multi meta-models

variance reduction

conflicting gradients

conclusion

Meta-learning with Many Tasks

James T. Kwok

MLA 2023



| intro | |
|----------|---|
| •0000000 | 2 |

Outline

- introduction
- weak meta-model: multiple meta-models
- large variance in meta-learning: variance reduction
- conflicting gradients in meta-learning: soft improvement function
- conclusion

intro o●ooooooo multi meta-models

variance reduction

conflicting gradients

conclusion

Data, Data, Data







deep networks need huge amount of labeled samples

some applications have few labeled samples (e.g., aviation turbulence)



James T. Kwok

MLA 2023



Meta-learning with Many Tasks

| intro | |
|-----------|--|
| 000000000 | |

variance reduction

conflicting gradients

conclusion

Limited Data

few-shot learning



how to learn quickly with limited data?

like humans, take advantage of prior experiences

James T. Kwok MLA 2023 Meta-learning with Many Tasks

idea

- extract meta-knowledge from learned tasks
- utilize meta-knowledge to learn new tasks more quickly



- meta-learner learns generic information (meta-knowledge) across source tasks
- base learner takes meta-knowledge as prior knowledge, then generalizes for the new task using task-specific information advantages
 - data efficiency, compute efficiency, lifelong learning

conflicting gradients

conclusion 00

Model-Agnostic Meta-Learning (MAML) Algorithm

MAML

- learns a meta-initialization w (shared among all tasks) from historical tasks
- fine-tune on new task (using a few gradient updates)

bilevel optimization

$$\begin{split} \min_{\boldsymbol{w}} & \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \mathcal{L}((\boldsymbol{w}, \boldsymbol{\theta}^{i}(\boldsymbol{w})); \mathcal{D}^{i}_{\mathsf{vld}}) \\ \text{s.t.} & \boldsymbol{\theta}^{i}(\boldsymbol{w}) = \arg\min_{\boldsymbol{\theta}} \mathcal{L}((\boldsymbol{w}, \boldsymbol{\theta}); \mathcal{D}^{i}_{\mathsf{tr}}) \end{split}$$

- ullet outer level: finds a suitable meta-initialization $oldsymbol{w}$
- ullet inner level: adapts $oldsymbol{w}$ to produce task-specific $oldsymbol{ heta}^i(oldsymbol{w})$



Lots of Tasks

Example (5-way few-shot learning)

minilmageNet dataset

- meta-training set has 64 classes
- number of meta-training tasks: $\binom{64}{5} \approx 7.6 \times 10^6$

tieredImageNet dataset

- meta-training set has 351 classes
- total number of meta-training tasks: $\binom{351}{5} \approx 4.3 imes 10^{10}$

Example (recommender system)

- each user is a task
- number of Amazon Prime subscribers: 150 million as of 2020

intro 000000●00 multi meta-models

variance reduction

conflicting gradients

conclusion

Problem: Weak Meta-Model

complex environment



diverse task model parameters \rightarrow one meta-model is not sufficient

intro 0000000●0 multi meta-models

variance reduction

conflicting gradients

conclusion

Problem: Large Variance

large variance

1 data variance: limited data samples for each task

2 task variance: task sampling from task distribution

large variance \rightarrow slow convergence

multi meta-models

variance reduction

conflicting gradients

conclusion

Problem: Conflicting Gradients

conflicting gradients

• meta-learning: minimize average loss in outer level

$$\min_{\boldsymbol{w}} \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \mathcal{L}((\boldsymbol{w}, \boldsymbol{\theta}^{i}(\boldsymbol{w})); \mathcal{D}_{v|d}^{i})$$

- different task gradients point in different directions
- conflicting gradients: significantly influenced by a small subset of tasks



ullet
ightarrow poor performance

 variance reduction

conflicting gradients

conclusion

One Meta-Model is not Sufficient



MUSML (MUltiple Subspaces for Meta-Learning)



multiple subspaces

 each subspace: one type of meta-knowledge

more challenging in meta-learning

- bilevel optimization
- unseen tasks with limited samples

[Subspace learning for effective meta-learning (ICML 2022)]

| intro | |
|-------|--|
| | |

variance reduction

conflicting gradients

conclusion

MUSML



base learner

- task parameters $oldsymbol{w}_{ au}$'s for tasks au's
- lie in K subspaces $\{\mathbb{S}_1, \dots, \mathbb{S}_K\}$; with basis $\{S_1, \dots, S_K\}$
- in each subspace \mathbb{S}_k , search for a linear combination $\pmb{v}^\star_{\tau,k}$ to form \pmb{w}_{τ}

$$\mathbf{v}_{ au,k}^{\star} = \arg\min_{\mathbf{v}_{ au} \in \mathbb{R}^m} \mathcal{L}(\mathcal{D}_{ au}^{tr}; \mathbf{S}_k \, \mathbf{v}_{ au})$$

meta-learner

 \bullet learn meta-parameters $\{S_1,\ldots,S_{\mathcal{K}}\}$

multi meta-models

variance reduction

conflicting gradients

conclusion

One-Hot Subspace Selection



• not efficient: only one subspace is updated at each step



multi meta-models conflicting gradients conclusion variance reduction

MUSML Algorithm

1: for
$$t = 0, 1, ..., T - 1$$
 do
2: sample a task τ with \mathcal{D}_{τ}^{tr} and \mathcal{D}_{τ}^{vl} ;
3: base learner:
4: for $k = 1, ..., K$ do
5: initialize $\mathbf{v}_{\tau,k}^{(0)} = \mathbf{v}^{(0)}$;
6: for $t' = 0, 1, ..., T_{in} - 1$ do
7: $\mathbf{v}_{\tau,k}^{(t'+1)} = \mathbf{v}_{\tau,k}^{(t')} - \alpha \nabla_{\mathbf{v}_{\tau,k}^{(t')}} \mathcal{L}(\mathcal{D}_{\tau}^{tr}; S_{k,t}\mathbf{v}_{\tau,k}^{(t')})$;
8: end for
9: $\mathbf{v}_{\tau,k} \equiv \mathbf{v}_{\tau,k}^{(T_{in})}$;
10: $o_{\tau,k} = \mathcal{L}(\mathcal{D}_{\tau}^{tr}; S_{k,t}\mathbf{v}_{\tau,k})$;
11: end for
12: meta-learner:
13: $\mathcal{L}_{vl} = \sum_{k=1}^{K} \frac{\exp(-o_{\tau,k}/\gamma_{t})}{\sum_{k'=1}^{K} \exp(-o_{\tau,k'}/\gamma_{t})} \mathcal{L}(\mathcal{D}_{\tau}^{vl}; S_{k,t}\mathbf{v}_{\tau,k})$;
14: $\{S_{1,t+1}, ..., S_{k,t+1}\} = \{S_{1,t}, ..., S_{k,t}\} - \eta_{t} \nabla_{\{S_{1,t}, ..., S_{k,t}\}} \mathcal{L}_{vl}$;
15: end for
16: Return $S_{1,T}, ..., S_{K,T}$.

multi meta-models

variance reduction

conflicting gradients

conclusion

Bound on Expected Excess Risk

$$\mathcal{R}(\mathcal{S}) \leq \mathcal{R}^{\star} + \rho \sqrt{m} \mathbb{E}_{\tau'} \mathbb{E}_{\mathcal{D}_{\tau'}^{tr}} \| \mathbf{v}_{\tau',k_{\tau'}} - \mathbf{v}_{\tau',k_{\tau'}}^{\star} \| \\ + \rho \mathbb{E}_{\tau'} \mathbb{E}_{\mathcal{D}_{\tau'}^{tr}} \operatorname{dist}(\mathbf{w}_{\tau'}^{\star}, \mathbb{S}_{k_{\tau'}}) + K \sqrt{\frac{\nu^2 + 12\rho\nu(1 + m\alpha\delta)^{T_{in}}}{2N_{tr}}}$$

• dist
$$(\boldsymbol{w}_{\tau'}^{\star}, \mathbb{S}_{k_{\tau'}}) \equiv \min_{\boldsymbol{w} \in \mathbb{S}_{k_{\tau'}}} \| \boldsymbol{w} - \boldsymbol{w}_{\tau'}^{\star} \|$$
: distance between $\boldsymbol{w}_{\tau'}^{\star}$
and $\mathbb{S}_{k_{\tau'}}$

• minimum risk

- ② distance between approximate minimizer $v_{\tau',k_{\tau'}}$ and exact minimizer $v^{\star}_{\tau',k_{\tau'}}$
- ${f 0}$ approximation error of ${m w}^{\star}_{ au'}$ using the learned subspaces
- complexity of subspaces (m and K)

Few-shot Classification Experiments

data sets

- Meta-Dataset-BTAF: consists of 4 image classification datasets: Bird, Texture, Aircraft, Fungi
- Ø Meta-Dataset-ABF: consists of Aircraft, Bird, and Fungi
- Meta-Dataset-CIO: consists of CIFAR-FS, mini-ImageNet, Omniglot

architecture

• Conv4 backbone + Prototype classifier

meta-learning baselines

- one globally-shared meta-model: MAML, ProtoNet, ANIL, BMG
- multiple meta-models arranged in hierarchy/graph/clusters
 - Dirichlet process mixture model (DPMM)
 - hierarchically structured meta-learning (HSML)
 - automated relational meta-learning (ARML)
 - task similarity aware meta-learning (TSA-MAML, TSA-ProtoNet)

multi meta-models

variance reduction

conflicting gradients

conclusion

5-way 5-shot Classification Accuracy

| | Meta-Dataset-BTAF | Meta-Dataset-ABF | Meta-Dataset-ClO |
|--------------|-------------------|------------------|------------------|
| MAML | 57.78 | 63.86 | 74.46 |
| ProtoNet | 62.29 | 65.62 | 76.51 |
| ANIL | 58.57 | 64.43 | 74.61 |
| BMG | 60.10 | 65.80 | 77.46 |
| DPMM | 63.00 | 66.26 | 76.63 |
| TSA-MAML | 63.20 | 68.17 | 76.89 |
| HSML | 62.39 | 64.17 | 75.54 |
| ARML | 63.95 | 64.52 | 76.12 |
| TSA-ProtoNet | 63.57 | 68.77 | 77.27 |
| MUSML | 66.18 | 71.10 | 77.83 |

• MUSML is more accurate

multi meta-models

variance reduction

conflicting gradients

conclusion

Task Assignment to Learned Subspaces

Meta-Dataset-BTAF (5-way 5-shot)



• tasks from the same dataset are assigned to the same subspace

multi meta-models

variance reduction

conflicting gradients

conclusion

Cross-Domain 5-way 5-shot Classification

-

$\mathit{Meta-Dataset-BTAF} ightarrow \mathit{Meta-Dataset-CIO}$

| MAML | 64.25 |
|--------------|-------|
| ProtoNet | 66.13 |
| ANIL | 65.19 |
| BMG | 66.98 |
| DPMM | 66.73 |
| TSA-MAML | 66.85 |
| HSML | 65.18 |
| ARML | 65.37 |
| TSA-ProtoNet | 66.92 |
| MUSML | 67.41 |
| | |

multi meta-models

variance reduction

conflicting gradients

conclusion

Improving Existing Meta-learning Algorithms

- MUSML can be used with any meta-learning algorithm
- 5-way 5-shot classification accuracy

| Meta-Dataset-BTAF | Meta-Dataset-ABF | Meta-Dataset-ClO |
|-------------------|---|--|
| 58.93 | 64.19 | 75.95 |
| 65.72 | 69.15 | (1.48 |
| 60.02 66.10 | 64.51 69.23 | 76.13 77.96 |
| | Meta-Dataset-BTAF 58.93 65.72 60.02 66.10 | Meta-Dataset-BTAF Meta-Dataset-ABF 58.93 64.19 65.72 69.15 60.02 64.51 66.10 69.23 |

multi meta-models

variance reduction

conflicting gradients

conclusion

Meta-Learning for Prompt Learning in LLM

large language models



[from "A Survey of Large Language Models", Zhao et al, 20223]



pre-training

 \bullet train on a large-scale corpora \rightarrow a very capable LLM fine-tuning on downstream tasks

Large Unlabeled Corpus

• adapt the pre-trained LLM according to specific goals

fine-tune the whole model? expensive for large models (e.g., GPT-3 contains 100+ billion parameters)

Labeled Corpus

| int ro 000000000 | multi meta-models 000000000000000000000000000000000000 | variance reduction | conflicting gradients | conclusion 00 |
|---------------------|---|--------------------|-----------------------|------------------|
| Prompt | Tuning | | | |
| | | 1 11 | | |

- freeze the pre-trained model
- learn a continuous-valued prompt wrapped into the input embedding



multi meta-models

variance reduction

conflicting gradients

conclusion

Meta-Learning the Prompt

learning the prompt can be sensitive to initialization

MetaPrompting [Hou et al., 2022]

• learn a meta-initialization for all task-specific prompts



| intro | |
|-------|--|
| | |

Problems

learn only one single meta-initialized prompt \rightarrow hard to find such a prompt when the tasks are complex

need to tune the whole LM \rightarrow expensive

5-way 5-shot classification

> with or without MLM tuning



variance reduction

conflicting gradients

conclusion

Meta-Learning a Prompt Pool

use a pool of prompts to extract more task knowledge

[Effective structured-prompting by meta-learning and representative verbalizer (ICML 2023)]

- more flexible
 - \rightarrow allows better adaptation to complex tasks
 - \rightarrow no need to fine-tune the whole LM
- ullet only the prompt pool is tuned o parameter-efficient
 - $\bullet\,$ e.g., 1800 $\times\,$ fewer parameters than MetaPrompting
- also proposed a new verbalizer

| intro | |
|-------|--|
| | |

MetaPrompter

- prompt pool: has K learnable prompts
 - $\{(\boldsymbol{k}_i, \boldsymbol{\theta}_i) : i = 1, \dots, K\}$, with key \boldsymbol{k}_i and value $\boldsymbol{\theta}_i$
 - shared meta-knowledge learned by meta-learning algorithms (e.g., MAML)
- compute attention between input x and the K prompts
- instance-dependent prompt
 - weighted combinations of all the prompts in the pool via attention

$$\boldsymbol{\theta}_{\mathbf{x}}(\mathsf{K}, \boldsymbol{\Theta}) = \sum_{i=1}^{K} a_i(\mathbf{x}) \boldsymbol{\theta}_i$$

| intro | multi meta-models | variance reduction | conflicting gradients | conclusi |
|-------|---|--------------------|-----------------------|----------|
| | 000000000000000000000000000000000000000 | | 000000000000000 | |
| | | | | |

Experiments

- 6 topic classification data sets
- LM: pre-trained BERT

5-way 5-shot classification

| | #param (×10⁰) | 20News | Amazon | HuffPost | Reuters | HWU64 | Liu54 |
|---------------|---------------|--------|--------|----------|---------|-------|-------|
| HATT | 0.07 | 55.00 | 66.00 | 56.30 | 56.20 | - | - |
| DS | 1.73 | 68.30 | 81.10 | 63.50 | 96.00 | - | - |
| M LA DA | 0.73 | 77.80 | 86.00 | 64.90 | 96.70 | - | - |
| ConstrastNet | 109.52 | 71.74 | 85.17 | 65.32 | 95.33 | 92.57 | 93.72 |
| MetaPrompting | 109.52 | 85.67 | 84.19 | 72.85 | 95.89 | 93.86 | 94.01 |
| MetaPrompter | 0.06 | 88.57 | 86.36 | 74.89 | 97.63 | 95.30 | 95.47 |

MetaPrompter

- better than both prompt-based and non-prompt-based baselines
- much more parameter-efficient than MetaPrompting

| intro | |
|-------|--|
| | |

variance reduction

conflicting gradients

conclusion

Visualization

samples from each target class prefer prompts whose tokens are related to that class

 samples from cocoa tend to use the 4th and 7th prompts (whose tokens are close to words like cocoa, chocolate)



| prompt id | nearest tokens |
|-----------|--|
| 1 | copper, steel, trading, gas, fx, aluminum, earn, coffee |
| 2 | gross, ship, index, money, gold, tin, iron, retail |
| 3 | product, cpi, industrial, acquisitions, jobs, supplying, orange, sugar |
| 4 | cocoa, production, grain, livestock, wholesale, cotton, bop, crude |
| 5 | oil, national, rubber, nat, interest, price, reserves, regional |
| 6 | nat, wholesale, sugar, golden, reserves, drinks, production, product |
| 7 | chocolate, sugar, cheat, orange, trade, fx, cash, acquiring |
| 8 | aluminum, livestock, cpc, tin, shops, wheat, petrol, supply |

variance reduction

conflicting gradients

conclusion

Integration with Other Meta-Learning Algorithms

- similar performance gain when used with any meta-learning algorithm
- example: BMG
- 5-way 5-shot classification meta-testing accuracy

| | 20News | Amazon | HuffPost | Reuters | HWU64 | Liu54 |
|-------------------|--------|--------|----------|---------|-------|-------|
| MetaPrompting+BMG | 85.71 | 83.47 | 73.92 | 96.27 | 93.31 | 93.04 |
| MetaPrompter+BMG | 87.91 | 86.45 | 74.99 | 98.01 | 95.41 | 94.52 |

Variance Reduction

two sources of gradient variance

- data sampling for each task
- Sampling of tasks

variance reduction: reduce variance in the stochastic gradients

single-level optimization

- classic methods: SVRG, SAG, SDCA
- faster convergence than SGD both theoretically and empirically

bilevel optimization

- SUSTAIN, MRBO/VRBO, RSVRB, VR-BiAdam
- faster than bilevel algorithms without variance reduction

Removing Batch Gradient

batch gradient has to be computed occasionally \rightarrow expensive

STORM [Cutkosky & Orabona (2019)]

- does NOT need batch gradient
- compute stochastic gradients at two points \pmb{w}_t, \pmb{w}_{t-1} with the same mini-batch ξ_t

```
1: Input: w_0, step-size \{\eta_t\}, decay parameter \{\gamma_t\}.

2: c_0 = \nabla \ell(w_0; \xi_0)

3: w_1 = w_0 - \eta_0 c_0

4: for t = 1 to T - 1 do

5: sample \xi_t

6: c_t = \nabla \ell(w_t; \xi_t) + (1 - \gamma_t)(c_{t-1} - \nabla \ell(w_{t-1}; \xi_t))

7: w_{t+1} = w_t - \eta_t c_t

8: end for
```

- cf. SGD: $\boldsymbol{c}_t = \nabla \ell(\boldsymbol{w}_t; \xi_t)$
- cf. momentum: $\boldsymbol{c}_t = \gamma_t \nabla \ell(\boldsymbol{w}_t; \xi_t) + (1 \gamma_t) \boldsymbol{c}_{t-1}$

multi meta-models

variance reduction

conflicting gradients

conclusion

Variance Reduction for Meta-Learning

straightforward approach

- use existing variance reduction methods
- replace the original gradients by their variance-reduced counterparts

meta-learning as bilevel optimization

- requires storing all task-specific parameters → prohibitive storage for large number of tasks or huge task models
- requires a large number of inner steps for convergence
 - meta-learning often uses only a small number of steps to avoid overfitting the limited data

multi meta-models

variance reduction

conflicting gradients

conclusion

Variance Reduction in Single-Level Meta-Learning

1: for
$$t = 0$$
 to $T - 1$ do
2: sample tasks $\mathcal{I}_t \subset \mathcal{I}$
3: for $i \in \mathcal{I}_t$ do
4: $u_0^i = w_t$
5: for $k = 0$ to $K - 1$ do
6: obtain task i samples $\xi_{k,t}^i$
7: $u_{k+1}^i = u_k^i - \alpha \nabla \ell(u_k^i, \xi_{k,t}^i)$
8: end for
9: $c_t^i = \frac{1}{K\alpha}(w_t - u_K^i)$
10: end for
11: $c_t = \frac{1}{|\mathcal{I}_t|} \sum_{i \in \mathcal{I}_t} c_t^i$
12: $w_{t+1} = w_t - \eta_t c_t$
13: end for

Reptile

- uses average gradient for adaptation
- VFML [Wang et al. (2021)]
 - integrates STORM with Reptile

problems

- lacks theoretical properties
- inferior empirical performance

multi meta-models

variance reduction

conflicting gradients

conclusion

Proposed Family of Methods

[Efficient variance reduction for meta-learning (ICML 2022)]

- can be integrated with various meta-learning algorithms
 - Reptile → VR-Reptile (Variance-Reduced Reptile)
 - MAML \rightarrow VR-MAML
 - FOMAML \rightarrow VR-FOMAML
 - $BMG \rightarrow VR\text{-}BMG$

multi meta-models

variance reduction

conflicting gradients

conclusion

Reptile and VR-Reptile

1: for
$$t = 0$$
 to $T - 1$ do
2: sample tasks $\mathcal{I}_t \subset \mathcal{I}$
3: for $i \in \mathcal{I}_t$ do
4: $u_0^i = w_t$
5: for $k = 0$ to $K - 1$ do
6: obtain task i samples $\xi_{k,t}^i$
7: $u_{k+1}^i = u_k^i - \alpha \nabla \ell(u_k^i, \xi_{k,t}^i)$
8: end for
9: $c_t^i = \frac{1}{K\alpha}(w_t - u_K^i)$
10: end for
11: $c_t = \frac{1}{|\mathcal{I}_t|} \sum_{i \in \mathcal{I}_t} c_t^i$
12: $w_{t+1} = w_t - \eta_t c_t$
13: end for

| 1: | initialization |
|-----|---|
| 2: | for $t=1$ to $T-1$ do |
| 3: | sample tasks ${\mathcal I}_t \subset {\mathcal I}$ |
| 4: | for $i \in \mathcal{I}_t$ do |
| 5: | $u_0^i = w_t$ |
| 6: | $\mathbf{v_0^i} = \mathbf{w_{t-1}}$ |
| 7: | for $k=0$ to $K-1$ do |
| 8: | obtain task i samples $\xi^i_{k,t}$ |
| 9: | $oldsymbol{u}_{k+1}^i = oldsymbol{u}_k^i - lpha abla \ell(oldsymbol{u}_k^i; oldsymbol{\xi}_{k,t}^i)$ |
| 10: | $\mathbf{v}_{k+1}^i = \mathbf{v}_k^i - lpha abla \ell(\mathbf{v}_k^i; \mathbf{\xi}_{k,t}^i)$ |
| 11: | end for |
| 12: | $	ilde{oldsymbol{d}}_{t-1}^i = rac{1}{Klpha} (oldsymbol{w}_{t-1} - oldsymbol{v}_K^i)$ |
| 13: | $	ilde{m{c}}_t^i = rac{1}{Klpha} (m{w}_t - m{u}_K^i)$ |
| 14: | end for |
| 15: | $\widetilde{\pmb{d}}_{t-1} = rac{1}{ \mathcal{T}_t } \sum_{i \in \mathcal{T}_t} \widetilde{\pmb{d}}_{t-1}^i$ |
| 16: | $\tilde{c}_t =$ |
| | $rac{1}{ \mathcal{I}_t } \sum_{i \in \mathcal{I}_t} \tilde{\boldsymbol{c}}_t^i + (1 - \gamma_t) (\tilde{\boldsymbol{c}}_{t-1} - \tilde{\boldsymbol{d}}_{t-1})$ |
| 17: | $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{\tilde{c}}_t$ |
| 18: | end for |

| intro | |
|----------|--|
| 00000000 | |

variance reduction

conflicting gradients

conclusion

Remarks

1: initialization 2: for t = 1 to T - 1 do 3: sample tasks $\mathcal{I}_t \subset \mathcal{I}$ 4: for $i \in \mathcal{I}_t$ do 5: $\mathbf{u}_0' = \mathbf{w}_t$ 6: $\mathbf{v}_0' = \mathbf{w}_{t-1}$ 7: for k = 0 to K - 1 do 8: obtain task *i* samples $\xi_{k,t}^{i}$ 9: $\mathbf{u}_{k+1}^{i} = \mathbf{u}_{k}^{i} - \alpha \nabla \ell(\mathbf{u}_{k}^{i}; \boldsymbol{\xi}_{k}^{i})$ $\mathbf{v}_{k+1}^{i} = \mathbf{v}_{k}^{i} - \alpha \nabla \ell(\mathbf{v}_{k}^{i}; \xi_{k}^{i})$ 10: 11:end for $\tilde{d}_{t-1}^{i} = \frac{1}{\kappa_{r}} (w_{t-1} - v_{K}^{i})$ 12: 13: $\tilde{\boldsymbol{c}}_t^i = \frac{1}{\kappa_s} (\boldsymbol{w}_t - \boldsymbol{u}_K^i)$ 14. end for $\tilde{\boldsymbol{d}}_{t-1} = \frac{1}{|\mathcal{T}_t|} \sum_{i \in \mathcal{T}_t} \tilde{\boldsymbol{d}}_{t-1}^i$ 15: $\tilde{\boldsymbol{c}}_t = rac{1}{|\mathcal{I}_t|} \sum_{i \in \mathcal{I}_t} \tilde{\boldsymbol{c}}_t^i + (1 - 1)$ 16: $(\gamma_t)(\tilde{\boldsymbol{c}}_{t-1} - \tilde{\boldsymbol{d}}_{t-1})$ 17: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \tilde{\mathbf{c}}_t$ 18: end for

- all γ_t 's = 1 \rightarrow VR-Reptile reduces to Reptile
- $K = 1 \rightarrow VR$ -Reptile reduces to STORM
- space-efficient: does NOT need to store task-specific $\boldsymbol{u}_{K}^{i}/\boldsymbol{v}_{K}^{i}$ and $\tilde{\boldsymbol{c}}_{t}^{i}/\tilde{\boldsymbol{d}}_{t-1}^{i}$
 - cf. direct application of bilevel variance reduction methods → requires storing all task-specific parameters

| intro | |
|----------|---|
| 00000000 | c |

variance reduction

conflicting gradients

conclusion 00

Convergence

Reptile

$$\frac{1}{T}\mathbb{E}\left[\sum_{t=0}^{T-1} \|\nabla \tilde{\mathcal{L}}(\boldsymbol{w}_t)\|^2\right] \leq \frac{\sqrt{2}G_1/\eta_0}{\sqrt{T}} + \frac{\sqrt{2}G_1/\eta_0}{T}$$

VR-Reptile

$$\frac{1}{T}\mathbb{E}\left[\sum_{t=0}^{T-1} \|\nabla \tilde{\mathcal{L}}(\boldsymbol{w}_t)\|^2\right] \leq \frac{4\tilde{M}G_2}{T^{2/3}} + \frac{65}{7} \cdot \frac{\tilde{M}G_2}{T}$$

• faster convergence rate

Experiments: Convergence

- 5-shot 5-way on Meta-Dataset (Bird, Texture, Aircraft, Fungi)
- accuracy with number of outer-loop iterations



James T. Kwok MLA 2023

Meta-learning with Many Tasks

conclusion

multi meta-models

variance reduction

conflicting gradients

conclusion

Reduce Stochastic Variance

variance of weight update \tilde{c}_t relative to squared norm

$$\mathbb{E} \| ilde{m{c}}_t - \mathbb{E}[ilde{m{c}}_t] \|^2 / \| \mathbb{E}[ilde{m{c}}_t] \|^2$$



- variance \downarrow
- \bullet Reptile has larger variance \rightarrow variance reduction also has larger improvements on Reptile

Few-Shot Classification on mini-ImageNet

architecture

• Conv4 backbone + Prototype classifier

meta-learning baselines

- standard methods: MAML/FOMAML/Reptile/BMG/DRS/ProtoNet
- straightforward STORM variants: replace stochastic gradients by variance-reduced counterparts by STORM
- VFML (Reptile+STORM)
- ANIL and variants with bilevel optimization variance reduction methods (SUSTAIN/MRBO/VRBO)

| intro | |
|----------|---|
| 00000000 | ċ |

variance reduction

conflicting gradients

conclusion

Accuracy on mini-ImageNet

| | var reduction | single/bilevel | 1-shot 5-way | 5-shot 5-way |
|---|--|--|---|---|
| MAML FOMAML | × × | single single | 48.7±1.8 48.1±1.8 | 63 1±0 9 63 2±0 9 |
| Reptile BMG | ×× | single single | 50.0 ± 0.3 50.7 ± 0.5 | 66.0±0.6 65.6±0.6 |
| DRS ProtoNet | ××× | single single | 24.5 ± 0.8 49.4 \pm 0.8 | 30.4 ± 0.6 68.2 ± 0.7 |
| MAML+STORM FOMAML+STORM Reptile+STORM BMG+STORM DRS+STORM | $\begin{array}{c} \checkmark\\ \checkmark\\ \checkmark\\ \checkmark\\ \checkmark\\ \checkmark\\ \checkmark\end{array}$ | sing e sing e sing e sing e sing e | $\begin{array}{c} 47.9 \pm 1.4 \\ 48.0 \pm 1.6 \\ 49.9 \pm 0.3 \\ 46.7 \pm 0.6 \\ 24.7 \pm 1.1 \end{array}$ | $\begin{array}{c} 61.6 \pm 1.2 \\ 63.4 \pm 1.1 \\ 66.2 \pm 0.3 \\ 60.9 \pm 0.8 \\ 30.3 \pm 0.7 \end{array}$ |
| VR-MAML VR-FOMAML VR-Reptile VR-BMG | \checkmark \checkmark \checkmark | single single single single | $49.2 \pm 1.4 \\ 48.3 \pm 1.2 \\ 50.4 \pm 0.4 \\ 51.4 \pm 0.3$ | 63.6±0.8 63.4±0.6 67.6±0.8 68.4±0.6 |
| VFML | \checkmark | single | 49.6±0.5 | 66.2±0.8 |
| ANIL ANIL+SUSTAIN ANIL+MRBO ANIL+VRBO | \sim \sim \sim \sim \sim \sim | bil ev el bil ev el bil ev el bil ev el | $\begin{array}{r} 46.9 \pm 0.4 \\ 47.0 \pm 0.4 \\ 47.2 \pm 0.5 \\ 47.2 \pm 0.4 \end{array}$ | |

• integrating with any of the variance reduction methods leads to better performance

multi meta-models

variance reduction

conflicting gradients

conclusion

5-shot 5-way Classification Accuracy on Meta-Dataset

| | var reduction | single/bilevel | Bird | Texture | Aircraft | Fungi |
|---------------|---------------|----------------|-------|---------|----------|-------|
| MAML | × | single | 74.56 | 45.68 | 69.06 | 53.68 |
| FOMAML | × | single | 73.64 | 42.82 | 66.38 | 52.18 |
| Reptile | × | single | 74.60 | 43.26 | 66.46 | 52.88 |
| BMG | × | single | 74.52 | 43.74 | 66.64 | 53.02 |
| DRS | × | single | 53.34 | 33.28 | 41.06 | 37.64 |
| ProtoNet | × | single | 74.22 | 49.86 | 71.38 | 53.94 |
| MAML+STORM | \checkmark | single | 74.86 | 45.26 | 68.48 | 53.72 |
| FOMAML+STORM | | single | 73.72 | 42.78 | 66.22 | 52.28 |
| Reptile+STORM | | single | 75.24 | 44.60 | 67.48 | 52.54 |
| BMG+STORM | \checkmark | single | 73.16 | 42.32 | 66.46 | 51.74 |
| DRS+STORM | \checkmark | single | 53.48 | 33.42 | 41.12 | 37.72 |
| VR-MAML | \checkmark | single | 75.06 | 46.18 | 68.36 | 53.86 |
| VR-FOMAML | | single | 74.28 | 43.28 | 66.98 | 52.16 |
| VR-Reptile | \checkmark | single | 76.48 | 46.94 | 71.62 | 54.24 |
| VR-BMG | \checkmark | single | 76.56 | 47.28 | 71.48 | 54.38 |
| VFML | \checkmark | single | 74.38 | 44.48 | 65.64 | 52.76 |
| ANIL | × | bilevel | 73.68 | 41.96 | 68.74 | 52.84 |
| ANIL+SUSTAIN | \checkmark | bilevel | 73.74 | 42.12 | 68.82 | 52.78 |
| ANIL+MRBO | | bilevel | 73.78 | 42.18 | 68.78 | 52.86 |
| ANIL+VRBO | | bilevel | 73.88 | 42.22 | 68.74 | 52.82 |

multi meta-models

variance reduction

conflicting gradients

conclusion 00

Conflicting Gradients

(meta-learning)
$$\min_w \sum_{ au=1}^m \mathcal{L}_ au(w_ au)$$
 s.t. $w_ au = w_ au^*(w)$

conflicting gradients

significantly influenced by a small subset of tasks



Meta-Learning as Multi-Objective Optimization (MOO)

$$MOO) \qquad \min_{w} (\mathcal{L}_1(w_1^*(w)), \dots, \mathcal{L}_m(w_m^*(w)))$$

each task is an objective

gradient-based MOO solvers

- e.g., MGDA, PCGard, CAGard
- in each iteration, find a descent direction common to all objectives

Example (MGDA (multiple-gradient descent algorithm))

direction: $\sum_{\tau=1}^{m} \gamma_{\tau}^* \nabla_x f_{\tau}(x)$

$$\{\gamma_{\tau}^*\} = \arg\min_{\{\gamma_{\tau}\}} \left\| \sum_{\tau=1}^m \gamma_{\tau} \nabla_x f_{\tau}(x) \right\|^2 \text{ s.t. } \sum_{\tau=1}^m \gamma_{\tau} = 1, \gamma_{\tau} \ge 0$$

| intro | multi meta-models | variance reduction | conflicting gradients | conclusion |
|-----------|-------------------------|--------------------|-----------------------|------------|
| 000000000 | 00000000000000000000000 | | 00●00000000000000 | 00 |
| Problem | | | | |

$$\{\gamma_{\tau}^*\} = \arg\min_{\{\gamma_{\tau}\}} \left\| \sum_{\tau=1}^m \gamma_{\tau} \nabla_x f_{\tau}(x) \right\|^2 \text{ s.t. } \sum_{\tau=1}^m \gamma_{\tau} = 1, \gamma_{\tau} \ge 0$$

requires gradients from all objectives

meta-learning: computing all task gradients can be very expensive

Example (5-way few-shot classification on Minilmagenet)

total number of meta-training tasks: $\binom{64}{5} \approx 7 \times 10^6$

multi meta-models

variance reduction

conflicting gradients

conclusion

Straightforward Solution

use a mini-batch B of objectives in each iteration

$$\min_{\{\gamma_{\tau}\}} \left\| \sum_{\tau \in \mathcal{B}} \gamma_{\tau} \nabla_{\tau} (\mathcal{L}_{\tau}(w_{\tau}^{*}(w))) \right\|^{2} \text{ s.t. } \sum_{\tau \in \mathcal{B}} \gamma_{\tau} = 1, \gamma_{\tau} \ge 0$$

problem: does not converge to Pareto front

- two tasks/objectives (mini-batch size of 1)
- (batch) MGDA vs mini-batch MGDA



multi meta-models

variance reduction

 conclusion

Improvement Function in MOO Problem

[Enhancing meta-learning via multi-objective soft improvement functions (ICLR 2023)]

$$\min_{x}[f_1(x),\ldots,f_m(x)]$$

improvement function

$$H(x,x') = \max_{\tau=1\dots,m} \left\{ f_{\tau}(x) - f_{\tau}(x') \right\}$$

 $x^* = \arg \min_x H(x, x^*)$ for a Pareto stationary point x^*

• to find x^* : descent on H

$$x^{s+1} = x^s + \beta d^s$$

$$d^* = \arg\min_d H(x^s + d, x^s) + \frac{\lambda'}{2} \|d\|^2$$

• $s \to \infty$, x^s is Pareto stationary

intro multi meta-models variance reduction con

conclusion

Meta-Learning Context

$$\min_{w}(\mathcal{L}_1(w_1^*(w)),\ldots,\mathcal{L}_m(w_m^*(w)))$$

improvement function

$$H(w, w') = \max_{\tau=1...,m} \left\{ \mathcal{L}_{\tau}(w_{\tau}^*(w)) - \mathcal{L}_{\tau}(w_{\tau}^*(w')) \right\}$$

=
$$\max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\mathcal{L}_{\tau}(w_{\tau}^*(w)) - \mathcal{L}_{\tau}(w_{\tau}^*(w')) \right]$$

• π : probability density function on au update

$$\begin{aligned} w^{s+1} &= w^s + \beta d^* \\ d^* &= \arg\min_d \left(\max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\mathcal{L}_{\tau} (w^*_{\tau}(w^s + d)) - \mathcal{L}_{\tau} (w^*_{\tau}(w^s)) \right] \right) + \frac{\lambda'}{2} \|d\|^2 \end{aligned}$$

multi meta-models

variance reduction

conflicting gradients

conclusion

Soft Improvement Function

take first-order approximation

$$\min_{d} \left(\max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\mathcal{L}_{\tau} (w_{\tau}^{*}(w^{s} + d)) - \mathcal{L}_{\tau} (w_{\tau}^{*}(w^{s})) \right] + \frac{\lambda'}{2} \|d\|^{2} \right)$$

$$= \max_{\pi} \underbrace{\left(\min_{d} \mathbb{E}_{\tau \sim \pi} \left[\mathcal{L}_{\tau} (w_{\tau}^{*}(w^{s} + d)) - \mathcal{L}_{\tau} (w_{\tau}^{*}(w^{s})) \right] + \frac{\lambda'}{2} \|d\|^{2} \right)}_{\text{closed-form solution:}} \frac{\frac{-1}{2\lambda'} \|\mathbb{E}_{\tau \sim \pi} \nabla_{w} \mathcal{L}_{\tau} (w_{\tau}^{*}(w))\|_{w = w^{s}} \|^{2}}$$

expectation $\mathbb{E}_{\tau\sim\pi}$

- sample tasks from uniform distribution U, then weight sampled task au with $r(au) \equiv \pi(au)/U(au)$
- ullet parameterize r as a neural network $r_{ heta}$ with parameter heta

$$\max_{\theta} \mathcal{K}(\theta) \equiv \frac{-1}{2\lambda'} \left[\mathbb{E}_{\tau \sim U} r_{\theta}(\tau) \nabla_{w} \mathcal{L}_{\tau}(w_{\tau}^{*}(w)) |_{w = w^{s}} \right]^{2} - \frac{\lambda''}{2} (\mathbb{E}_{\tau \sim U} r_{\theta}(\tau) - 1)^{2}$$

• the last term: enforcing the constraint $\mathbb{E}_{ au \sim U} r_{ heta}(au) = 1$

| intro | multi meta-models | variance reduction | conflicting gradients | conclusio |
|-------|-------------------|--------------------|-----------------------|-----------|
| | | | 000000000000000 | |
| | | | | |

$$\max_{\theta} \mathcal{K}(\theta) \equiv \frac{-1}{2\lambda'} \left[\mathbb{E}_{\tau \sim U} r_{\theta}(\tau) \nabla_{w} \mathcal{L}_{\tau}(w_{\tau}^{*}(w)) |_{w=w^{s}} \right]^{2} - \frac{\lambda''}{2} \left(\mathbb{E}_{\tau \sim U} r_{\theta}(\tau) - 1 \right)^{2}$$

 $\mathbb{E}_{ au \sim oldsymbol{U}}$: still needs to access all tasks

| int ro 000000000 | multi meta-models | variance reduction | conflicting gradients ooooooooooooooo | conclusion 00 |
|----------------------------|-------------------|--------------------|--|------------------|
| Mini-Bato | ch Again! | | | |

• B: mini-batch of k tasks

$$\underbrace{\frac{-1}{2\hat{\lambda}'} \left\| \frac{1}{|B|} \sum_{\tau \in B} r_{\theta}(\tau) \nabla_{w} \mathcal{L}_{\tau}(w_{\tau}^{*}(w)) \right\|_{w=w^{s}} \right\|^{2} - \frac{\hat{\lambda}''}{2} \left(\frac{1}{|B|} \sum_{\tau \in B} r_{\theta}(\tau) - 1 \right)^{2}}_{\tilde{K}_{B}(\theta)}}$$

 $\mathcal{K}(heta)$ can be approximated by using mini-batches

when
$$k \ll m$$
, $\left(\mathcal{K}(\theta) - rac{1}{|\mathcal{B}|} \sum_{B \in \mathcal{B}} \tilde{\mathcal{K}}_B(\theta)\right)^2 \leq rac{\mathcal{G}_1}{k\lambda'} + rac{\mathcal{G}_2}{k}\lambda''$

• in experiments: $m\geq 10^6$, $k\approx 10^2
ightarrow {
m diff} \le 0.2$

intro multi meta-models variance reduction conflicting gradients conclusion

Soft Improvement Multi-Objective Meta Learning (SIMOL)

update the distribution

$$\theta^{s+1} = \theta^s + \beta' \nabla_\theta \tilde{K}_B(\theta^s)$$

update the meta-model

$$egin{array}{rcl} w^{s+1}&=&w^s+eta ilde{d}^*\ \widetilde{d}^*&=&-rac{1}{\lambda'|B|}\sum_{ au\in B}
abla_w\mathcal{L}_ au(w^*_ au(w))|_{w=w^s} \end{array}$$

theoretically converge to an ϵ -Pareto stationary point

| intro mu | ultimeta-models \ | variance reduction | conflicting gradients |
|----------|-------------------|--------------------|---|
| | | | 000000000000000000000000000000000000000 |

Experiment

- two objectives (mini-batch size of 1)
- MGDA vs mini-match MGDA vs SIMOL



| intro 000000000 | multi meta-models 00000000000000000000000 | variance reduction | conflicting gradients 000000000000000000000000000000000000 | conclusion |
|--------------------|--|--------------------|---|------------|
| Few/shot | Regression | | | |

- synthetic data: $y = a_{ au} \sin(x + b_{ au})$
- 160,000 meta-training tasks, 1,000 meta-testing tasks
- meta-learner / re-weighting network: MLP with 2 FC layers

| | overal | | w orst - 10% | |
|---|---|---|--|--|
| MAML | 5- sh ot | 2-shot | 5-shot | 2-sh ot |
| min average loss mini-batch MGDA mini-batch CAGrad SIMOL | $\begin{array}{c} 0.43 \pm 0.11 \\ 0.60 \pm 0.02 \\ 1.90 \pm 0.24 \\ 0.34 \pm 0.04 \end{array}$ | $\begin{array}{c} 1.70 \pm 0.11 \\ 1.73 \pm 0.10 \\ 1.82 \pm 0.44 \\ 1.24 \pm 0.08 \end{array}$ | $ \begin{vmatrix} 2.13 \pm 0.21 \\ 2.62 \pm 0.10 \\ 8.18 \pm 0.63 \\ 1.69 \pm 0.18 \end{vmatrix} $ | $\begin{array}{c} 7.75 \pm 0.48 \\ 7.04 \pm 0.51 \\ 8.23 \pm 2.33 \\ \textbf{5.66} \pm 0.33 \end{array}$ |

- SIMOL consistently outperforms SGD, MGDA, and CAGrad in terms of both the overall and worst-10% MSEs
- mini-batch MGDA and CAGrad are not good

variance reduction

conflicting gradients

conclusion

Few-Shot Image Classification (minilmageNet)

accuracy

| | | overall | | worst- 10% | |
|--------------------|---|---|---|---|--|
| | | 1-sh ot | 5-sh ot | 1-shot | 5-sh ot |
| (MAML) | min average loss mini-batch MGDA mini-batch CAGrad SIMOL | $\begin{array}{ } 49.24 \pm 0.78 \\ 46.08 \pm 0.78 \\ 44.67 \pm 0.75 \\ 50.62 \pm 1.39 \end{array}$ | $62.13 \pm 0.72 \\ 60.15 \pm 0.41 \\ 60.05 \pm 0.67 \\ 65.83 \pm 0.86$ | $\begin{array}{c} 13.33 \pm 1.07 \\ 10.60 \pm 1.33 \\ 11.33 \pm 1.12 \\ 14.99 \pm 1.72 \end{array}$ | $\begin{array}{c} 41.71 \pm 1.02 \\ 39.67 \pm 0.55 \\ 40.01 \pm 0.88 \\ \textbf{44.81} \pm 0.58 \end{array}$ |
| (MAML variants) | Reptile FOMAML IMAML Meta-MP TS-MAML MTL | $\begin{array}{c} 47.07 \pm 0.26 \\ 45.53 \pm 1.58 \\ 49.30 \pm 1.88 \\ 48.51 \pm 0.92 \\ 48.44 \pm 0.91 \\ 49.87 \pm 0.41 \end{array}$ | $\begin{array}{c} 62.74 \pm 0.37 \\ 61.02 \pm 1.12 \\ 59.77 \pm 0.41 \\ 64.15 \pm 0.92 \\ 65.52 \pm 0.68 \\ 65.81 \pm 0.33 \end{array}$ | - - 13.64 ± 1.45 | 43.42 ± 0.47 |

• SIMOL consistently outperforms all the baselines in terms of both overall and worst-10% accuracies

multi meta-models

variance reduction

conflicting gradients

conclusion

Few-Shot Image Classification (tieredImageNet)

accuracy

| | | overall | | worst- 10% | |
|--------------------|---|--|---|---|---|
| | | 1-sh ot | 5-sh ot | 1-shot | 5-sh ot |
| (MAML) | min average loss mini-batch MGDA mini-batch CAGrad SIMOL | $ \begin{vmatrix} 50.58 \pm 1.44 \\ 22.92 \pm 1.04 \\ 49.04 \pm 0.93 \\ 51.42 \pm 1.50 \end{vmatrix} $ | $\begin{array}{c} 69.33 \pm 0.74 \\ 53.41 \pm 0.74 \\ 65.43 \pm 0.73 \\ \textbf{70.13} \pm 0.74 \end{array}$ | $ \begin{vmatrix} 11.60 \pm 1.96 \\ 7.12 \pm 2.11 \\ 11.40 \pm 1.97 \\ 12.00 \pm 1.95 \end{vmatrix} $ | $\begin{array}{c} 46.95 \pm 1.14 \\ 32.79 \pm 0.88 \\ 42.63 \pm 0.95 \\ 47.51 \pm 1.46 \end{array}$ |
| (MAML variants) | Reptile FOMAML IMAML Meta-MP TS-MAML MTL | $ \begin{vmatrix} 49.12 \pm 0.43 \\ 45.53 \pm 1.58 \\ 38.54 \pm 1.37 \\ 50.14 \pm 1.37 \\ 48.82 \pm 0.88 \\ 51.02 \pm 0.46 \end{vmatrix} $ | $\begin{array}{c} 65.99 \pm 0.42 \\ 61.02 \pm 1.12 \\ 60.24 \pm 0.76 \\ 68.30 \pm 0.91 \\ 67.82 \pm 0.72 \\ 66.47 \pm 0.39 \end{array}$ | - - - 13.60 ± 1.59 | - - - 49.45 ± 0.58 |

multi meta-models

variance reduction

conflicting gradients

conclusion

Using Batch MGDA and CAGrad

per-epoch running time (minilmageNet)

| standard MAML | SIMOL | batch MGDA | batch CAGrad |
|---------------|----------------------|------------|--------------|
| 2.0 sec | 2.3 <mark>sec</mark> | 6.9 days | 5.6 days |

- SIMOL has comparable per-epoch running time as MAML
- batch MGDA and CAgrad are much more computationally expensive (around 432,000 times slower)

Conclusion

how to learn with a lot of tasks in meta-learning?

use multiple meta-models into meta-learning

- subspace learning of task model parameters
- meta-learning a prompt pool

variance reduction to accelerate convergence of meta-learning

• does not need to store all task-specific parameters

soft improvement function in MOO to handle conflicting gradients

 scalable gradient-based solver with theoretical guarantees to Pareto-optimality

empirically, outperform existing meta-learning algorithms

| intro | multi meta-models | variance reduction | conflicting gradients | conclusion |
|----------|---|--------------------|-----------------------|------------|
| 00000000 | 000000000000000000000000000000000000000 | 0000000000000 | 000000000000000 | 00 |

References

- W. Jiang, J. Kwok, Y. Zhang. Subspace learning for effective meta-learning. ICML, 2022
- W. Jiang, Y. Zhang, J. Kwok. Effective structured-prompting by meta-learning and representative verbalizer. ICML, 2023.
- H. Yang, J. Kwok. Efficient variance reduction for meta-learning. ICML, 2022
- R. Yu, W. Chen, X. Wang, J. Kwok. Enhancing meta-learning via multi-objective soft improvement functions, ICLR, 2023

