Δ-encoder: an effective sample synthesis method for few-shot object recognition

NIPS 2018

Eli Schwartz^{*1,2}, Leonid Karlinsky^{*1}, Joseph Shtok¹, Sivan Harary¹, Mattias Marder¹, Abhishek Kumar¹, Rogerio Feris¹, Raja Giryes² and Alex M. Bronstein³

¹IBM Research AI ²School of Electrical Engineering, Tel-Aviv University, Tel-Aviv, Israel ³Department of Computer Science, Technion, Haifa, Israel



01 Few-shot learning

02 The framework

03 Experiments

Few-shot learning

01	

.

Table 3: Summary of the datasets used in our experiments					
Dataset	Fine grained	Image size	Total # images	Seen classes	Unseen classes
<i>mini</i> ImageNet [42] CIFAR-100 [22] Caltech-256 Object Category [15] Caltech-UCSD Birds 200 (CUB) [47] Attribute Pascal & Yahoo (aPY) [9] Scene UNderstanding (SUN) [49] Animals with Attributes 2 (AWA2) [48]	× × × ✓ × × ✓	Medium Small Large Large Large Large Large	$60K \\ 60K \\ 30K \\ 12K \\ 14K \\ 14K \\ 37K$	80 80 156 150 20 645 40	20 20 50 50 12 72 10

With many instances

Only k instances for k-shot



Few-shot learning by metric learning

Learn an embedding into a metric space where some simple (usually L2) metric is then used to classify instances of new categories

Few-shot meta-learning (learning-to-learn)

The meta-learned classifiers are optimized to be easily fine-tuned on new fewshot tasks using the provided small training data.

Generative and augmentation-based few-shot approaches

Either generative models are trained to synthesize new data based on few examples, or additional examples are obtained by some other form of transfer learning from external data.

Intuition

02

The offset in feature space between a pair of same-class examples conveys information on a valid deformation. The deformation called " Δ "



Method

02

A model to capture the deformation between examples in the same class.

A model use the deformation AND an example X to generate a new example Y. The expected results is : The generated example Y is a deformation of the given example X.



Auto-encoder

02

Δ-encoder



Problem:

Z has the information of X. How can you guarantee the reconstructed X'_s is relied on Y?







Testing

N-way *k*-shot : draw *N* random unseen categories, and draw *k* random samples from each category.



we use our trained network to synthesize a total of 1024 samples per category based on those k examples. This is followed by training a simple linear N-class classifier over those $1024 \cdot N$ samples, and finally, the calculation of the few-shot classification accuracy on a set of M real (query) samples from the tested N categories.

Experiment



Table 1: 1-shot/5-shot 5-way	accuracy results
------------------------------	------------------

Method	<i>mini</i> ImageNet	CIFAR100	Caltech-256	CUB	Average
Nearest neighbor (baseline)	44.1 / 55.1	56.1 / 68.3	51.3 / 67.5	52.4 / 66.0	51.0/64.2
MACO [19]	41.1 / 58.3	-	-	60.8 / 75.0	-
Meta-Learner LSTM [33]	43.4 / 60.6	-	-	40.4 / 49.7	-
Matching Nets 42	46.6 / 60.0	50.5 / 60.3	48.1 / 57.5	49.3 / 59.3	48.6 / 59.3
MAML [10]	48.7 / 63.1	49.3 / 58.3	45.6 / 54.6	38.4 / 59.1	45.5 / 58.8
Prototypical Networks [38]	49.4 / 68.2	-	-	-	-
SRPN 29	55.2 / 69.6	-	-	-	-
RELATION NET [40]	57.0/71.1	-	-	-	-
DEML+Meta-SGD 50	58.5 / 71.3 °	61.6 / 77.9 *	62.2 / 79.5 [°]	66.9 / 77.1 [°]	62.3 / 76.4
Dual TriNet 4	58.1 / 76.9 [†]	63.4 / 78.4 †	63.8 / 80.5 †	69.6 / 84.1 *	63.7 / 80.0
Δ -encoder	58.7 / 73.6	65.9 / 80.1	63.9 / 84.7	69.9 / 82.5	64.6 / 80.2

♦ Model also trained on an external large-scale dataset

† Using word embedding trained on large corpus and applied to the label name

 \star Using human annotated class attributes

Feature extractor backbone only trained on the subset of training categories of the target dataset



Table 2: 1-shot/5-shot 5-way accuracy with ImageNet model features (trained on disjoint categories)

Method	AWA2	APY	SUN	CUB
Nearest neighbor (baseline)	65.9 / 84.2	57.9 / 76.4	72.7 / 86.7	58.7 / 80.2
Prototypical Networks	80.8 / 95.3	69.8 / 90.1	74.7 / 94.8	71.9 / 92.4
Δ -encoder	90.5 / 96.4	82.5 / 93.4	82.0 / 93.0	82.2 / 92.6

Feature extractor backbone is a VGG16 backbone pre-trained on ImageNet. The unseen test categories were verified to be disjoint from the ImageNet categories



Are we synthesizing non-trivial samples?







03

Generated samples for 12-way one-shot. The two-dimensional embedding was produced by t-SNE. Best viewed in color.



Are we synthesizing non-trivial samples?



THANKS