

CS-E4890: Deep Learning Learning with few labeled examples

Alexander Ilin

- Deep learning is data hungry. To learn to classify handwritten digits, we need thousands of samples, to learn to classify natural images we need millions of images.
- Suppose we have a custom classification task, for example, we need to classify images to custom classes (not covered by imageNet). What can we do?
 - Collect a lot of training examples and label them.
 - Time consuming and expensive.
 - Sometimes collecting new examples is impossible.
 - Transfer learning
 - Semi-supervised learning
 - Self-supervised learning
 - Few-shot learning (meta learning)

Transfer learning

- Features that are useful for some tasks can be useful for other tasks in the same domain.
- For image classification tasks, it is common to fine-tune the last layer of a deep neural network pre-trained on imageNet (there is a bunch of them in PyTorch).
- For example, we can take the AlexNet (Krizhevsky, 2012) pre-trained on ImageNet and fine-tune the last two layers.



- Example results: Classification of images from the Caltech-101 dataset (Donahue et al., 2013):
 - Yang et al. (2009): a method employing a combination of five traditional hand-engineered image features followed by a multi-kernel based classifier.
 - DeCAF6: Features from the sixth layer of AlexNet pre-trained on ImageNet dataset.



Semi-supervised learning

• Semi-supervised classification: few labeled examples, many unlabeled examples.



Source: (Tarvainen and Valpola, 2017)

- Semi-supervised learning: few labeled examples, many unlabeled examples.
- Semi-supervised learning is possible when the knowledge on p(x) that one gains through the unlabeled data carries information that is useful in the inference of p(y | x).
- In the hypothetical example on the right, modeling the distribution of unlabeled data can improve the classification accuracy.



The labels can be propagated to the unlabeled data in the same cluster yielding better classification accuracy.

- Semi-supervised learning: few labeled examples, many unlabeled examples.
- Semi-supervised learning is possible when the knowledge on p(x) that one gains through the unlabeled data carries information that is useful in the inference of p(y | x).
- In the hypothetical example on the right, modeling the distribution of unlabeled data can improve the classification accuracy.



The labels can be propagated to the unlabeled data in the same cluster yielding better classification accuracy.

- Semi-supervised learning: few labeled examples, many unlabeled examples.
- Semi-supervised learning is possible when the knowledge on p(x) that one gains through the unlabeled data carries information that is useful in the inference of p(y | x).
- In the hypothetical example on the right, modeling the distribution of unlabeled data can improve the classification accuracy.



The labels can be propagated to the unlabeled data in the same cluster yielding better classification accuracy.

- Semi-supervised learning: few labeled examples, many unlabeled examples.
- Semi-supervised learning is possible when the knowledge on p(x) that one gains through the unlabeled data carries information that is useful in the inference of p(y | x).
- In the hypothetical example on the right, modeling the distribution of unlabeled data can improve the classification accuracy.



The labels can be propagated to the unlabeled data in the same cluster yielding better classification accuracy.

- Semi-supervised learning: few labeled examples, many unlabeled examples.
- Semi-supervised learning is possible when the knowledge on p(x) that one gains through the unlabeled data carries information that is useful in the inference of p(y | x).
- In the hypothetical example on the right, modeling the distribution of unlabeled data can improve the classification accuracy.



The labels can be propagated to the unlabeled data in the same cluster yielding better classification accuracy.

- The architecture resembles a ladder (or a U-net).
- The bottom-up pass produces label **y** for a given input **x**.
- For labeled examples, we can compute the standard classification (e.g., cross-entropy) loss using the network output $\hat{\mathbf{y}}$ and the correct label \mathbf{y} .



- The inputs x are corrupted by noise during training (e.g, we never use clean images as inputs).
- The top-down pass tries to reconstruct the clean (without noise) input **x**.
- For all examples (both labeled and unlabeled), we compute the denoising cost at the bottom:

denoising cost = $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$

• The minimized cost is the sum of the classification and denoising costs.



Ladder networks: Intuitions

- Intuition: In order to reconstruct the clean image from a noisy image, one has to understand what features are commonly present in images, that is we need to model the data distribution p(x) inside the network.
- Therefore, denoising is an auxiliary task that helps model p(x) and hopefully develop features useful for the primary classification task.
- The label itself cannot contain enough information to reconstruct the input. We need skip connections to pass low-level details from the bottom-up pass.
- Note: corruption and denoising happens on multiple levels.



□-model (Laine and Aila, 2016)

- The problem with the Ladder network: We need to model the whole distribution $p(\mathbf{x})$ which can contain lots of details irrelevant for the classification task.
- Laine and Aila (2016) proposed a simplification of the Ladder networks that does not contain the top-down pass.
- There are two copies of the same network performing computations for x₁ and x₂, which is the same training example changed with different transformations.
- The cost for unlabeled data is

consistency cost =
$$\|\mathbf{z}_1 - \mathbf{z}_2\|^2 = \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|^2$$

 In the Π-model, the gradients propagate through both networks, z is the input of the softmax (logit).



- The intuition: we do not know the correct output for unlabeled data but we know that the output should not change for a transformed input.
- Since the introduction of the Π-model, the majority of the semi-supervised methods have been based on optimizing consistency between different transformations of the same training examples.
- The idea resembles siamese networks (Bromley et al., 1993).



 Instead of using two copies of the same network, one of the networks (teacher) is obtained by computing exponential moving average of the weights of the other (student) network:

$$oldsymbol{ heta}_t^\prime = \gamma oldsymbol{ heta}_{t-1}^\prime + (1-\gamma) oldsymbol{ heta}_t$$

• The same consistency cost is minimized on both labeled and unlabeled data:

consistency cost = $\left\| f(\mathbf{x}_1, \boldsymbol{\theta}_t) - f(\mathbf{x}_2, \boldsymbol{\theta}_t') \right\|^2$

• The teacher is more accurate than the student. The gradients propagate only through the student.



	Student	Teacher	Teacher label	CIFAR-10 accuracy
Algorithm	augment.	augment.	post-processing	4K labels (from here)
Fully supervised Wide ResNet with 50K labels	-	-	-	94.6
П-Model (Laine and Aila, 2016)	Weak	Weak	-	87.84
VAT (Miyato et al., 2017)	Adversarial	-	-	88.64
Mean Teacher (Tarvainen and Valpola, 2017)	Weak	Weak	-	93.72
UDA (Xie et al., 2019)	Strong	Weak	Sharpening	94.73
MixMatch (Berthelot et al., 2019)	Weak	Weak	Sharpening	93.76
ReMixMatch (Berthelot et al., 2019)	Strong	Weak	Sharpening	94.86
FixMatch (Sohn et al., 2020)	Strong	Weak	Pseudo-labeling	95.69

- Weak augmentations:
 - Translation, flip, Gaussian noise.
 - MixMatch uses MixUp (Zhang et al., 2018).
- Adversarial: Use adversarial examples for data transformation.
- Strong augmentations: uniformly sample all image processing transformations from the Python Image Library (PIL).

	Student	Teacher	Teacher label	CIFAR-10 accuracy
Algorithm	augment.	augment.	post-processing	4K labels (from here)
Fully supervised Wide ResNet with 50K labels	-	-	-	94.6
Π-Model (Laine and Aila, 2016)	Weak	Weak	-	87.84
VAT (Miyato et al., 2017)	Adversarial	_	-	88.64
Mean Teacher (Tarvainen and Valpola, 2017)	Weak	Weak	-	93.72
UDA (Xie et al., 2019)	Strong	Weak	Sharpening	94.73
MixMatch (Berthelot et al., 2019)	Weak	Weak	Sharpening	93.76
ReMixMatch (Berthelot et al., 2019)	Strong	Weak	Sharpening	94.86
FixMatch (Sohn et al., 2020)	Strong	Weak	Pseudo-labeling	95.69

• Label sharpening: Use (small) temperature parameter au to compute the targets:

$$p_i^{ ext{teacher}} = rac{ \exp(z_i/ au) }{\sum_j \exp(z_j/ au) }$$

• Pseudo-labeling: use "hard" labels (i.e., the arg max of the model's output) and only retain the teacher's labels whose largest class probability fall above a predefined threshold.

	Student	Teacher	Teacher label	CIFAR-10 accuracy
Algorithm	augment.	augment.	post-processing	4K labels (from here)
Fully supervised Wide ResNet with 50K labels	-	-	-	94.6
П-Model (Laine and Aila, 2016)	Weak	Weak	-	87.84
VAT (Miyato et al., 2017)	Adversarial	-	-	88.64
Mean Teacher (Tarvainen and Valpola, 2017)	Weak	Weak	-	93.72
UDA (Xie et al., 2019)	Strong	Weak	Sharpening	94.73
MixMatch (Berthelot et al., 2019)	Weak	Weak	Sharpening	93.76
ReMixMatch (Berthelot et al., 2019)	Strong	Weak	Sharpening	94.86
FixMatch (Sohn et al., 2020)	Strong	Weak	Pseudo-labeling	95.69

• Mean Teacher: Use exponential moving average of parameters to get the teacher model

$$\boldsymbol{\theta}_t' = \gamma \boldsymbol{\theta}_{t-1}' + (1-\gamma) \boldsymbol{\theta}_t$$

- MixMatch: The target is computed as the average prediction obtained for K augmentations.
- ReMixMatch: Distribution alignment encourages the marginal distribution of predictions on unlabeled data to be close to the marginal distribution of ground-truth labels.

Self-supervised learning

- The assumption that is made in semi-supervised learning: The unlabeled examples belong to the same classes. This can be difficult to assure in many practical applications.
- Can we lean useful representations in a completely unsupervised manner?
- Self-supervised learning: Invent an auxiliary task that can be learned in an unsupervised manner and use the developed features for the downstream task.
- In order to succeed, the auxiliary task should be relevant for the downstream task.

Examples of self-supervised learning models



Ladder networks



Discriminative unsupervised feature learning (Dosovitskiy et al., 2014)

- One of the early works on self-supervised training.
- A convolutional neural network is pre-trained using an artificially created learning task.
 - *N* patches of size 32 × 32 are sampled from different images at varying positions and scale.
 - Each patch is transformed multiple times using (a composition of elementary) transformations.
 - The task is to classify a transformed image patch to one of the *N* classes that correspond to the original *N* patches.
- The features produced by the CNN are used as inputs of a support vector machine classifier.

sampled patches



transformed patches



elementary transformations used: translation, scaling, rotation, contrast (raise S and V components of the HSV color representation), color (change the H component of the HSV representation)

Contrastive Predictive Coding (van den Oord et al., 2018)

- The goal is to learn representations that encode the underlying shared information between different parts of the (high-dimensional) signal. At the same time we want to discard low-level information and noise that is more local.
- When predicting further in the future, the amount of shared information becomes much lower, and the model needs to infer more global structure. These 'slow features' that span many time steps are often more interesting (e.g., phonemes and intonation in speech, objects in images, or the story line in books.)



Contrastive Predictive Coding (CPC): The model

- A non-linear encoder g_{enc} maps the input sequence of observations x_t to a sequence of latent representations $z_t = g_{enc}(x_t)$, potentially with a lower temporal resolution.
- Next, an autoregressive model g_{ar} summarizes all $z \le t$ in the latent space and produces a context latent representation $c_t = g_{ar}(z \le t)$.



- The task: Given the context c_t, select the correct future code z_{t+k} = g_{enc}(x_{t+k}) after k steps among N alternatives {z_{t+k}, z_{τ1}, ..., z_{τN-1}}.
- The alternatives z_τ = g_{enc}(x_τ) are selected as encoder outputs produced for inputs x_τ randomly selected from the dataset (for example, from the same input sequence).
- The loss is the categorical cross-entropy of selecting the correct encoding:

$$\mathcal{L} = -\log rac{\exp(z_{t+k}^ op W_k c_t)}{\sum\limits_j \exp(z_{ op}^ op W_k c_t)}$$

 The logits produced by the classifier are postulated to have the form z^T_τ W_kc_t.



- The quality of developed representations are tested by training a classifier using representations c_t as features on the phone classification task:
 - linear classifier: 64.6
 - MLP classifier: 72.5
- CPC captures both speaker identity and speech contents.



t-SNE visualization of audio (speech) representations for a subset of 10 speakers (out of 251). Every color represents a different speaker.

Method	ACC
Phone classification	
Random initialization	27.6
MFCC features	39.7
CPC	64.6
Supervised	74.6
Speaker classification	
Random initialization	1.87
MFCC features	17.6
CPC	97.4
Supervised	98.5

Classification accuracy on audio data. Phone classification: 41 classes Speaker classification: 251 classes. random initialization: random g_{enc} and • Contrastive Predictive Coding for images:



Top-5 ACC
48.3
53.1
59.2
62.5
69.3
73.6

ImageNet top-5 unsupervised classification results.

• The quality of developed representations are tested by training a linear classifier using RNN outputs c_t as input features.

A Simple Framework for Contrastive Learning (SimCLR) (Chen et al., 2020a, 2020b)

- SimCLR can be seen as the siamese networks processing two different augmentations of the same image (as in the Π -model, Mean Teacher) combined with the CPC contrastive loss.
 - Sample a minibatch of *N* examples.
 - Augment each example with two different transformations, which results in 2*N* data points.
 - Process each example with a deep neural network $\mathbf{z} = g(f(\mathbf{x}))$.
 - The training task is similar to CPC: For each image in the minibatch, we need to select the other image that was produced with the other transformation of the same original image among 2N 1 alternatives

$$l_{i,j} = -\log rac{\exp(ext{sim}(\mathbf{z}_i, \mathbf{z}_j / au))}{\sum\limits_{k=1}^{2N} \mathbb{1}_{[k
eq i]} \exp(ext{sim}(\mathbf{z}_i, \mathbf{z}_k) / au)}$$



SimCLR (Chen et al., 2020a)

$$l_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j / \tau))}{\sum\limits_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k) / \tau)}$$

• The cosine similarity is used:

$$sim(\mathbf{u}, \mathbf{v}) = \mathbf{u}^\top \mathbf{v} / \|\mathbf{u}\| \|\mathbf{v}\|$$

- As representation **h**, we take the output of an intermediate layer (two layers before the output).
 - The intuition: The formulated task of contrastive learning needs development of the right features but the downstream task is likely to be different and to require a different post-processing head.



• Three augmentations are sequentially applied:



original



1) random cropping followed by resize to the original size



2) random color distortions



3) random Gaussian blur

- Random cropping and color distortions give the largest boost in performance.
- It is important to apply both transformations for good performance. Why?

• Three augmentations are sequentially applied:



original



1) random cropping followed by resize to the original size



2) random color distortions



3) random Gaussian blur

- Random cropping and color distortions give the largest boost in performance.
- It is important to apply both transformations for good performance. Why?
 Because with one transformation only it is relatively easy to find matching pairs of images.

Method	Architecture	Param (M)	Top 1	Top 5
Methods using R	esNet-50:			
Local Agg.	ResNet-50	24	60.2	-
MoCo	ResNet-50	24	60.6	-
PIRL	ResNet-50	24	63.6	-
CPC v2	ResNet-50	24	63.8	85.3
SimCLR (ours)	ResNet-50	24	69.3	89.0
Methods using o	ther architectures	:		
Rotation	RevNet-50 $(4 \times)$) 86	55.4	-
BigBiGAN	RevNet-50 (4×)) 86	61.3	81.9
AMDIM	Custom-ResNet	626	68.1	-
CMC	ResNet-50 $(2\times)$	188	68.4	88.2
MoCo	ResNet-50 $(4 \times)$	375	68.6	-
CPC v2	ResNet-161 (*)	305	71.5	90.1
SimCLR (ours)	ResNet-50 $(2\times)$	94	74.2	92.0
SimCLR (ours)	ResNet-50 $(4 \times)$	375	76.5	93.2

ImageNet accuracies of linear classifiers trained on representations learned with different self-supervised methods

• The whole base network is fine-tuned on the few labeled data without regularization.

		Label fraction			
Method	Architecture	1%	10%		
		To	p 5		
Supervised baseline	ResNet-50	48.4	80.4		
Methods using other labe	l-propagation:				
Pseudo-label	ResNet-50	51.6	82.4		
VAT+Entropy Min.	ResNet-50	47.0	83.4		
UDA (w. RandAug)	ResNet-50	-	88.5		
FixMatch (w. RandAug)	ResNet-50	-	89.1		
S4L (Rot+VAT+En. M.)	ResNet-50 (4 \times)	-	91.2		
Methods using representa	tion learning only:				
InstDisc	ResNet-50	39.2	77.4		
BigBiGAN	RevNet-50 $(4 \times)$	55.2	78.8		
PIRL	ResNet-50	57.2	83.8		
CPC v2	ResNet-161(*)	77.9	91.2		
SimCLR (ours)	ResNet-50	75.5	87.8		
SimCLR (ours)	ResNet-50 $(2 \times)$	83.0	91.2		
SimCLR (ours)	ResNet-50 $(4\times)$	85.8	92.6		

ImageNet accuracy of models trained with few labels

SimCLRv2 (Chen et al., 2020b)

- SimCLRv2 improves the performance in the semi-supervised scenario by the following procedure:
 - 1. Pre-train SimCLR in an unsupervised way.
 - 2. Fine-tine the model using the few labeled examples (this is same as in SimCLRv1).
 - 3. Perform "knowledge distillation" using both labeled and unlabeled data.
- This last step is essentially equivalent to consistency-based semi-supervised training.
 - The teacher is fixed to the network obtained after step 2.
 - Only weak augmentations are used.
 - The student can have a smaller architecture.



Bootstrap your own latent (BYOL) (Grill et al., 2020)

- BYOL can be seen as an extension of the Mean Teacher to the fully unsupervised scenario:
 - The teacher network is obtained by computing exponential moving average of the weights of the student.
 - The two networks process two different transformations of the same example.
 - The objective is

consistency cost = $\left\| q(\mathbf{z}) - \mathbf{z}' \right\|^2$

where $\mathbf{z} = f(\mathbf{x}, \boldsymbol{\theta}_t)$, $\mathbf{z}' = f(\mathbf{x}', \boldsymbol{\theta}'_t)$

- Differences:
 - Use strong augmentations (same set of augmentations as in SimCLR).
 - Use an extra predictor MLP q(z) in the student network.



consistency cost = $\|q(\mathbf{z}) - \mathbf{z}'\|^2$

- The representations can collapse: the network can learn to produce the same output z for any input x, which would yield a zero consistency cost.
- In practice, this does not happen for the given architecture. In fact, the representations collapse if the predictor network q(z) is removed.



Method	Top-1	Top-5	Method	Architecture	Param.	Top-1	То
Local Agg.	60.2	-	SimCLR [8]	ResNet-50 ($2\times$)	94M	74.2	9
PIRL [35]	63.6	-	CMC [11]	ResNet-50 $(2\times)$	94M	70.6	8
CPC v2 [32]	63.8	85.3	BYOL (ours)	ResNet-50 $(2\times)$	94M	77.4	9
CMC [11]	66.2	87.0	CPC v2 [32]	ResNet-161	305M	71.5	9
SimCLR [8]	69.3	89.0	MoCo [9]	ResNet-50 ($4 \times$)	375M	68.6	
MoCo v2 [37]	71.1	-	SimCLR [8]	ResNet-50 $(4\times)$	375M	76.5	9
InfoMin Aug. [12]	73.0	91.1	BYOL (ours)	ResNet-50 $(4\times)$	375M	78.6	9
BYOL (ours)	74.3	91.6	BYOL (ours)	ResNet-200 $(2\times)$	250M	79.6	94

(a) ResNet-50 encoder.

(b) Other ResNet encoder architectures.

ImageNet accuracies of linear classifiers trained on representations learned with different self-supervised methods • The whole base network is fine-tuned on the few labeled data with spatial augmentations, i.e., random crops with resize and random flips.

Method	Top	b -1	Top	5 -5	Method	Architecture	Param.	Top	b -1	Top	p-5
	$1\%^{-1}$	10%	$1\%^{-1}$	10%				1%	10%	1%	10%
Supervised [77]	25.4	56.4	48.4	80.4	CPC v2 [32]	ResNet-161	305M	-	-	77.9	91.2
					SimCLR [8]	ResNet-50 $(2\times)$	94M	58.5	71.7	83.0	91.2
InstDisc	-	-	39.2	77.4	BYOL (ours)	ResNet-50 $(2\times)$	94M	62.2	73.5	84.1	91.7
PIRL [35]	-	-	57.2	83.8	SimCLR [8]	ResNet-50 $(4\times)$	375M	63.0	74.4	85.8	92.6
SimCLR [8]	48.3	65.6	75.5	87.8	BYOL (ours)	ResNet-50 $(4\times)$	375M	69.1	75.7	87.9	92.5
BYOL (ours)	53.2	68.8	78.4	89.0	BYOL (ours)	ResNet-200 (2×)	250M	71.2	77.7	89.5	93.7

(a) ResNet-50 encoder.

(b) Other ResNet encoder architectures.

ImageNet accuracy of models trained with few labels

Few-shot learning

• People can learn new concepts from just a single example:



- Yet machine learning algorithms typically require thousands of examples to perform with similar accuracy (Lake et al., 2015).
- Few-shot learning: How can we train an accurate model using a very small amount of training data?

- The problem of few-shot learning attracted a lot of attention after releasing the Omniglot challenge (Lake et al., 2015).
- The authors also proposed a (non-deep) model that represents concepts as simple programs that best explain observed examples under a Bayesian criterion (BPL).
- BPL achieves human-level performance on the one-shot classification task.

Α

One-shot classification



		ಜ		
١	Where	is an	other'	?
ല	لور	ബ	പ	ఔ
ਝ	ಖ	77	പ	ಝ
â	J	63	ె	ದೆ
ನ	സ	ಲ	മ്	പ്

Part A of Omniglot challenge: Two trials of one-shot classification, where a single image of a new character is presented (top) and the goal is to select another example of that character amongst other characters from the same alphabet (in the grid below). • There has been a great progress of deep learning methods in few-shot learning tasks but deep learning seems to be behind BPL and human level.

• Deep learning methods rely heavily on data augmentation.

	Original Within Within alphabet alphabet (minimal)		Augn Within alphabet	nented Between alphabet
background set				
# alphabets	30	5	30	40
# classes	964	146	3,856	4,800
2015 results				
Humans	$\leq 4.5\%$,)		
BPL	3.3%	4.2%		
Simple ConvNe	t 13.5%	23.2%		
Siamese Net			8.0%*	
2016-2018 resul	ts			
Prototypical No	et 13.7%	30.1%	6.0%	4.0%
Matching Net				6.2%
MAML				4.2%
Graph Net				$\mathbf{2.6\%}$
ARC			1.5%*	2.5%*
RCN	7.3%			
VHE	18.7%			4.8%

One-shot classification error rate

* results used additional data augmentation beyond class expansion

- Consider the task of *one-shot* classification: building a classifier from one training example.
- A siamese neural network (Bromley et al., 1993) is trained to compare a pair of examples to decide whether they belong to the same class or not (binary classification):
 - Twin networks which accept distinct inputs
 - The output is computed using a distance (e.g., Euclidean) between the highest-level feature representations of the twin networks.
 - The network is trained on pairs of positive (same class) and negative (distinct classes) examples.
- Works for one-shot learning, few-shot requires tweaking.



We wish to map from a (small) support set of k examples of image-label pairs
 S = {(x_i, y_i)}^k_{i=1} to a classifier c_S(x̂).

 $S \rightarrow c_S(\hat{x})$

- The classifier c_S(x̂) defines a probability distribution ŷ over classes in the support set for a query example x̂.
- Matching networks: Parameterize this mapping $S \rightarrow c_S(\hat{x})$ as a neural network.



Matching networks: The model

• The output of the classifier $c_S(\hat{x})$:

$$\hat{y} = \sum_{i=1}^{k} a(\hat{x}, x_i) y_i$$

with an attention mechanism a:

 $a(\hat{x}, x_i) = \frac{\exp(c(f(\hat{x}), g(x_i)))}{\sum_{j=1}^{k} \exp(c(f(\hat{x}), g(x_j)))}$

- f and g are parameterized as neural networks.
- Intuition: support samples whose representations $g(x_i)$ are close to the representation $f(\hat{x})$ of the query sample contribute more to the output.
- We tune the representations to work well in the few-shot learning scenario.



Matching networks: Episodic training

- *N*-way *K*-shot classification task: Each training example is a classification task with *N* classes and *K* training examples per class.
- One iteration of episodic training:
 - Select *N* classes by randomly sampling from the training set.
 - Select a support set by taking K random samples for each of the selected classes.
 - Select a query set (a few samples from the same classes as in the support set).
 - Do forward computations and compute the classification loss using the query samples.
 - Do backpropagation and update the parameters of the network.



Prototypical networks (Snell et al., 2017)

 Prototypical networks compute an *M*-dimensional representation c_k, or prototype, of each class through an embedding function f_θ. Each prototype is the mean vector of the embedded support points belonging to its class:

$$\mathbf{c}_k = rac{1}{|\mathcal{S}_k|} \sum_{(x_i, y_i) \in \mathcal{S}_k} f_{\boldsymbol{ heta}}(x_i)$$



• Then they produce a distribution over classes for a query point x based on a softmax over distances to the prototypes in the embedding space:

$$p(y = k \mid x) = \frac{\exp(-d(f_{\theta}(x), \mathbf{c}_k))}{\sum_{k'} \exp(-d(f_{\theta}(x), \mathbf{c}_{k'}))}$$

• In one-shot learning scenario, prototypical networks are equivalent to matching networks.

One iteration of episodic training:

- The support set is used to compute the prototypes.
- The query set is used to compute the loss.
- Compute the gradients of the loss and update the parameters θ of the embedding network.



- We want to train a classifier $y = f_{\theta}(\mathbf{x})$ to solve a new few-shot learning task.
- Learning can be done by performing a few iterations of gradient descent (GD).
- In the case of one iteration of GD:

$$oldsymbol{ heta}' \leftarrow oldsymbol{ heta}_0 - lpha
abla_{oldsymbol{ heta}} L((\mathbf{x}_1, \mathbf{y}_1), ..., (\mathbf{x}_k, \mathbf{y}_k))$$

- $\{(x_i, y_i)\}_{i=1}^k$ are the few training examples (support set)
- L is the loss function (for example, cross-entropy for classification tasks)
- α is the learning rate
- $heta_0$ is the vector of the initial values of the parameters
- The idea of meta-learning: we can learn initialization θ_0 and the learning rate α to minimize the loss (on the query set) after the GD-adaptation.

One iteration of episodic training:

- Use the support set to compute the loss and its gradient $\nabla_{\theta} L$.
- Compute adapted values θ' of the parameters (as part of computational graph) with one (or a few) iteration of gradient descent.
- Use the query set to compute the loss with the adapted parameters θ'.
- Compute the loss on the query set.
- Perform backpropagation and update parameters θ₀ and learning rate α.



- MAML requires computation of gradient through gradient, which can be computationally expensive.
- The first-order approximation (which stops gradient propagation through ∇_θL) works almost equally well.



Computational graph

• Simplification of MAML: Instead of backpropagating through the computational graph of MAML, we update the initial parameter values θ_0 towards the adapted parameter values θ_T with a small step ϵ .

```
for iteration 1,2,3,... do
Randomly sample a task T
Perform k > 1 steps of SGD on task T starting from \theta_0: \theta_0 \rightarrow \theta_T
Update: \theta_0 \leftarrow \theta_0 + \epsilon(\theta_T - \theta_0)
return \theta_0
```

• Reptile demo

Home assignment

• You need to implement prototypical networks (Snell et al., 2017).

$$\mathbf{c}_{k} = \frac{1}{|S_{k}|} \sum_{(x_{i}, y_{i}) \in S_{k}} f_{\theta}(x_{i})$$
$$p(y = k \mid x) = \frac{\exp(-d(f_{\theta}(x), \mathbf{c}_{k}))}{\sum_{k'} \exp(-d(f_{\theta}(x), \mathbf{c}_{k'}))}$$



• Papers cited in the lecture slides.