Meta Learning: Learn to learn

Hung-yi Lee

What does "meta" mean? meta-X = X about X

Source of image: https://medium.com/intuitionmachine/the-brute-force-method-of-deep-learning-innovation-58b497323ae5 (Denny Britz's graphic)

這門課的作業在做甚麼?



朋友覺得我在

我媽覺得我在

大眾覺得我在



指導教授覺得我在

我以為我在

事實上我在

感謝 沈昇勳 同學提供圖檔

Industry



Academia



Using 1000 GPUs to try 1000 sets of hyperparameters "Telepathize" (通靈) a set of good hyperparameters

Can machine automatically determine the hyperparameters?

Machine Learning 101



Using $\boldsymbol{\theta}$ to represent the unknown parameters.



Machine Learning 101



sum over examples

Introduction of Meta Learning

What is Meta Learning?



What is *learnable* in a learning algorithm?



What is *learnable* in a learning algorithm?







 θ^{1*} parameters of the classifier learned by F_{ϕ} using the training examples of task 1



Evaluate the classifier on testing set



Ground Truth









In typical ML, you compute the loss based on training examples Task 1 In meta, you compute the loss based on testing examples f_{**0}1*</sub>** Hold on! You use testing examples during training??? apple prediction orange Compute difference apple orange

Testing Examples



apple orange

Ground Truth

Task 1In typical ML, you compute the
loss based on training examples
In meta, you compute the loss
based on testing examples
of training tasks.

Testing Examples







- Loss function for learning algorithm $L(\phi) = \sum_{n=1}^{\infty} l^n$
- Find ϕ that can minimize $L(\phi) \qquad \phi^* = \arg \min_{\phi} L(\phi)$
- Using the optimization approach you know If you know how to compute $\partial L(\phi)/\partial \phi$

Gradient descent is your friend.

n=1

```
What if L(\phi) is not differentiable?
```

Reinforcement Learning / Evolutionary Algorithm

Now we have a learned "learning algorithm" F_{ϕ^*}



ML v.s. Meta

Goal

Machine Learning ≈ find a function f

Dog-Cat Classification



 $= f \dots$

Meta Learning

≈ find a function F that finds a function f

 $\begin{array}{c} \text{Learning} \\ \text{Algorithm} \end{array} F$



Machine Learning Training Data **One task** Meta Learning cat dog Train **Training tasks** Task 1 Test Train Apple & apple apple orange orange Orange Task 2 Test Train Car & Bike bike bike car car

(in the literature of "learning to compare")

Support set

Query set





Loss





Meta Learning v.s ML

- What you know about ML can usually apply to meta learning
 - Overfitting on training tasks
 - Get more training tasks to improve performance
 - Task augmentation
 - There are also hyperparameters when learning a learning algorithm
 - Development task 😳

What is the objects act of the objects act of the objects for any context set of the o

int("please select exaction

mirror_mod = modifier_ob mirror object to mirror mirror_mod.mirror_object

peration == "MIRROR_X": irror_mod.use_x = True irror_mod.use_y = False

X mirror to the select ject.mirror_mirror_x" ror X"

ontext): oxt.active_object is not



Review: Gradient Descent

Learning to initialize

Model-Agnostic Meta-Learning (MAML)



Chelsea Finn, Pieter Abbeel, and Sergey Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks", ICML, 2017

• Reptile



https://arxiv.org/abs/1803.02999

How to train your Dragon MAML

Strided MAML vs Strided MAML++



Antreas Antoniou, Harrison Edwards, Amos Storkey, How to train your MAML, ICLR, 2019



Testing Task

dog



Pre-training (Self-supervised Learning)



Trained by proxy tasks (fill-in the blanks, etc.)



Isn't it domain adaptation / transfer learning?

Task 1



cat





cat

dog

find good init



Pre-training (more typical ways)



Use data from different tasks to train a model

Also known as multi-task learning (baseline of meta)

MAML v.s. Pre-training

https://youtu.be/vUwOA3SNb_E



這就是 "meta 業配"



MAML is good because

• ANIL (Almost No Inner Loop)



Aniruddh Raghu, Maithra Raghu, Samy Bengio, Oriol Vinyals, Rapid Learning or Feature Reuse? Towards Understanding the Effectiveness of MAML, ICLR, 2020

More about MAML

- More mathematical details behind MAML
 - https://youtu.be/mxqzGwP_Qys
- First order MAML (FOMAML)
 - https://youtu.be/3z997JhL9Oo
- Reptile
 - https://youtu.be/9jJe2AD35P8



Optimizer

Marcin Andrychowicz, et al., Learning to learn by gradient descent by gradient descent, NIPS, 2016





$$\widehat{\phi} = \arg\min_{\phi} L(\phi) \qquad \nabla_{\phi} L(\phi) =?$$
Network
Architecture

- Reinforcement Learning
 - Barret Zoph, et al., Neural Architecture Search with Reinforcement Learning, ICLR 2017
 - Barret Zoph, et al., Learning Transferable Architectures for Scalable Image Recognition, CVPR, 2018
 - Hieu Pham, et al., Efficient Neural Architecture Search via Parameter Sharing, ICML, 2018

An agent uses a set of actions to determine the network architecture.

 ϕ : the agent's parameters

 $-L(\phi)$

Reward to be maximized



Within-task Training



- <u>Reinforcement Learning</u>
 - Barret Zoph, et al., Neural Architecture Search with Reinforcement Learning, ICLR 2017
 - Barret Zoph, et al., Learning Transferable Architectures for Scalable Image Recognition, CVPR, 2018
 - Hieu Pham, et al., Efficient Neural Architecture Search via Parameter Sharing, ICML, 2018
- Evolution Algorithm
 - Esteban Real, et al., Large-Scale Evolution of Image Classifiers, ICML 2017
 - Esteban Real, et al., Regularized Evolution for Image Classifier Architecture Search, AAAI, 2019
 - Hanxiao Liu, et al., Hierarchical Representations for Efficient Architecture Search, ICLR, 2018



• DARTS Hanxiao Liu, et al., DARTS: Differentiable Architecture Search, ICLR, 2019





Data Augmentation



Yonggang Li, Guosheng Hu, Yongtao Wang, Timothy Hospedales, Neil M. Robertson, Yongxin Yang, DADA: Differentiable Automatic Data Augmentation, ECCV, 2020

Daniel Ho, Eric Liang, Ion Stoica, Pieter Abbeel, Xi Chen, Population Based Augmentation: Efficient Learning of Augmentation Policy Schedules, ICML, 2019 Ekin D. Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, Quoc V. Le, AutoAugment: Learning Augmentation Policies from Data, CVPR, 2019

Sample Reweighting

• Give different samples different weights



Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, Deyu Meng, Meta-Weight-Net: Learning an Explicit Mapping For Sample Weighting, NeurIPS, 2019 Mengye Ren, Wenyuan Zeng, Bin Yang, Raquel Urtasun, Learning to Reweight Examples for Robust Deep Learning, ICML, 2018

Beyond Gradient Descent

Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, Raia Hadsell, Meta-Learning with Latent Embedding Optimization, ICLR, 2019

This is a Network. Its parameter is ϕ (Invent new learning algorithm! Not gradient descent anymore) Training Training Data Data

 $\boldsymbol{\theta}^*$

How about? Until now cat cat Learning Learning + Classification **H*** Algorithm (Function F) (Function F) dog cat dog cat **Testing Data Training Data Testing Data** Training Data https://youtu.be/yyKaACh_j3M Learning to compare

(metric-based approach)

https://youtu.be/yyKaACh_j3M https://youtu.be/scK2EIT7klw https://youtu.be/semSxPP2Yzg https://youtu.be/ePimv_k-H24



Applications

Few-shot Image Classification

• Each class only has a few images.



- N-ways K-shot classification: In each task, there are N classes, each has K examples.
- In meta learning, you need to prepare many N-ways K-shot tasks as training and testing tasks.

Omniglot

https://github.com/brendenlake/omniglot

- 1623 characters
- Each has 20 examples



ゴムジアルドリーカマネるもロベッシッシュえひにんた 3 th 8 50 V 1221ᢓᡜᡛᡦᡛᡱᡅಐ᠊᠋ぉ᠄᠈ᡆ᠆᠆᠆᠋᠋ᡣᡄ᠓ᡃ᠋ᡣ᠙᠔᠙᠙᠙᠋᠋᠋᠋ᠴ᠇ᡕᠻᠺ᠕ᠢᢃ᠋ᠲ║ᢃᠲ ╡╉ҘҨ҄ҋѡӹѡぉ҄ѯ҆҆҆ҽѿҥҵҧҧѠ҄҄ѡҧӷӀ҄҄҄҄҄҄҄҄ӹҝҧҲ <u>ਵਿਚਅਉ</u>ਬੈਂਡੇਡੇಫੇಫੇਸ਼ਦਸ਼ਦਦਦਾ ወጀ ጭ ஸි வு ァ ァ ッ ፵ ፵ ዓ ፋ ェ ር ዀ ጄ ዓ 8 ቮ ኹ 日日日日のや、ロノチョ町の日であるしているのでのまで、ロンジのスンンシ BOJEBOJ8 a montre preservisor a contre servisor a contre preservisor a contre servisor a contre servis 5 4 主 Los, ∞ N N N A G @ S f D & v L 20 5 G Z U L W H G d W Y J K, VDJPXYYNOZEBJ= 20MUTTOVPLCUUUU ア m m m 、 NHX P 1 米 ベタ 米 広 チ m m m m m m s E E K 4 B N W W や U ビ A A A A

Omniglot

Demo: https://openai.com/blog/reptile/



- Split your characters into training and testing characters
 - Sample N training characters, sample K examples from each sampled characters → one training task
 - Sample N testing characters, sample K examples from each sampled characters → one testing task

	(A) Learning to initialize	(B) Learning to compare	(C) Other
Sound Event Detection	(Shi et al., 2020)	(Shimada et al., 2020a) (Chou et al., 2019) (Wang et al., 2020) (Shimada et al., 2020b) (Shi et al., 2020)	Network architecture search: (Li et al., 2020)
Keyword Spotting	(Chen et al., 2020a)	(Huh et al., 2020)	Net2Net: (Veniat et al., 2019) Network architecture search: (Mazzawi et al., 2019) Network architecture search: (Mo et al., 2020)
Text Classification	(Dou et al., 2019) (Bansal et al., 2019)	(Yu et al., 2018) (Tan et al., 2019) (Geng et al., 2019) (Sun et al., 2019)	Learning the learning algorithm: (Wu et al., 2019)
Voice Cloning			Learning the learning algorithm: (Chen et al., 2019b) (Serrà et al., 2019)
Sequence Labelng	(Wu et al., 2020)	(Hou et al., 2020)	
Machine Translation	(Gu et al., 2018) (Indurthi et al., 2020)		
Speech Recognition	(Hsu et al., 2020) (Klejch et al., 2019) (Winata et al., 2020a) (Winata et al., 2020b)		Learning to optimize: (Klejch et al., 2018) Network architecture search: (Chen et al., 2020b) (Baruwa et al., 2019)
Knowledge Graph	(Obamuyide and Vlachos, 2019) (Bose et al., 2019) (Lv et al., 2019) (Wang et al., 2019)	(Ye and Ling, 2019) (Chen et al., 2019a) (Xiong et al., 2018) (Gao et al., 2019)	
Dialogue / Chatbot	(Qian and Yu, 2019) (Madotto et al., 2019) (Mi et al., 2019)		Learning to optimize: (Chien and Lieow, 2019)
Parsing	(Guo et al., 2019) (Huang et al., 2018)		
Word Embedding	(Hu et al., 2019)	(Sun et al., 2018)	
Multi-model		(Eloff et al., 2019)	Learning the learning algorithm: (Surís et al., 2019)

http://speech.ee. ntu.edu.tw/~tlkag k/meta_learning_ table.pdf