \alpha_1)\mathbf{c}_p) + (1 - \alpha_2)\mathbf{c}_d$ end for **Step 2** Predict class of test sample  $\mathbf{x}_{te}$ Construct k'-nearest-neighbor graph with  $n_b$  base prototypes and  $n_{nov}$  novel prototypes as nodes. Find initial probability vector  $\mathbf{u}_0$  using distance of  $\mathbf{x}_{te}$  to all the prototypes. Construct Markov chain and obtain most probable base class. Construct Markov chain and obtain most probable novel class. Use nearest neighbor to obtain the most probable class among the two.

datasets. For training purposes, we used the prototypes of the base categories and the few-shot samples from the novel categories. For evaluation purposes, test samples from both base and novel categories were used. Recognition performance is reported as class-wise averaged accuracy. This ensures that major classes do not dominate the performance and minor classes containing less number of samples are not ignored. It is noted that performing cross-validation is impossible since we do not have access to enough data to be held out as a validation set and therefore results were reported by fixing the hyper-parameters. For evaluation, we used the following models: (NA) The No-adaptation baseline which consists of just using nearest neighbor on the few-shot sample mean; (M1) It uses nearest neighbor on the estimated novel-class prototypes; (M2) It uses the Markov-chain-based manifold distance on the few-shot sample mean; (Oracle) It assumes access to novel-class prototypes and uses nearest neighbor for prediction; (M1+M2) involves computing the manifold distance on the estimated class prototypes.

For the first set of experiments, we used the ImageNet dataset with 800 base and 200 novel categories and studied the effect of changing the number of shots per novel category. The results were taken over 10 trials and reported in Table 1. From the results, it is seen that M1 improves the recognition performance over the no-adaptation baseline but the difference diminishes as the number of shots increases. This is because for the novel categories, the few-shot mean becomes closer to the prototype location as the number of shots increases. Also, the contribution of M2 over the baseline or over M1 is incremental. This

can be attributed to the fact that the ResNet-101 features are not trained using the manifold-based distance and there is a mis-match between the training and testing evaluation measures. The standard error reduces with the increasing number of shots because of reduced variance in the few-shot mean and eventually the estimated prototype over the trials. We repeated the same experiment for the CUB-200 dataset, the results of which are reported in Table 2. In this case, we have 100 base and 50 novel categories, all of which are fine-grained. As a result, the recognition performance is poorer compared to ImageNet, even though CUB-200 has lesser number of categories. Still, the observed recognition performance has a pattern similar to that of the ImageNet dataset. However, there is no reduction in the standard error with increasing shots. This can be attributed to larger overlap between the fine-grained classes of CUB-200. For the

**Table 1.** Average accuracy results over 10 trials of the ImageNet dataset with 800 base and 200 novel categories as the number of shots per novel category is changed. Standard error is shown in the parentheses. The hyper-parameter setting is  $r = 20, q = 20, k' = 3, \alpha_1 = 0.9, \alpha_2 = 0.7$ .

	1 shot	2 shot	5 shot	10 shot	20 shot
NA	$64.31 \ (0.05)$	$67.60\ (0.05)$	71.09(0.03)	72.24(0.02)	72.89(0.01)
M1	66.58(0.05)	69.67 (0.05)	71.62(0.03)	$72.31 \ (0.02)$	72.91(0.01)
M1+M2	$66.91 \ (0.05)$	69.88(0.05)	$72.05\ (0.03)$	72.72(0.02)	72.97(0.02)
M2	$65.21 \ (0.05)$	67.98(0.05)	$71.33\ (0.02)$	72.07(0.02)	72.60(0.01)
Oracle	73.3	73.3	73.3	73.3	73.3

**Table 2.** Accuracy results over 10 trials of the CUB-200 dataset with 100 base and 50 novel categories as the number of shots per novel category is changed. The hyperparameter setting is  $r = 20, q = 20, k' = 5, \alpha_1 = 0.5, \alpha_2 = 0.5$ .

	1 shot	2 shot	5 shot	10 shot	20 shot
NA	43.40 (0.12)	45.28(0.13)	51.19(0.18)	55.70 (0.20)	58.16(0.16)
M1	45.91(0.23)	48.80 (0.20)	51.90(0.16)	55.94(0.18)	57.63(0.14)
M1+M2	46.13(0.28)	49.01 (0.22)	52.13(0.11)	55.57(0.17)	58.67(0.12)
M2	43.81 (0.11)	45.86(0.15)	51.45(0.15)	55.91 (0.18)	58.31(0.14)
Oracle	60.51	60.51	60.51	60.51	60.51

next set of experiments, we considered the performance change on the ImageNet dataset for the 1-shot setting as the numbers of base and novel categories are varied. We considered two such scenarios. In the first case, the total number of categories was fixed at 1000 while the proportion of base categories was changed.

Table 3. Accuracy results over 10 trials of the ImageNet d.	ataset for the 1-shot setting
as the ratio of number of base categories to the total numb	per of categories is changed.
(x-b, y-n) implies x base and y novel categories.	

	0.1 (100-b, 900-n)	0.2 (200-b,800-n)	0.4 (400-b, 600-n)	0.6 (600-b, 400-n)
NA	38.29 (0.30)	39.59(0.29)	46.50 (0.18)	55.28 (0.11)
M1	47.70 (0.28)	49.65(0.30)	54.65(0.23)	$60.71 \ (0.12)$
M1+M2	47.89(0.28)	50.33(0.29)	55.13(0.23)	$61.34\ (0.11)$
M2	38.36(0.30)	39.68(0.29)	46.57(0.19)	$55.31 \ (0.11)$
Oracle	73.3	73.3	73.3	73.3

This setting considers less number of base categories compared to novel categories and it has rarely been studied in previous work. The results of this setting are reported in Table 3. From the table, it can be seen that M1 improves over NA by a large margin (9 points) especially when the number of base categories is very less (ratio of 0.1). This is alluded to our assumption that all the class prototypes have a structural arrangement on a manifold. Therefore, the use of this structure is especially beneficial in the few-class regime. However, the difference between M1 and NA decreases mainly due to more base categories and lesser difficult novel categories. After that, we considered the experimental setting where the total number of categories was varied but the proportion of base and novel categories was kept the same at 4:1, the results of which are reported in Table 4. The results show that the upper bound of the recognition performance, i.e., the oracle performance, decreases with an increase in the number of categories. This is because classification becomes more difficult as the number of categories increases. As expected, M1 performs better compared to NA and the contribution of M2 is incremental.

Our proposed few-shot learning setting is new and therefore we do not have previous work to compare and benchmark against. However, we can study whether our approach can improve existing relevant work. Prototypical networks [13] consider the mean of the few-shot samples to represent class prototypes. Therefore, there is a possibility of obtaining the class prototypes using our manifold-based approach and further improving the performance. Accordingly, we tested the contribution of M1 and M2 over prototypical networks on *mini-ImageNet*, which is a subset of the ImageNet dataset. The results are shown in Table 5. In the table, K-way N-shot implies that K novel categories are sampled per testing episode with N samples per category. From the results, it is clear that M1 improves the performance, however M2 declines it. This is mainly because of the discrepancy between the Euclidean distance metric used during training ProtoNets and the manifold-based distance metric used during testing.

#### 3.3 Parameter Sensitivity Studies

In this section, we study the effect of hyper-parameters on the recognition performance. We only report results of sensitivity with respect to r,  $\alpha_1$  and  $\alpha_2$  in

	50 (40-b, 10-n)	100 (80-b, 20-n)	200 (160-b, 40-n)	500 (400-b, 100-n)
NA	81.76 (0.54)	75.98(0.34)	70.25(0.16)	67.43 (0.10)
M1	86.10(0.45)	79.92(0.36)	73.50(0.16)	69.59 (0.08)
M1+M2	86.43 (0.47)	$80.61 \ (0.51)$	74.09(0.34)	$70.41 \ (0.08)$
M2	82.52(0.53)	75.48(0.39)	70.88(0.23)	67.44(0.10)
Oracle	92.20	86.76	80.85	76.87

**Table 4.** Accuracy results over 10 trials of the ImageNet dataset for the 1-shot setting as the total number of categories is changed but keeping the ratio of base categories to novel categories as 4:1.

Table 5. Few-shot classification accuracies on the miniImageNet dataset averaged over 600 test episodes for different ways and shots. 95% confidence intervals are shown in the parentheses.

	5-way 1-shot	5-way 5-shot	20-way 1-shot	20-way 5-shot
ProtoNet	47.21(0.69)	63.62(0.61)	20.51 (0.46)	35.20(0.59)
${\rm ProtoNet}{+}{\rm M1}$	48.79(0.51)	$65.67 \ (0.56)$	21.93(0.62)	$35.66\ (0.53)$
${\rm ProtoNet}{+}{\rm M2}$	41.36(0.47)	57.48(0.43)	$16.94\ (0.57)$	31.52(0.64)

Fig. 4. We found our recognition performance to be negligibly sensitive to q and k'. This suggests that the location of the novel-class prototype estimate only depends on the number of subspaces (r) rather than its dimensionality (q + 1). Similarly, the Markov-chain-based prediction does not depend on the number of nearest neighbors k' used for connecting the graph. In Fig. 4(a), the number of subspaces (r) is varied, keeping rest of the hyper-parameters the same. This is done for both the ImageNet (denoted as (I)) and the CUB-200 dataset (denoted as (C)). From the plot, it is seen that the performance increases as the number of subspaces increases but then decreases after reaching a peak. Initially, more linear subspaces in the neighborhood are required for estimating the local structure of the non-linear manifold but becomes irrelevant after more suppspaces.

Next, we studied the effects of  $\alpha_1$  and  $\alpha_2$  on the recognition performance for the 1- and 5-shot settings of the ImageNet and the CUB-200 datasets as reported in Figs. 4(b) and 4(c), respectively. Accordingly, we obtained an optimal performance when  $\alpha_1$  or  $\alpha_2$  is between 0 and 1. From Eq. (2), it suggests that the location of the novel-class prototype is within the space bounded by the fewshot sample mean ( $\mathbf{x}_n$ ), the subspace projection ( $\mathbf{c}_p$ ) and the weighted mean of the nearby class prototypes ( $\mathbf{c}_d$ ). For higher number of shots, the maxima seems to move towards the right; that is, closer to  $\alpha_1$ ,  $\alpha_2$  values of 1. This implies more contribution from the few-shot class mean as visible from Eq. (2). This is intuitive because as the number of shots increases, we expect the few-shot sample mean to converge to the class prototype. In fact, for shots of 10 and higher, we obtained the maxima at  $\alpha_1 = \alpha_2 = 0.9$  on both datasets. For low values of  $\alpha_1$ ,  $\alpha_2$  (less contribution of the few-shot sample mean), we observed a dip in



Fig. 4. (a) Effect of the number of subspaces r on recognition performance for both ImageNet (I) and CUB-200 (C). Effect of  $\alpha_1$  and  $\alpha_2$  on 1-shot and 5-shot recognition performance for (b) ImageNet and (c) CUB-200. Legends of (c) hold for (b) as well.  $\alpha_1$  in parenthesis suggests that  $\alpha_1$  is varied while  $\alpha_2 = 1$  and vice versa. All results are over 10 trials.

performance, even sometimes worse than the NA baseline. This suggests that the contribution of the few-shot mean is important in estimating the novel-class prototype. From the plot, we see that the setting  $\alpha_1 = 0$ ,  $\alpha_2 = 1$  produces much better performance as compared to  $\alpha_1 = 1$ ,  $\alpha_2 = 0$ . According to Eq. (2), it means that the contribution of the projection ( $\mathbf{c}_p$ ) is more important compared to contribution of nearby prototypes ( $\mathbf{c}_d$ ).

## 4 Conclusions

In this paper, we have introduced a new setting in few-shot learning that assumes access to only the base-class prototypes. To address this problem, we used the structural arrangement of the class prototypes on a manifold, firstly to estimate the novel-class prototypes and secondly to induce an absorbing Markov-chain for test-time prediction. From the experiments, it is evident that our proposed method improved over the no-adaptation baseline but there is still a lot of room for improvement to reach the oracle-level performance. Therefore, our results serve as a benchmark for future researchers to work upon. Acknowledgments. This work was supported in part by the National Science Foundation under Grant IIS-1813935. Any opinion, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation. We also gratefully acknowledge the support of NVIDIA Corporation for the donation of a TITAN XP GPU used for this research.

# References

- 1. Amodei, D., et al.: Deep speech 2: end-to-end speech recognition in English and Mandarin. In: International Conference on Machine Learning, pp. 173–182 (2016)
- Andrychowicz, M., et al.: Learning to learn by gradient descent by gradient descent. In: Advances in Neural Information Processing Systems, pp. 3981–3989 (2016)
- Basri, R., Jacobs, D.W.: Lambertian reflectance and linear subspaces. IEEE Trans. Pattern Anal. Mach. Intell. 25(2), 218–233 (2003)
- Cho, K., et al.: Learning phrase representations using RNN encoder-decoder for statistical machine translation. In: Proceedings of Conference on Empirical Methods in Natural Language Processing, pp. 1724–1734 (2014)
- Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. arXiv preprint arXiv:1703.03400 (2017)
- Hariharan, B., Girshick, R.: Low-shot visual recognition by shrinking and hallucinating features. In: Proceedings of IEEE International Conference on Computer Vision (ICCV), Venice, Italy (2017)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778 (2016)
- Jie, L., Tommasi, T., Caputo, B.: Multiclass transfer learning from unconstrained priors. In: Proceedings of IEEE International Conference on Computer Vision, pp. 1863–1870 (2011)
- Krizhevsky, A., Šutskever, I., Hinton, G.E.: ImageNet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems, pp. 1097–1105 (2012)
- Kuzborskij, I., Orabona, F., Caputo, B.: Scalable greedy algorithms for transfer learning. Comput. Vis. Image Under. 156, 174–185 (2017)
- Mehrotra, A., Dukkipati, A.: Generative adversarial residual pairwise networks for one shot learning. arXiv preprint arXiv:1703.08033 (2017)
- 12. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: International Conference on Learning Representations (2017)
- Snell, J., Swersky, K., Zemel, R.: Prototypical networks for few-shot learning. In: Advances in Neural Information Processing Systems, pp. 4080–4090 (2017)
- Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P.H., Hospedales, T.M.: Learning to compare: relation network for few-shot learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018)
- Tommasi, T., Orabona, F., Caputo, B.: Learning categories from few examples with multi model knowledge transfer. IEEE Trans. Pattern Anal. Mach. Intell. 36(5), 928–941 (2014)
- Vinyals, O., Blundell, C., Lillicrap, T., Kavukcuoglu, K., Wierstra, D.: Matching networks for one shot learning. In: Advances in Neural Information Processing Systems, pp. 3630–3638 (2016)
- 17. Wang, Y.X., Girshick, R., Herbert, M., Hariharan, B.: Low-shot learning from imaginary data. In: Computer Vision and Pattern Recognition (CVPR) (2018)