

# Understanding meta-learning for fast adaptation

Kavya Ranjan Saxena  
Dr. Vipul Arora

Department of Electrical Engineering  
Indian Institute of Technology, Kanpur

January 4, 2023



# Flow of the tutorial

- 1 Introduction
  - Supervised learning
- 2 Different techniques
  - Fine-tuning
  - Few-shot learning
  - Meta-learning
- 3 Meta-learning : Theory + Hands-on
  - Theory
  - Example
  - Hands-on

# Introduction

# Supervised learning

- $\hat{y} = F_{\theta}(x)$  <sup>12</sup>
- Define loss :  $\mathcal{L}(y, \hat{y})$
- Train it using data  $\mathcal{D} = \{(x, y)\}$  (input-output pairs)

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \sum_{(x,y) \in \mathcal{D}} \mathcal{L}(y, \hat{y}) \quad (1)$$

- Evaluate on test data (accuracy, mse, etc.)

---

<sup>1</sup>Bishop, C. M., Nasrabadi, N. M. (2006). Pattern recognition and machine learning (Vol. 4, No. 4, p. 738). New York: springer

<sup>2</sup>Goodfellow, I., Bengio, Y., Courville, A. (2016). Deep learning. MIT press

# Limitations of Supervised learning

- **Problem 1:** Needs a large amount of labelled data (input-output) pairs
  - annotations may be difficult to obtain
    - manual labelling
    - scarce data
    - annotations may be computationally too expensive
  - annotations may be changing over time
- **Problem 2:** Even though large amount of labelled data is available, the model performs well on target data if data distribution of target data is same as source data i.e  $P(X_t) = P(X_s)$

- Semi-supervised learning
  - self-supervised
  - unsupervised
- Few-shot learning
- Model adaptation
  - transfer learning
  - meta learning

Different techniques

Model pre-training:

- Supervised training on related tasks
  - Large networks trained on large labeled datasets
  - e.g., VGG16<sup>3</sup> for images classification
    - No. of trainable parameters 100M
    - 1.3M images
    - 1000 classes

---

<sup>3</sup>Simonyan and Zisserman. "Very deep convolutional networks for large-scale image recognition", ICLR 2015



# Fine-tuning

Supervised model pre-training:

- $\hat{y} = F_{\phi}(F_{\theta}(x))$
- Train  $F_{\phi}(F_{\theta}(x))$  for a supervised task with big data
- Discard  $F_{\phi}$  and use  $F_{\phi'}(F_{\theta})$  for the new task
- Freeze  $F_{\theta}$  and update  $F_{\phi'}$
- Update the complete model end-to-end

# Few-shot learning

- To overcome the problem of scarcity of large amount of labeled data
- Classifies data from target domain when few labeled training samples in target domain
- N-way K-shot classification problem
- How to solve this problem?
- Meta-Learning!
- Meta Learning-
  - Metric-based : Matching networks, Prototypical networks, Relation Networks
  - Gradient-based : Model agnostic meta learning<sup>4</sup>

---

<sup>4</sup>Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In International conference on machine learning, pp. 1126-1135. PMLR, 2017.

# Matching Networks

First method to solve few-shot learning problem

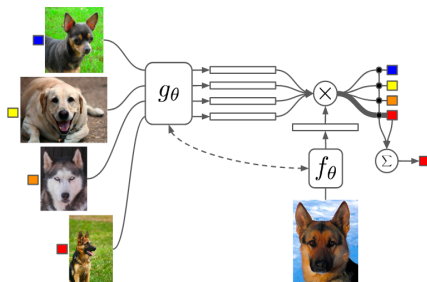


Figure: Matching Networks<sup>5</sup>

---

<sup>5</sup>Vinyals, Oriol, Charles Blundell, Timothy Lillicrap, and Daan Wierstra. "Matching networks for one shot learning." Advances in neural information processing systems 29 (2016).

# Prototypical Networks

Similar to matching networks, but there are small differences!

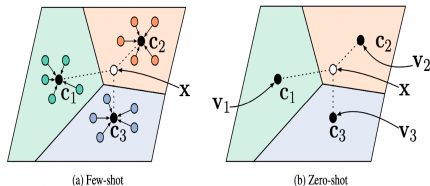


Figure 1: Prototypical Networks in the few-shot and zero-shot scenarios. **Left:** Few-shot prototypes  $c_k$  are computed as the mean of embedded support examples for each class. **Right:** Zero-shot prototypes  $c_k$  are produced by embedding class meta-data  $v_k$ . In either case, embedded query points are classified via a softmax over distances to class prototypes:  $p_\phi(y = k | \mathbf{x}) \propto \exp(-d(f_\phi(\mathbf{x}), c_k))$ .

Figure: Prototypical Networks<sup>6</sup>

---

<sup>6</sup>Snell, Jake, Kevin Swersky, and Richard Zemel. "Prototypical networks for few-shot learning." Advances in neural information processing systems 30 (2017).

# Relation Networks

Almost similar to prototypical networks, but there are small differences!

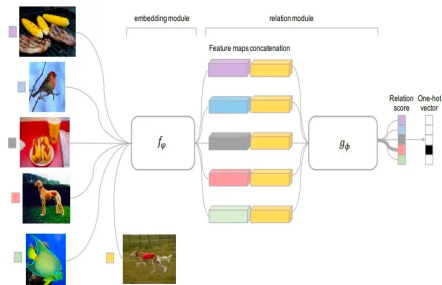


Figure: Relation Networks<sup>7</sup>

<sup>7</sup>Sung, Flood, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M. Hospedales. "Learning to compare: Relation network for few-shot learning." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1199-1208. 2018.

# Meta-learning

What is meta-learning?

- **learning-to-learn** algorithm
  - Learn optimal model initialization
  - Learn optimal learning method
  - Learn optimal hyperparameters

# Meta-learning

What is meta-learning?

- **learning-to-learn** algorithm
  - **Learn optimal model initialization**
  - Learn optimal learning method
  - Learn optimal hyperparameters

Meta-learning : Theory + Hands-on



# Framework

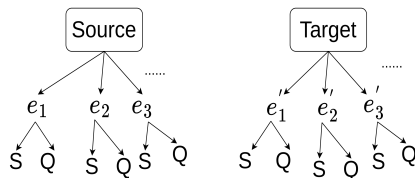
Algorithm to learn optimal model initialization : **Model agnostic meta learning**

- $y = F_{\theta_e}(x)$
- $\theta_e$  is initialized with  $\theta_0$
- Train  $\theta_e$  on S

$$\theta_e \leftarrow \theta_e - \alpha \frac{\partial \mathcal{L}(S)}{\partial \theta_e}$$

- Train  $\theta_0$  on Q

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{\partial \mathcal{L}(Q)}{\partial \theta_0}$$



S : Support Set ( $X_S, Y_S$ )

Q : Query Set ( $X_Q, Y_Q$ )

# Framework

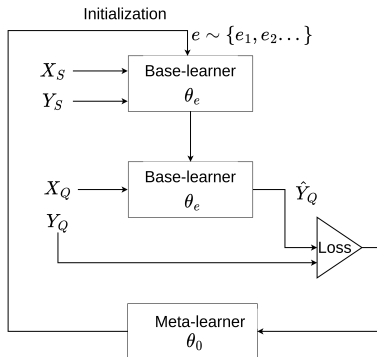


Figure: Meta-training framework

# Framework

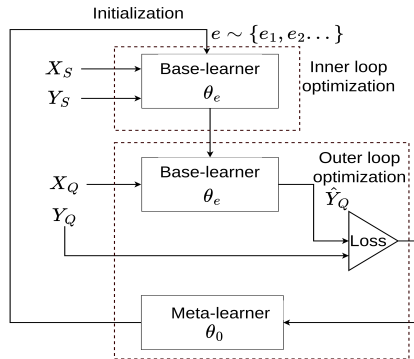


Figure: Meta-training framework

# Framework

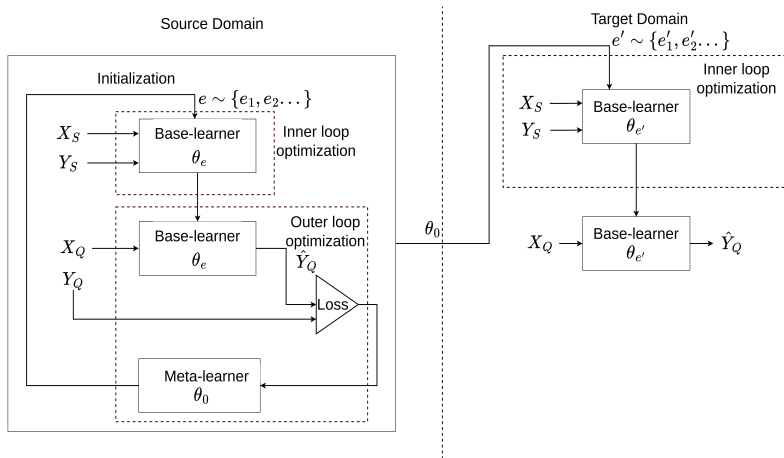


Figure: Overall Meta-learning framework

# Meta-learning theory

---

## Algorithm 1 Meta Training Algorithm

---

**Require:** Different tasks:  $e = \{e_1, e_2, e_3, \dots\}$

**Require:**  $\alpha, \beta$ : learning rates

- 1: Initialize  $\theta_e$  with  $\theta_0$
  - 2: **for** all episodes  $e$  **do**
  - 3:   Divide  $e$  into support set  $S$  and query set  $Q$
  - 4:   Update  $\theta_e$  on support set  $S$  with gradient descent:  
    
$$\theta_e \leftarrow \theta_e - \alpha \frac{\partial \mathcal{L}(S)}{\partial \theta_e}$$
  - 5:   Update  $\theta_0$  on query set  $Q$  with gradient descent:  
    
$$\theta_0 \leftarrow \theta_0 - \alpha \frac{\partial \mathcal{L}(Q)}{\partial \theta_0}$$
  - 6: **end for**
-

---

**Algorithm 2** Meta Testing Algorithm

---

**Require:** Different tasks:  $e' = \{e'_1, e'_2, e'_3, \dots\}$

**Require:**  $\alpha$ : learning rate

- 1: Initialize  $\theta_{e'}$  with  $\theta_0$  obtained from **Algorithm 1**
  - 2: **for** all episodes  $e'$  **do**
  - 3:   Divide  $e'$  into support set  $S$  and query set  $Q$
  - 4:   Update  $\theta_{e'}$  on support set  $S$  with gradient descent:  
    
$$\theta_{e'} \leftarrow \theta_{e'} - \alpha \frac{\partial \mathcal{L}(S)}{\partial \theta_{e'}}$$
  - 5:   Test the updated  $\theta_{e'}$  on query set  $Q$
  - 6: **end for**
-

Submitted to ICASSP 2023

## DEEP DOMAIN ADAPTATION FOR POLYPHONIC MELODY EXTRACTION

*Kavya Ranjan Saxena      Vipul Arora*

Department of Electrical Engineering, Indian Institute of Technology Kanpur, India

### ABSTRACT

Extraction of the predominant pitch from polyphonic audio is one of the fundamental tasks in the field of music information retrieval and computational musicology. To accomplish this task using machine learning, a large amount of labeled audio data is required to train the model that predicts the pitch contour. But a classical model pre-trained on data

learning-based domain adaptation approaches. Meta-learning has been recently used to improve the performance of few-shot learning problems [7][8]. The common approaches to meta-learning are metric-based [7], model-based [8], and optimization-based [9] learning that improve the learning speed [10]. To the best of our knowledge, no such work on meta-learning-based domain adaptation for polyphonic

## IEEE Sensors Letters 2021

Sensor Signal Processing

---

### Few-shot calibration of low-cost air pollution (PM<sub>2.5</sub>) sensors using meta-learning

Kalpit Yadav<sup>1\*</sup>, Vipul Arora<sup>1</sup>, Mohit Kumar<sup>2</sup>, Sachchida Nand Tripathi<sup>2,3</sup>, Vidyanand Motiram Motghare<sup>4</sup>, and Karansingh A. Rajput<sup>4</sup>

<sup>1</sup>Department of Electrical Engineering, Indian Institute of Technology Kanpur, India

<sup>2</sup>Centre for Environmental Science and Engineering, Indian Institute of Technology Kanpur, India

<sup>3</sup>Department of Civil Engineering, Indian Institute of Technology Kanpur, India

<sup>4</sup>Maharashtra Pollution Control Board, India

**Abstract**—Low-cost particulate matter sensors are transforming air quality monitoring because they have greater mobility as compared to reference monitors. Calibration of these low-cost sensors requires training data from co-deployed reference monitors. Machine Learning based calibration gives better performance than conventional techniques, but requires a large amount of training data from the sensor, to be calibrated, co-deployed with a reference monitor. In this work, we propose novel transfer learning methods for quick calibration of sensors with minimal co-deployment with reference monitors. Transfer learning utilizes a large amount of data from other sensors along with a limited amount of data from the target sensor. Our experimentation finds the proposed Model-Agnostic-Meta-Learning (MAML) based transfer learning method to be significantly more effective over other competitive baselines, reducing the calibration errors by 32% and 15% relative to the raw observations and the best baseline, respectively.

**Index Terms**—Air quality, low-cost sensor calibration, few-shot learning, MAML, machine learning



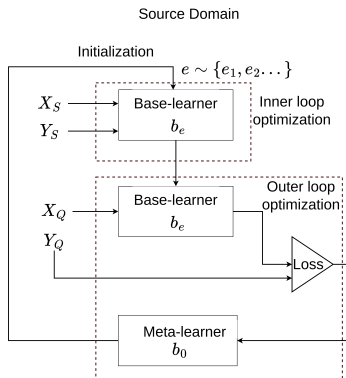
# Linear Regression

Consider a linear regression problem and apply meta-learning to it.

- Example:  $\hat{y} = f_{\theta}(x) = \sigma(a + bx)$
- $\theta = [a, b]$  are the model parameters
- Consider the loss function as  $\mathcal{L} = (y - \hat{y})^2$ , where  $y$  is the target value and  $\hat{y}$  is the predicted output.
- Let us derive equations for meta-learning the initializer  $b_0$  for weight  $b$ .

# Linear Regression

Modifying the meta-training framework for our regression problem.



# Linear Regression

- Inner-loop optimization

$$b_e \leftarrow b_e - \alpha \nabla_{b_e} \mathcal{L}_S(f(b_e)) \quad (2)$$

$$b_e \leftarrow b_e - \alpha \frac{\partial \mathcal{L}_S}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b_e} \quad (3)$$

# Linear Regression

$$\frac{\partial \mathcal{L}_S}{\partial \hat{y}} = -2(y - \hat{y}) \quad (4)$$

$$\frac{\partial y}{\partial b_e} = \frac{\partial \sigma(a_e + b_e x)}{\partial b_e} \quad (5)$$

$$\frac{\partial y}{\partial b_e} = \hat{y}(1 - \hat{y})x \quad (6)$$

# Linear Regression

Putting eq. 4 and eq. 6 in eq. 3, we get:

$$b_e \leftarrow b_e + 2\alpha(y - \hat{y})\hat{y}(1 - \hat{y})x \quad (7)$$

# Linear Regression

- Outer-loop optimization

$$b_0 \leftarrow b_0 - \alpha \nabla_{b_0} \mathcal{L}_Q(f(b_0)) \quad (8)$$

$$b_0 \leftarrow b_0 - \alpha \frac{\partial \mathcal{L}_Q}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b_e} \frac{\partial b_e}{\partial b_0} \quad (9)$$

# Linear Regression

$$\frac{\partial \mathcal{L}_Q}{\partial \hat{y}} = -2(y - \hat{y}) \quad (10)$$

$$\frac{\partial \hat{y}}{\partial b_e} = \hat{y}(1 - \hat{y})x \quad (11)$$

# Linear Regression

But, how to compute:

$$\frac{\partial b_e}{\partial b_0} = ? \quad (12)$$

No direct relation between  $b_e$  and  $b_0$ !



# Linear Regression

Develop a relation between  $b_e$  and  $b_0$  as :

$$b_e = b_0 - \gamma \frac{\partial \mathcal{L}(f(b_0))}{\partial b_0} \quad (13)$$

So, the derivative of the above equation w.r.t  $b_0$  is given by:

$$\frac{\partial b_e}{\partial b_0} = 1 - \gamma \frac{\partial^2 \mathcal{L}(f(b_0))}{\partial b_0^2} \quad (14)$$

# Linear Regression

Solving eq. 14, we get:

$$\frac{\partial b_e}{\partial b_0} = 1 - \gamma \frac{\partial}{\partial b_0} \frac{\partial \mathcal{L}}{\partial b_0} \quad (15)$$

$$\frac{\partial b_e}{\partial b_0} = 1 - \gamma \frac{\partial}{\partial b_0} - 2(y - \hat{y})(\hat{y}(1 - \hat{y})x) \quad (16)$$

# Linear Regression

$$\begin{aligned} \frac{\partial b_e}{\partial b_0} = & 1 + 2\gamma x [\hat{y}(1 - \hat{y}) \frac{\partial(y - \hat{y})}{\partial b_0} + \\ & (y - \hat{y})(1 - \hat{y}) \frac{\partial \hat{y}}{\partial b_0} + (y - \hat{y})\hat{y} \frac{\partial(1 - \hat{y})}{\partial b_0}] \end{aligned} \quad (17)$$

On solving eq.17 we get,

$$\frac{\partial b_e}{\partial b_0} = 1 - 2\gamma x^2 (y - \hat{y}) [2\hat{y} - 3\hat{y}^2 + 2y\hat{y} - y] \quad (18)$$

So, putting all values of eq. 10,eq. 11 and eq. 18 in eq. 9, we get:

$$b_0 \leftarrow b_0 + 2\beta(y - \hat{y})\hat{y}(1 - \hat{y})[1 - 2\gamma x^2(y - \hat{y})(2\hat{y} - 3\hat{y}^2 + 2y\hat{y} - y)] \quad (19)$$

# Advantages of Meta-learning

- **higher model prediction accuracy**
  - optimizing learning algorithms
  - helping learning algorithms better adapt to changes in conditions
- **faster, cheaper training process**
  - learning from fewer examples
  - increase speed of learning processes by reducing necessary experiments
- **more generalized models**

- Source Dataset : MNIST data (first 8 classes)
- Target Dataset : MNIST data (remaining 2 classes)
- Model used : CNNs

# Fine-tuning : Hands-on

- $\hat{y} = F_{\phi}(F_{\theta}(x))$
- Train  $F_{\phi}(F_{\theta}(x))$  for a supervised task with big data
- Discard  $F_{\phi}$  and use  $F_{\phi'}(F_{\theta})$  for the new task
- Freeze  $F_{\theta}$  and update  $F_{\phi'}$
- Update the complete model end-to-end

`https:  
//github.com/madhavlab/2023_codscomad_metalearning`



# Model-agnostic meta-learning: Hands-on

Jupyter Notebook!