

*Review Series of Recent Deep Learning Papers:*  
Parameter Prediction Paper: Learning  
Feed-Forward One-Shot Learners

Luca Bertinetto, João F. Henriques, Jack Valmadre, Philip H. S. Torr,  
Andrea Vedaldi  
NIPS 2016

Reviewed by : Arshdeep Sekhon

<sup>1</sup>Department of Computer Science, University of Virginia  
<https://qdata.github.io/deep2Read/>

August 25, 2018

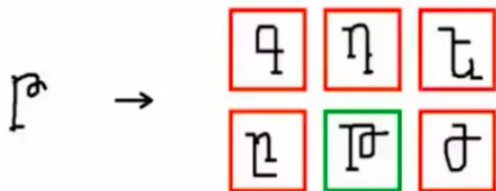
## One Shot Learning

Learn a concept (classifier) from a single example

# Learning Feed-Forward One-Shot Learners

## One Shot Learning

Learn a concept (classifier) from a single example

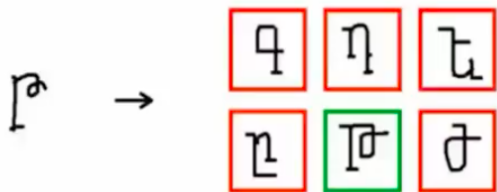


Task: Identify a character from the Armenian Alphabet

# Learning Feed-Forward One-Shot Learners

## One Shot Learning

Learn a concept (classifier) from a single example



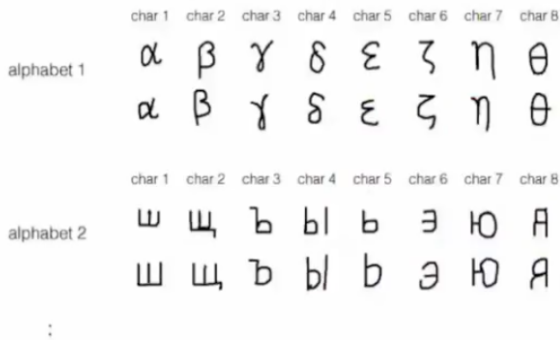
Task: Identify a character from the Armenian Alphabet

The model hasn't seen the alphabet during training

# Learning Feed-Forward One-Shot Learners

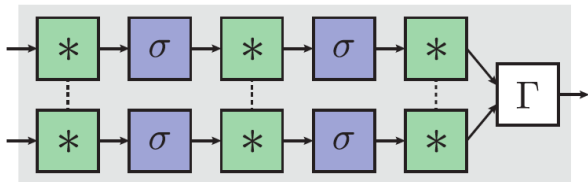
## One Shot Learning

Learn a concept (classifier) from a single example



Offline Training

# Siamese Architecture



$$\min_{W'} \frac{1}{n} \sum_{i=1}^n \text{Loss}(\langle \varphi(x_i; W), \varphi(z_i; W) \rangle, l_i) \quad (1)$$

# Dynamic Parameter Prediction and Learning to Learn Approach

## The Task

Given one exemplar, recognize other instances of the same class.

# Dynamic Parameter Prediction and Learning to Learn Approach

## The Task

Given one exemplar, recognize other instances of the same class.

## Dynamic Parameter Prediction: Online Phase

**Learnet**

**predictor**



# Dynamic Parameter Prediction and Learning to Learn Approach

## The Task

Given one exemplar, recognize other instances of the same class.

## Dynamic Parameter Prediction: Online Phase

### **Learnet**

Given a single example  $z$  predict parameters for predictor

$$W(z) = \omega(z, W') \quad (2)$$

### **predictor**

$$y = \varphi(x, W(z)) \quad (3)$$

# Dynamic Parameter Prediction and Learning to Learn Approach

## The Task

Given one exemplar, recognize other instances of the same class.

## Dynamic Parameter Prediction: Online Phase

### Learnet

Given a single example  $z$  predict parameters for predictor

$$W(z) = \omega(z, W') \quad (2)$$

### predictor

$$y = \varphi(x, W(z)) \quad (3)$$

*'Learning to Learn' approach*

# Standard Discriminative Learning vs One shot Discriminative learning

## Standard Discriminative Learning

- $$\min_W \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, W), l_i) \quad (4)$$

# Standard Discriminative Learning vs One shot Discriminative learning

## Standard Discriminative Learning

- 

$$\min_W \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, W), l_i) \quad (4)$$

- Also add regularization
- However, not enough for one shot learning still.

# Standard Discriminative Learning vs One shot Discriminative learning

## Standard Discriminative Learning

- 

$$\min_W \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, W), l_i) \quad (4)$$

- Also add regularization
- However, not enough for one shot learning still.
- **Solution: Learning to learn**
- Inject prior information in the task

# Standard Discriminative Learning vs One shot Discriminative learning

## Standard Discriminative Learning

- 

$$\min_W \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, W), l_i) \quad (4)$$

- Also add regularization
- However, not enough for one shot learning still.
- **Solution: Learning to learn**
- Inject prior information in the task
- Prior Information about the learning domain is introduced in the offline phase.

The Learnet

$$W = \omega(z_i, W') \quad (5)$$

# Training the Learnet

The Learnet

$$W = \omega(z_i, W') \quad (5)$$

Task: Find optimal parameters  $W'$  of the learner.



# Training the Learnet

The Learnet

$$W = \omega(z_i, W') \quad (5)$$

Task: Find optimal parameters  $W'$  of the learnet.

old standard discriminative objective function

$$\min_W \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, W), \ell_i) \quad (6)$$

$\ell_i$  is true label of  $x_i$

# Training the Learnet

The Learnet

$$W = \omega(z_i, W') \quad (5)$$

Task: Find optimal parameters  $W'$  of the learnet.

old standard discriminative objective function

$$\min_W \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, W), \ell_i) \quad (6)$$

$\ell_i$  is true label of  $x_i$

New One Shot discriminative objective function

$$\min_{W'} \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, \omega(z_i, W')), \ell_i) \quad (7)$$

$\ell_i$  is positive if  $(x_i, z_i)$  are of the same class.

## New One Shot discriminative objective function

$$\min_{W'} \frac{1}{n} \sum_{i=1}^n \text{Loss}(\varphi(x_i, \omega(z_i, W')), \ell_i) \quad (8)$$

$\ell_i$  is positive if  $(x_i, z_i)$  are of the same class.

Training Data: labeled sample pairs  $(x_i, \ell_i)$  and  $(z_i, \ell_i)$   
triplets  $(x_i, z_i, \ell_i)$

# The challenge

A fully connected linear layer:

$$y = Wx + b \quad (9)$$

$x \in \mathbb{R}^d$ , outputs  $y \in \mathbb{R}^k$ , weights  $W \in \mathbb{R}^{d \times k}$  and biases  $b \in \mathbb{R}^k$

# The challenge

A fully connected linear layer:

$$y = Wx + b \quad (9)$$

$x \in \mathbb{R}^d$ , outputs  $y \in \mathbb{R}^k$ , weights  $W \in \mathbb{R}^{d \times k}$  and biases  $b \in \mathbb{R}^k$

## The Learnet

$$W = \omega(z_i, W') \quad (10)$$

$$y = \omega(z)x + b(z) \quad (11)$$

$$\omega : \mathbb{R}^m \rightarrow \mathbb{R}^{d \times k} \quad (12)$$

# The challenge

A fully connected linear layer:

$$y = Wx + b \quad (9)$$

$x \in \mathbb{R}^d$ , outputs  $y \in \mathbb{R}^k$ , weights  $W \in \mathbb{R}^{d \times k}$  and biases  $b \in \mathbb{R}^k$

## The Lernet

$$W = \omega(z_i, W') \quad (10)$$

$$y = \omega(z)x + b(z) \quad (11)$$

$$\omega : \mathbb{R}^m \rightarrow \mathbb{R}^{d \times k} \quad (12)$$

Assuming the Lernet is also a linear layer:

$$\omega(z) = W'z \quad (13)$$

# The challenge

A fully connected linear layer:

$$y = Wx + b \quad (9)$$

$x \in \mathbb{R}^d$ , outputs  $y \in \mathbb{R}^k$ , weights  $W \in \mathbb{R}^{d \times k}$  and biases  $b \in \mathbb{R}^k$

## The Lernet

$$W = \omega(z_i, W') \quad (10)$$

$$y = \omega(z)x + b(z) \quad (11)$$

$$\omega : \mathbb{R}^m \rightarrow \mathbb{R}^{d \times k} \quad (12)$$

Assuming the Lernet is also a linear layer:

$$\omega(z) = W'z \quad (13)$$

Lernet needs to learn  $d \times k \times m$  parameters.  $d=k=100$  and an exemplar with 100 features, total lernet parameters: **1 million.**

# The solution: Reducing output space of the learnet

## Factorized Linear Layers

- 1 Inspired by SVD

$$Wx = U \text{diag}(s) V^T x$$



# The solution: Reducing output space of the learnet

## Factorized Linear Layers

- 1 Inspired by SVD

$$Wx = U \text{diag}(s) V^T x$$

- 2

$$W(z).x = M' \text{diag}(\omega(z)) Mx \quad (14)$$

$$M \in \mathbb{R}^{d \times d} \text{ and } M' \in \mathbb{R}^{d \times k}$$

Offline Phase: Learn constant basis U and V

Online Phase(Test time): predict weights of diagonal transform

# The solution: Reducing output space of the learnet

## Factorized Linear Layers

- 1 Inspired by SVD

$$Wx = U \text{diag}(s) V^T x$$

- 2

$$W(z).x = M' \text{diag}(\omega(z)) Mx \quad (14)$$

$$M \in \mathbb{R}^{d \times d} \text{ and } M' \in \mathbb{R}^{d \times k}$$

Offline Phase: Learn constant basis  $U$  and  $V$

Online Phase (Test time): predict weights of diagonal transform

- 3 Now, the learnet needs to predict just  $d$  parameters.  $\omega(z) : \mathbb{R}^m \rightarrow \mathbb{R}^d$

A convolutional layer:

$$y = W * x + b \quad (15)$$

$x \in \mathbb{R}^{r \times c \times d}$ ,  $W \in \mathbb{R}^{f \times f \times d \times k}$ ,  $y \in \mathbb{R}^{r' \times c' \times k}$  d: the number of input channels, f: filter size, k output channels

A convolutional layer:

$$y = W * x + b \quad (15)$$

$x \in \mathbb{R}^{r \times c \times d}$ ,  $W \in \mathbb{R}^{f \times f \times d \times k}$ ,  $y \in \mathbb{R}^{r' \times c' \times k}$  d: the number of input channels, f: filter size, k output channels

The number of parameters to be predicted by learnet are  $f^2 dk$ .

# Extending to CNNs

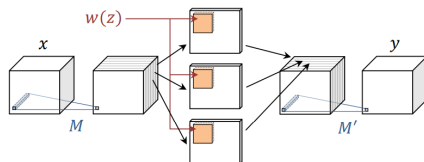
Factorize:

$$y = M' * w(z) *_d M * x + b(z) \quad (16)$$

$$M \in \mathbb{R}^{1 \times 1 \times d \times d}, M' \in \mathbb{R}^{1 \times 1 \times d \times k}, w(z) \in \mathbb{R}^{1 \times f \times f \times d}$$

$*_d$  does independent filtering of  $d$  channels:

$x *_d y$  is the convolution of corresponding channels in  $x$  and  $y$ .



Factorized Convolutional layer

# Extending to CNNs

Factorize:

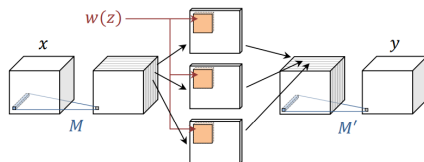
$$y = M' * w(z) *_d M * x + b(z) \quad (16)$$

$$M \in \mathbb{R}^{1 \times 1 \times d \times d}, M' \in \mathbb{R}^{1 \times 1 \times d \times k}, w(z) \in \mathbb{R}^{1 \times f \times f \times d}$$

$*_d$  does independent filtering of  $d$  channels:

$x *_d y$  is the convolution of corresponding channels in  $x$  and  $y$ .

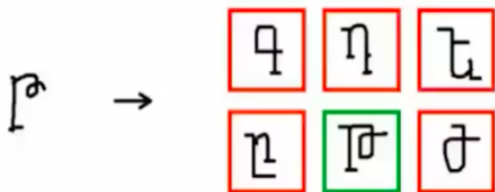
The number of elements to be predicted by learnt now are :  $f^2 d$



Factorized Convolutional layer

# An example: Character Recognition in Alphabets

Test Phase(Online):



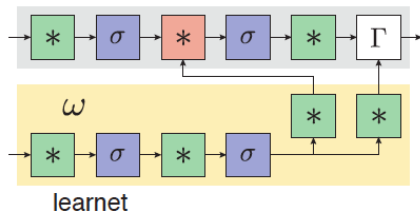
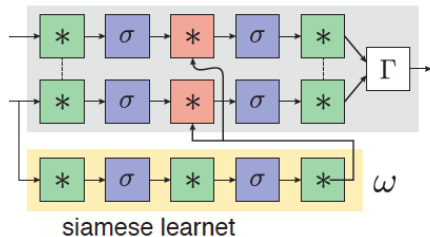
# An example: Character Recognition in Alphabets

Training Phase(Offline):





# Architectures



Model	Error Percentage
Siamese (shared)	41.8
Siamese (unshared)	34.6
Siamese (unshared, factorized)	33.6
Siamese Learnet (shared)	31.4
Learnet	28.6