Meta Learning and Its Applications to Human Language Processing

Hung-yi Lee, Ngoc Thang Vu, Shang-Wen (Daniel) Li









Why Meta Learning?

What does "meta" mean?

meta-X = X about X



Hung-yi Lee (NTU)





Meme



Meta Meme

Why Meta Learning?



What does "meta" mean?

Hung-yi Lee (NTU)

meta-X = X about X

- Meta Learning is "learning about learning", or "learn to learn".
- How to reduce the requirement of annotated data for language and speech processing?
- Can machines find learning algorithms that only require limited labeled data?

Why Meta Learning?



- Thang Vu A professor at the Institute for Natural Language Processing (IMS), University of Stuttgart, Germany
- My reasons:
 - Better understanding on results which I published during my PhD on *mutlilingual speech recognition, e.g.* [1]
 - A potential solution for low resource settings in speech and language processing applications

[1] Ngoc Thang Vu, Wojtek Breiter, Florian Metze, Tanja Schultz. Initialization Schemes for Multilayer Perceptron Training and their Impact on ASR Performance using Multilingual Data. Interspeech, 2012

Meta-learning for ConvAl

how long does it take to go to <mark>enloe</mark> hospital?

It will take 20 minutes for driving





- ConvAl
- Impact on industry
 - many to-B, to-C, and to-D applications
- Language, domain, and task expansion
 - Meta-learning is a scalable solution for industry need

Machine Learning 101



Using θ to represent the learnable parameters.



Machine Learning 101



loss:
$$l(\theta) = \sum_{k=1}^{K} d_k$$
 sum over
examples

$$\hat{\theta} = \arg\min_{\theta} l(\theta)$$

done by gradient descent

 $f_{\widehat{\theta}}$ is the function learned by learning algorithm from data

Introduction of Meta Learning

What is Meta Learning?



What is *learnable* in a learning algorithm?



What is *learnable* in a learning algorithm?







 $\hat{\theta}^1$: parameters of the classifier learned by F_{ϕ} using the training examples of task 1



Evaluate the classifier on testing set











In typical ML, you compute the loss based on training examples Task 1 In meta, you compute the loss based on testing examples Hold on! You use testing examples during training???



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 $L(\phi)$ $\hat{\phi} = \arg\min_{\phi} L(\phi)$
- 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" $F_{\hat{\phi}}$



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





Machine Learning


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



Step 2 – Loss Function





How to compute $\nabla_{\phi} l$ (ⁿ is ignored here)

 ϕ_i : the i-th parameter of ϕ





Can be computationally intensive ...





Here within-task training only uses <u>one step</u> ...

$$\hat{\theta} = \phi - \varepsilon \nabla_{\phi} l(\phi)$$

- Fast ... Fast ... Fast ...
- Good to train a model with one step. \bigcirc
- In few-shot learning scenario, we have limited training examples in each task.
 → only update once to prevent overfitting
- Within-task training is different in across-task training and testing

In across-task **training**, within-task training <u>updates once</u> In across-task **testing**, within-task training <u>updates many times</u>

$$\frac{\partial l}{\partial \phi_i} = \sum_j \frac{\partial l}{\partial \hat{\theta}_j} \frac{\partial \hat{\theta}_j}{\partial \phi_i}$$

Here within-task training only uses <u>one step</u> ...

$$\hat{\theta} = \phi - \varepsilon \nabla_{\phi} \, l(\phi)$$



How to train your Dragon MAML

Strided MAML vs Strided MAML++



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

One "Initialization" fits all?

• Task conditioning: Give different tasks different initialization.



Huaxiu Yao, Ying Wei, Junzhou Huang, Zhenhui Li, Hierarchically Structured Meta-learning, ICML, 2019

Initialization of "Learn to initialize"



Turtles all the way down?

- MAML learns the initialization parameter ϕ by gradient descent
- What is the initialization parameter ϕ^0 for ϕ ?
 - Gradient descent
 - Learn to initialize
 - Learn to learn to initialize

MAML is good because

• Rapid Learning or Feature Reuse?



Aniruddh Raghu, Maithra Raghu, Samy Bengio, Oriol Vinyals, Rapid Learning or Feature Reuse? Towards Understanding the Effectiveness of MAML, ICLR, 2020 Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, Chelsea Finn, "Meta-Learning without Memorization", ICLR, 2020

Ignoring Training Examples?



correct answers \rightarrow orange apple

For the models simply learn what apple & orange look like (not learn to learn), they cannot obtain correct answers.

If you want more update steps, but don't like intensive computation.



- iMAML: Aravind Rajeswaran, Chelsea Finn, Sham Kakade, Sergey Levine, Meta-Learning with Implicit Gradients, NeurIPS, 2019
- **Reptile:** Alex Nichol, Joshua Achiam, John Schulman, On First-Order Meta-Learning Algorithms, arXiv, 2018
- Pan Zhou, Xiao-Tong Yuan, Huan Xu, Yan Shuicheng, Jiashi FengEfficient Meta Learning via Minibatch Proximal Update, NeurIPS, 2019
- Luca Bertinetto, João F. Henriques, Philip H.S. Torr, Andrea Vedaldi, Meta Learning with Differentiable Closed-Form Solver, ICLR, 2019

More Approaches



Learning Optimizer

Step 1 – What is learnable?







Sachin Ravi, et al., Optimization as a Model for Few-Shot Learning, ICLR, 2017



(a) Forget gate values for 1-shot meta-learner

(b) Input gate values for 1-shot meta-learner

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



- 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
- 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

Learning as a Network?

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)



 $\widehat{ heta}$



Learning to Compare

Training

Meta Learning

Training tasks



(in the literature of "learning to compare")

Training

Meta Learning

Training tasks



Testing

Meta Learning



Learning to Compare

- What is the learned *learning algorithm* in this case?
- Think about <u>non parametric models</u> such as k-nearest neighbors
 - All training data are stored \implies no learning needed
 - Performance depends on the distance/similarity metrics
- 'Learning to compare' algorithms
 - learn such models
 - do not have the within-task training
 - make the metrics *trainable* across tasks
First Example: Siamese Network

Koch, Zemel, Salakhutdinov, 2015



First Example: Siamese Network

Koch, Zemel, Salakhutdinov, 2015





Frame It as a Meta Learning Setting Network Test Train Yes Training Test Yes Train Network Tasks Train Test No Network Yes Testing Train Test or **Tasks**

No

Yes

No

Matching Network

Vinyals, Blundell, Lillicrap, Kavukcupglu, Wierstra, 2017



Prototypical Network



Relation Network

Sung, Yang, Zhang, Xiang, Torr, Hospedales, 2018



Meta Learning vs. Multi-task Learning vs. Transfer Learning

Meta Learning vs. Multi-task Learning

- Both use training data from many different tasks but have different objectives
- Meta learning aims at improving the accuracies of future tasks while multi-task learning optimizes the accuracies on all existing tasks
- The more tasks, the better the meta model, while multi-task learning methods might have problems with a large number of tasks

Meta Learning vs. Transfer Learning

- The goals are similar: improving accuracies on future new tasks
- While meta learning focuses on improving the training algorithms for future tasks, transfer learning aims at re-using knowledge learnt from previous tasks
- Meta learning assumes the same distribution between training tasks and testing tasks while transfer learning does not assume it between previous tasks and future tasks



Speech Recognition

Fast Adapt to Unseen Languages

Training Tasks

Testing Tasks



Fast Adapt to Unseen Languages

Learning to Initialize (MAML)

Jui-Yang Hsu, Yuan-Jui Chen, Hung-yi Lee, META LEARNING FOR END-TO-END LOW-RESOURCE SPEECH RECOGNITION, ICASSP, 2020

- Data from the IARPA BABEL project
 - Training tasks: Bengali, Tagalog, Zulu, Turkish, Lithuanian, Guarani
 - Testing tasks: Vietnamese, Swahili, Tamil, Kurmanji
- ASR model: CTC

CER (%)	Vietnamese	Swahili	Tamil	Kurmanji
Rand Init	71.8	47.5	69.9	64.3
Pre-train	59.7	48.8	65.6	62.6
MAML	50.1	42.9	58.9	57.6

Fast Adapt to Unseen Languages

Learning to Initialize (MAML)

Jui-Yang Hsu, Yuan-Jui Chen, Hung-yi Lee, META LEARNING FOR END-TO-END LOW-RESOURCE SPEECH RECOGNITION, ICASSP, 2020

- Data from the IARPA BABEL project
 - Training tasks: Bengali, Tagalog, Zulu, Turkish, Lithuanian, Guarani
 - Testing tasks: Vietnamese, Swahili, Tamil, Kurmanji
- ASR model: CTC

CER (%)	Vietnamese	Swahili	Tamil	Kurmanji
Rand Init	71.8	47.5	69.9	64.3
Pre-train	59.7	48.8	65.6	62.6
MAML	50.1	42.9	58.9	57.6

Fast Adapt to Unseen Languages Network Architecture Search (DARTS)

Yi-Chen Chen, Jui-Yang Hsu, Cheng-Kuang Lee, Hung-yi Lee, "DARTS-ASR: Differentiable Architecture Search for Multilingual Speech Recognition and Adaptation", INTERSPEECH, 2020

- Data from the IARPA BABEL project
 - Monolingual: Training and testing tasks are the same language
 - oracle setting in meta learning
 - *Multilingual*: Training and testing tasks are different languages
- ASR model: CTC



CER (%)	Vietnamese	Swahili	Tamil	Kurmanji
VGG-small	45.3	36.3	55.7	54.5
VGG-large	43.2	36.1	55.0	55.1
DARTS-ASR	40.9	32.3	45.9	53.5

(the network architecture can also be fine-tuned in the testing task)

VGG (not only useful in image, but speech related applications)



Fast Adapt to Unseen Accents



Fast Adapt to Unseen Accents

Learning to Initialize (MAML)

Genta Indra Winata, Samuel Cahyawijaya, Zihan Liu, Zhaojiang Lin, Andrea Madotto, Peng Xu, Pascale Fung, Learning Fast Adaptation on Cross-Accented Speech Recognition, INTERSPEECH, 2020

- Data from CommonVoice Dataset
- ASR model: Seq2seq



Fast Adapt to Unseen Speakers

Training Tasks

Testing Tasks



l can learn Speaker X better

Speaker Adaptive Training?

Yes. New approaches for speaker adaptive training.

Fast Adapt to Unseen Speakers

Learning to Initialize (MAML)

Ondřej Klejch, Joachim Fainberg, Peter Bell, Steve Renals, Speaker Adaptive Training using Model Agnostic Meta-Learning, ASRU, 2019

Huaxin Wu, Genshun Wan, Jia Pan, Speaker Code Based Speaker Adaptive Training Using Model Agnostic Meta-learning, INTERSPEECH, 2020

• Learning Optimizer (optimizer as RNN)

Ondřej Klejch, Joachim Fainberg, Peter Bell, Learning to adapt: a meta-learning approach for speaker adaptation, INTERSPEECH, 2018



More

Krsto Proroković, et al., Adaptation of an EMG-Based Speech Recognizer via Meta-Learning, GlobalSIP, 2019

• Speech recognition from EMG signal of facial muscles



- Session = recording between putting and removing the electrodes
- Adaptation is needed for each session
- Using learning to initialize (MAML)

More

Speech Translation

Sathish Indurthi, et al., Data Efficient Direct Speech-to-Text Translation with Modality Agnostic Meta-Learning, ICASSP 2020



Testing Task: Speech Translation

Code Switching

Genta Indra Winata, Samuel Cahyawijaya, Zhaojiang Lin, Zihan Liu, Peng Xu, Pascale Fung, Meta-Transfer Learning for Code-Switched Speech Recognition, ACL, 2020



Speaker Verification

(Sound Event Detection, Keyword Spotting)

Speaker Verification



Speaker Verification







This is "Learning to Compare"!

Framework

The speakers in <u>stages 2</u> and 3 are not seen in stage 1.

Testing Task (Episode)





Across-task Training: Generating training tasks





Across-task Training: Generating training tasks





Across-task Training: Generating training tasks



Also refer to generalized end-to-end (GE2E)

Li Wan, et al., Generalized End-to-End Loss for Speaker Verification, ICASSP, 2018

Sound Event Detection

Find this kind of sound event





Learning to Compare in Audio



"Attention" is used to gather useful information in audio clips

Source of image: https://arxiv.org/pdf/1812.01269.pdf
Attentional Similarity



Szu-Yu Chou, Kai-Hsiang Cheng, Jyh-Shing Roger Jang, Yi-Hsuan Yang, Learning to match transient sound events using attentional similarity for few-shot sound recognition, ICASSP, 2019

Attention Mechanism

Chowdhury, et al., Attention-Based Models for Text-Dependent Speaker Verification, ICASSP, 2018

NetVLAD

Xie et al., Utterance-level Aggregation For Speaker Recognition In The Wild, ICASSP, 2019

VLAD = Vector of Locally Aggregated Descriptors



More Speaker Verification

- More Learning to compare
 - Georg Heigold, et al., End-to-End Text-Dependent Speaker Verification, ICASSP, 2016
 - Li Wan, et al., Generalized End-to-End Loss for Speaker Verification, ICASSP, 2018
 - Jixuan Wang, et al., Centroid-based deep metric learning for speaker recognition, ICASSP, 2019
 - Tom Ko, , et al., Prototypical Networks for Small Footprint Text-Independent Speaker Verification, ICASSP, 2020
 - Seong Min Kye, et al., Meta-Learning for Short Utterance Speaker Recognition with Imbalance Length Pairs, INTERSPEECH 2020
 - Joon Son Chung, et al., In defence of metric learning for speaker recognition, INTERSPEECH 2020

More Speaker Verification

- Also architecture search
 - Shaojin Ding, et al., AutoSpeech: Neural Architecture Search for Speaker Recognition, INTERSPEECH, 2019
- Learn to initialize
 - Jiawen Kang, Ruiqi Liu, Lantian Li, Dong Wang, Thomas Fang Zheng, Domain-Invariant Speaker Vector Projection by Model-Agnostic Meta-Learning, INTERSPEECH, 2020

More Sound Event Detection

- More Learning to compare
 - Pranay Manocha, et al., Content-based Representations of audio using Siamese neural networks, ICASSP, 2018
 - Kazuki Shimada, et al., Metric Learning with Background Noise Class for Few-shot Detection of Rare Sound Events, ICASSP 2020
 - Yu Wang, et al., Few-Shot Sound Event Detection, ICASSP, 2020
 - Bowen Shi, Ming Sun, Krishna C. Puvvada, Chieh-Chi Kao, Spyros Matsoukas, Chao Wang, Few-Shot Acoustic Event Detection Via Meta Learning, ICASSP 2020
- Also architecture search
 - Jixiang Li, et al., Neural Architecture Search on Acoustic Scene Classification, INTERSPEECH, 2020

Keyword Spotting / Intent Detection

- Learn to compare
 - Ashish Mittal, Samarth Bharadwaj, Shreya Khare, Saneem Chemmengath, Karthik Sankaranarayanan, Brian Kingsbury, Representation based meta-learning for few-shot spoken intent recognition, INTERSPEECH, 2020
- Learning to initialize
 - Yangbin Chen, et al., An Investigation of Few-Shot Learning in Spoken Term Classification, INTERSPEECH, 2020
- Also architecture search
 - Tom Véniat, et al., Stochastic Adaptive Neural Architecture Search for Keyword Spotting, ICASSP, 2019
 - Hanna Mazzawi, et al., Improving Keyword Spotting and Language Identification via Neural Architecture Search at Scale, INTERSPEECH, 2019
 - Bo Zhang, et al., UTOKWS: KEYWORD SPOTTING WITH DIFFERENTIABLE ARCHITECTURE SEARCH, arXiv, 2020
 - Tong Mo, et al., Neural Architecture Search For Keyword Spotting, INTERSPEECH, 2020

Emotion Recognition

- Sample weighting
 - Takuya Fujioka, Takeshi Homma, Kenji Nagamatsu, META-LEARNING FOR SPEECH EMOTION RECOGNITION CONSIDERING AMBIGUITY OF EMOTION LABELS, INTERSPEECH 2020
- Learn to initialize
 - Ruichu Cai, Kaibin Guo, Boyan Xu, Xiaoyan Yang, Zhenjie Zhang, Meta Multi-task Learning for Speech Emotion Recognition, INTERSPEECH 2020

- Intent classification
- Named-entity recognition (NER), slot tagging
- Dialog state tracking
- Question answering
- Dialog generation



Intent classification

How are you today?	
Phatic	Pretty good. How can I help you?
I want a Hotel of moderate pric	e
BookHotel	Which city are you looking for?
I am thinking of Boston	
BookHotel	Sure, there are 5 hotels in Boston moderately priced. Is there any specific amenity you are looking for?
I would like to celebrate my birth	day
Out of Domain	Sorry, could you be more specific?
Which hotel provides lounge for p	arty
InquireAmenity	Hotel XYZ has that

NER, Slot tagging



Dialog state tracking



Question Answering



Dialog generation

How are you today?		
	Pretty good. How can I help you?	
I want a Hotel of moderate price		
	Which city are you looking for?	Context
I am thinking of Boston		Querv
Ū		

- Zero-shot OOD detection and few-shot ID classification
- OOD-resistant Prototypical Network

Intent Label	Example
Help_List	List what you can help me with.
	Watson, I need your help
Schedule_Appointment	Can you book a cleaning
	with my dentist for me?
	Can you schedule my den-
	tist's appointment?
End_Meeting	You can end the meeting
	now
	Meeting is over
	•••
OOD utterances	My birthday is coming!
	blah blah

Tan et al. "Out-of-Domain Detection for Low-Resource Text Classification Tasks" EMNLP 2019

• Across task training

Sample Intent *in, out,* and 4 *others* from D^{across_train} 5-way (in + 4)

Task 2

Task 1

Sample Intent *in, out,* and 4 *others* from D^{across_train}

Task N

...



$$\mathcal{L}_{in} = -\log \frac{\exp \alpha F(x_i^{in}, S_{l_i}^{in})}{\sum_{l'} \exp \alpha F(x_i^{in}, S_{l'}^{in})}$$

$$\mathcal{L}_{ood} = max[0, \max_{l}(F(x_j^{out}, S_l^{in}) - \mathcal{M}_1)]$$

$$\mathcal{L}_{gt} = max[0, \mathcal{M}_2 - F(x_i^{in}, S_{l_i}^{in}))]$$



- Dataset: Amazon review, conversation
- Metrics: intent classification error rate at EER

	Amazon review	Conversation
LSTM AutoEnc.	38.6	79.5
CNN	42.8	77.6
O-Proto (proposed)	29.1	40.8

 Meta-learning (Proto) > Supervised learning (CNN / / LSTM)

X	l am	thinking	of Boston
У	00	0	O B _{City}

$$p(\boldsymbol{y} \mid \boldsymbol{x}, \mathcal{S}) = \frac{1}{Z} \exp(\operatorname{TRANS}(\boldsymbol{y}) + \lambda \cdot \operatorname{EMIT}(\boldsymbol{y}, \boldsymbol{x}, \boldsymbol{S})),$$

$$Z = \sum_{\boldsymbol{y}' \in \boldsymbol{Y}} \exp(\operatorname{TRANS}(\boldsymbol{y}') + \lambda \cdot \operatorname{EMIT}(\boldsymbol{y}', \boldsymbol{x}, \boldsymbol{S})),$$

$$\operatorname{TRANS}(\boldsymbol{y}) = \sum_{i=1}^{n} f_T(y_{i-1}, y_i) \quad f_T(y_{i-1}, y_i) = p(y_i \mid y_{i-1}).$$

Hou et al. "Few-shot Slot Tagging with Collapsed Dependency Transfer and Label-enhanced Task-adaptive Projection Network" ACL 2020

• Abstract labels: O, sB, dB, sI and dI





compute transition over tasks





- SNIPS
 - 5 for across-task training, 1 for dev, 1 for across-task test
 - 5-shot

	We	Mu	Pl	Во	Se	Re	Cr	Avg
BiLSTM	25.2	39.8	46.1	74.6	53.5	40.4	25.1	43.5
SimBERT	53.5	54.1	42.8	75.5	57.1	55.3	32.4	53.0
TransferBERT	59.4	42.0	46.1	20.7	28.2	67.8	58.6	46.1
Meta-learning	67.8	56.0	46.0	72.2	73.6	60.2	66.9	63.2
Meta-learning + transition	74.7	56.7	52.2	78.8	80.6	69.6	67.5	68.6

- Meta-learning > transfer learning > no transfer
- Transition helps

Cross-lingual NER

- Goal: source -> target
- Previously: annotation projection
- Meta-learning for language-independent features
 - On top of multilingual BERT
 - Further improvement

Wu et al. "Enhanced Meta-Learning for Cross-lingual Named Entity Recognition with Minimal Resources " AAAI 2020

Cross-lingual NER

• Across task training $\begin{array}{c} Guarantee different \\ \hline \\ Train Examples x_{j} & Test x_{i} \\ \hline \\ Train Examples from D^{s} = \{x_{i}\}_{i=1}^{N} \end{array}$

MAML – model: M

Across task testing

 $\frac{Task \ 1...N}{for \ x_i \ in \ D^T}$



Cross-lingual NER

- English-en / German-de / Spanish-es / Dutch-nl, (CoNLL);
 French-fr, Europeana News; Chinese-zh , MSRA
- Results (F1-scores)

	es	nl	de	fr	zh	Avg
Ni, Dinu, and Florian (2017)	65.1	65.4	58.5	-	-	-
Mayhew, Tsai, and Roth (2017)	66.0	66.5	59.1	-	-	-
Xie et al. (2018)	72.4	71.3	57.8	-	-	-
Multilingual BERT	74.6	79.6	70.8	50.9	76.4	70.5
Meta-cross-lingual NER	76.8	80.4	73.2	55.3	77.9	72.7

• Meta-learning > multilingual BERT > annotation projection



(BookHotel, City = Boston)

Dialog frame



Algorithm 2 MAML algorithm **Input:** D_d^{train} ; D_d^{valid} ; α ; β . **Output:** Trained model M with MAML algorithm. 1: while not done do for each domain d do 2: Select a batch of size from D_d^{train} and 3: Support set $\leftarrow D_d^{valid}$ to get D_d^t and D_d^v ; $\leftarrow D_d^{valid}$ to get D_d^t and D_d^v ; Pre-update model with gradient descent: → Query set $M'_d \leftarrow M - \alpha \nabla_M L_d(M, D^t_d) \longrightarrow \mathsf{Pre-update}$ Compute $L_d(M'_d, D^v_d)$ using D^v_d ; 5: end for 6: 7: Update the current model M: $M \leftarrow M - \beta \nabla_M \sum_d L_d(M'_d, D^v_d) \longrightarrow$ Update meta-learn model M 8: end while 9: **return** meta-learned model M;

> Huang et al. "Meta-Reinforced Multi-Domain State Generator for Dialogue Systems" ACL 2020

- MultiWOZ
 - Restaurant, hotel, train (source)
 - Taxi / attrachtion (new)
- Results

New Domain (proportion)	Initialization	Joint Acc.	Slot Acc.
	Public BERT	60.6	73.3
Taxi (1%)	Supervised-learning	59.0	78.7
	Meta-learning	64.4	83.2
	Public BERT	27.9	63.4
Attraction (1%)	Supervised-learning	29.1	62.2
	Meta-learning	43.1	74.3

 Meta-learning > Supervised with source domain >= BERT (no pre-training with matched tasks)





Yan et al. "Multi-source Meta Transfer for Low Resource Multiple-Choice Question Answering" ACL 2020



17 Get all batches of data $\tau_i^t \sim p^t(\tau)$; 18 for all τ_i^t do 19 Evaluate $\nabla_{\theta} L_{\tau_i^t}(f(\theta))$ with respect to batch size; 20 Gradient for meta transfer learning: $\theta = \theta - \lambda \nabla_{\theta} L_{\tau_i^t}(f(\theta))$; 21 end

Adaptation

Name	DREAM	RACE	MCTEST	SemEval	SWAG
Туре	Dialogue	Exam	Story	Narrative Text	Scenario Text
Ages	15+	12-18	7+	-	-
Generator	Expert	Expert	Crowd.	Crowd.	AF./Crowd.
Level	High School/College	High/Middle School	Children	Unlimited	Unlimited
Choices	3	4	4	2	4
Samples	6,444	27,933	660	2,119	92,221
Questions	10,197	97,687	2,640	13,939	113,557

• Source: RACE + SWAG, target: McTEST

Method	Target labels	Acc.
Supervised-learning (BERT-base)	Yes	68.0
IMC [1]	Yes	76.6
Transfer-learning (BERT-base)	No	79.2
Meta-learning (BERT-base)	Νο	81.6

• Meta-learning > previous SOTA and simple transfer learning

• Target: DREAM/McTEST/SemEval, source: rest 4

Method	DREAM	MCTEST	SemEval
GPT+Strategies [2]	-	81.9	89.5
IMC [1]	-	76.6	-
XLNet [3]	72.0	-	-
Supervised-learning (BERT-base)	61.6	68.0	87.5
Supervised-learning (RoBERTa)	84.4	87.3	94.0
Meta-learning (BERT-base)	68.9	82.0	88.9
Meta-learning (RoBERTa)	85.6	88.8	94.2

Meta-learning > previous SOTA and supervised learning

Dialog generation

How are you today?		
	Pretty good. How can I help you?	+
I want a Hotel of moderate price		
	Which city are you looking for?	Context
I am thinking of Boston		Query
Dialog generation

- Problem is complicated
- Data collection is costly
- Domain adaptation
- Meta-learning

Qian et al. "Domain Adaptive Dialog Generation via Meta Learning " ACL 2019

Dialog generation

- MAML model: M
- Across task training

Task 1...N

dialog d

Support set

for domain D in source domains sample 1 dialog *d* from D

Test

dialog d

Query set





Train

$$h = \text{Encoder}(B_{t-1}, R_{t-1}, U_t)$$

- $B_t = BspanDecoder(h)$
- $R_t = \text{ResponseDecoder}(h, B_t, m_t)$

Dialog generation

- SimDial
 - Source: *restaurant, weather* and *bus*
 - Target: movie
 - (within-task) train, valid, test: 900, 100, 500 dialogs
 - Across-task training: train of sources
 - Across-task testing:
 - 9 (1%) dialogs from train of target -> support set
 - Test of target -> query set

Method	Entity F1
Transfer learning	64.0
ZSDG	52.6
Meta-learning	66.2

Personalized Dialog generation

Persona 2

I am an artist I have four children I recently got a cat I enjoy walking for exercise I love watching Game of Thrones

[PERSON 1:] Hi
[PERSON 2:] Hello ! How are you today ?
[PERSON 1:] I am good thank you , how are you.
[PERSON 2:] Great, thanks ! My children and I were just about to watch Game of Thrones.
[PERSON 1:] Nice ! How old are your children?
[PERSON 2:] I have four that range in age from 10 to 21. You?
[PERSON 1:] I do not have children at the moment.
[PERSON 2:] That just means you get to keep all the popcorn for yourself.
[PERSON 1:] And Cheetos at the moment!
[PERSON 2:] Good choice. Do you watch Game of Thrones?
[PERSON 1:] No, I do not have much time for TV.
[PERSON 2:] I usually spend my time painting: but, I love the show.

Given dialog context, persona, responses from counterpart (1), generate response (2) consistent with persona

Personalized Dialog generation

- MAML model: M
- Across task training



Madotto et al. "Personalizing Dialogue Agents via Meta-Learning " ACL 2019

Personalized Dialog generation

• Persona-chat dataset

	Coherence	Fluency	Consistency
Human	0.33	3.43	0.23
Supervised	-0.03	-	-
Supervised+Persona	0.07	3.05	0.01
Fine-tuning+Persona	0.00	3.10	0.04
Meta-learning+Persona	0.20	3.19	0.20

- Meta-learning > fine-tuning > no fine-tuning
- Using Persona is better

Conclusion

- Meta-learning in ConvAl
 - transferring knowledge across user intents, domains, and languages
 - Scalable solution for expansion

Meta Learning for NLP

Meta Learning in NLP

- Text classification
- Relation classification

 a special type of text classification
- Sequence to sequence related tasks such as machine translation and semantic parsing
- Knowledge graphs related tasks

Text Classification

• Problem settings:

Input: a piece of text, e.g. sentences



Output: a label

- Some examples in a sentiment analysis task:
 - This movie is great
 positive
- Modern NLP methods:
 - Leverage word embeddings (word
 vectors)
 - Input text → matrices

Diverse Few-Shot Text Classification with Multiple Metrics

- Argued that in previous work, low variants among tasks
 not realistic
 In a more realistic setting, tasks are diverse
- Key ideas and take-home messages:
 - Based on metrics based methods
 - Two steps: 1) tasks clustering; 2) metrics-based
 - Extend meta learning that allows combining multiple metrics depending on different task clusters

Mo Yu, Xiaoxiao Guo, Jinfeng Yi, Shiyu Chang, Saloni Potdar, Yu Cheng, Gerald Tesauro, Haoyu Wang, Bowen Zhou, Diverse Few-Shot Text Classification with Multiple Metrics, ACL 2018

Diverse Few-Shot Text Classification with Multiple Metrics

- How to cluster tasks:
 - Create a transfer performance matrix
 - Apply scores filtering and matrix completion
 - Apply spectral clustering





- How to combine decisions:
 - Linearly combine decisions from different task clusters
 - Linear coefficients are adaptable parameters

$$p(y|x) = \sum_{k} \alpha_{k} P(y|x; f_{k}).$$

Investigating Meta-Learning for Low-Resource NLU Tasks

- Applied MAML and Reptile on a bunch of tasks in GLUE dataset (four high resource tasks for training and four low resource tasks for testing)
- Key take-home messages:
 - Reptile worked best, outperformed other methods
 - Meta learning frameworks outperformed state-of-theart systems trained with a single task and multi-task learning systems
 - It is more effective with less training data

Zi-Yi Dou, Keyi Yu, Antonios Anastasopoulos, Investigating Meta-Learning Algorithms for Low-Resource Natural Language Understanding Tasks, EMNLP 2019

- Key ideas and take-home messages
 - Leverage dynamic routing algorithms (proposed in capsule network – Sabour et al 2017) to improve the generalization of the class representation
 - Leverage the Neural Tensor Network (Socher et al 2013) to compute the relation scores between queries and class vectors
 - Both steps are important and their combination works best

Ruiying Geng, Binhua Li, Yongbin Li, Xiaodan Zhu, Ping Jian, Jian Sun, Induction Networks for Few-Shot Text Classification, EMNLP, 2019







Hierarchical Attention Prototypical Networks for Few-Shot Text Classification

- Key ideas and take-home messages
 - Based on the prototypical network
 - Hierarchical attention architecture
 - Word level attention over words to obtain the sentence representation
 - Instance level attention over instances in the support set to form the prototypes
 - Feature level as proposed in Gao et al AAAI 2019 to improve the distance function

Shengli Sun, Qingfeng Sun, Kevin Zhou, Tengchao Lv, Hierarchical Attention Prototypical Networks for Few-Shot Text Classification, EMNLP 2019

Relation Classification

- As a special case of text classification
- Given a sentence with two marked entities, an NLP system should predict the semantic relation between these two entities
- Examples:
 - The Queen Consort [Jetsun Pema] gave birth to a son on 5 February 2016, [Jigme Namgyel Wangchuck].
 - Relation: *mother*

Model-Agnostic Meta-Learning for Relation Classification with Limited Supervision

- Key ideas and take-home messages
 - Applied MAML for relation classification tasks
 - Utilized two strong models and applied MAML objectives
 - MAML worked very well on their setups





Abiola Obamuyide, Andreas Vlachos, Model-Agnostic Meta-Learning for Relation Classification with Limited Supervision, ACL 2019 Hybrid Attention-Based Prototypical Networks for Noisy Few-Shot Relation Classification

- Key ideas and take-home messages
 - Special design for corrupted text inputs
 - Based on prototypical network
 - Novel method to compute the matching scores based on attention mechanism
 - Hybrid attention:
 - Instance level attention: improves robustness against noisy instances
 - Feature level attention: improves the distance function

Tianyu Gao, Xu Han, Zhiyuan Liu, Maosong Sun, Hybrid Attention-Based Prototypical Networks for Noisy Few-Shot Relation Classification, AAAI 2019

Hybrid Attention-Based Prototypical Networks for Noisy Few-Shot Relation Classification



Multi-Level Matching and Aggregation Network for Few-Shot Relation Classification

- Key ideas and take-home messages
 - Based on matching networks
 - Extend them to multi-level matching and aggregation
 - Local matching
 - Instances matching
 - Class matching

Zhi-Xiu Ye, Zhen-Hua Ling, Multi-Level Matching and Aggregation Network for Few-Shot Relation Classification, ACL 2019

Multi-Level Matching and Aggregation Network for Few-Shot Relation Classification



- 1) Encoder: use a CNN that convert a sentence and the positions of two entities to matrices
- 2) Local matching: use attention method to collect matching information between support instances and the query instance, then use max-pooling and average pooling to convert them to representation vectors for all the support instances and the query instance
- 3) Instance matching: use attention method to compute the prototype
- 4) Class matching: trainable matching scores between the query instance and prototypes

Seq2seq Tasks

- Problem settings:
 - Input: sequence of symbols, e.g. words in one language
 - Output: sequence of symbols, e.g. words in another language
- Examples:
 - Machine translation
 - How are you? \implies Wie geht es Dir?
 - Semantic parsing
 - Input text \implies SQL command

Meta-Learning for Low-Resource Neural Machine Translation

- Key ideas and take-home messages:
 - Apply MAML for multilingual neural translation
 - Propose to use the universal lexical representation to handle the problem of mismatch input and output
 - Utilize 18 translation pairs as training tasks to train the meta learner
 - Fine-tune the meta learnt models for five other translation pairs

Jiatao Gu, Yong Wang, Yun Chen, Kyunghyun Cho, Victor O.K. Li, Meta-Learning for Low-Resource Neural Machine Translation, EMNLP, 2018

Natural Language to Structured Query Generation via Meta-Learning

- Key ideas and take-home messages
 - Map a natural language question to a SQL query
 - Artificially generate pseudo tasks by sampling a batch of training data as a support set and one example as query
 - Design a *relevance function* to find similar examples
 - Relevance function is task dependent
 - E.g. in this paper, the relevance function depends on 1) the predicted SQL type of the input and 2) the input length
 - Apply MAML to train the meta learner

Po-Sen Huang, Chenglong Wang, Rishabh Singh, Wen-tau Yih, Xiaodong He, Natural Language to Structured Query Generation via Meta-Learning, NAACL 2018 Coupling Retrieval and Meta-Learning for Context-Dependent Semantic Parsing

- Key ideas and take-home messages
 - Given a natural language, generate a source code conditioned on the class environment
 - Similar setup as previous paper
 - Introduce a *context aware retriever* to dynamically collect examples from the training as supporting evidences
 - Apply MAML to train the meta learner

Daya Guo, Duyu Tang, Nan Duan, Ming Zhou, Jian Yin, Coupling Retrieval and Meta-Learning for Context-Dependent Semantic Parsing, ACL, 2019

Coupling Retrieval and Meta-Learning for Context-Dependent Semantic Parsing



The retriever finds top-K nearest examples based on the following distance:

$$distance = KL(p(z|x, c)||p(z|x', c'))$$
$$= KL(p(z_x|x)||p(z_x|x'))$$
$$+ KL(p(z_c|c)||p(z_c|c'))$$

Knowledge Graphs

- Problem settings:
 - Knowledge graph, i.e. a collection of triples (h, r, t)
 h: head entity; r: relation; t: tail entity
 - Tasks could be either (h, ?r?, t) or (h, r, ?t?)
- Examples:



Knowledge Graphs



Examples from Chen et al EMNLP 2019

One-Shot Relational Learning for Knowledge Graphs

- (h, r, ?t?) a ranking problem, i.e. search for the right t in a candidate pool C
- Key ideas and take-home messages:
 - Embedding function:
 - Entity embeddings and neighbor encoders
 - Matching scores:
 - Matching processor to compute similarity scores
 - Could be seen as applying matching network on tail entity ranking task

Wenhan Xiong, Mo Yu, Shiyu Chang, Xiaoxiao Guo, William Yang Wang, One-Shot Relational Learning for Knowledge Graphs, EMNLP 2018

One-Shot Relational Learning for Knowledge Graphs



Tackling Long-Tailed Relations and Uncommon Entities in Knowledge Graph Completion

- (h, r, ?t) with the same setup as Xiong et al.
- Key ideas and take-home messages:
 - Take into account uncommon entities, i.e. entities that appear only several times or absent
 - Apply matching network on tail entity ranking task
 - The novelty lies on the integration of text descriptions of entities and relations to compute the triple representation

Zihao Wang, Kwun Ping Lai, Piji Li, Lidong Bing, Wai Lam, Tackling Long-Tailed Relations and Uncommon Entities in Knowledge Graph Completion, EMNLP 2019 Relational Learning for Few-Shot Link Prediction in Knowledge Graphs

- (h, r, ?t) with the same setup as Xiong et al
- Key ideas and take-home messages:
 - Relation-Meta Learner: learn the relation embeddings between head and tail entities
 - Embedding Learner: evaluate the true value of entity pairs under specific relations and evaluate the within-task training loss
 - Outperformed Xiong et al paper

Mingyang Chen, Wen Zhang, Wei Zhang, Qiang Chen, Huajun Chen, Meta Relational Learning for Few-Shot Link Prediction in Knowledge Graphs, EMNLP 2019

Relational Learning for Few-Shot Link Prediction in Knowledge Graphs



- 1) Relation-Meta Learner: based on a multilayer perceptron that takes head and tail entity as inputs and outputs a vector presenting the relation from head and tail entity
- 2) Embedding Learner: outputs the L2 norm of the vector: *head + meta relation tail*
- 3) In the within-task training step: update the meta relation vector with the meta gradient
- **4)** In the querying step: use the embedding learner to output the score as mention in 2) but with an updated meta relation vector
Summary

- Meta learning framework
- Two concrete examples:
 - Learning to initialize
 - Learning to compare
- Applications of Meta learning:
 - Speech
 - Conversational AI
 - Natural language processing

Thanks for your attention!