

Optimization-Based Meta-Learning

CS 330

Course Reminders

Project group form due **tonight.**
(for assigning project mentors)

Homework 1 due **Wednesday**

Following up on some high-res feedback:

- optional reading materials posted on website
- coding environment set-up: detailed guidance on Azure for HW2 & HW3
- managing questions in lecture

Plan for Today

Recap

- Meta-learning problem & black-box meta-learning

Optimization Meta-Learning

} Part of Homework 2!

- Overall approach
- Compare: optimization-based vs. black-box
- Challenges & solutions
- Case study of land cover classification (time-permitting)

Goals for by the end of lecture:

- Basics of optimization-based meta-learning techniques (& how to implement)
- Trade-offs between black-box and optimization-based meta-learning

Problem Settings Recap

Multi-Task Learning

Solve multiple tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ at once.

$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$

Transfer Learning

Solve target task \mathcal{T}_b after solving source task \mathcal{T}_a
by *transferring* knowledge learned from \mathcal{T}_a

Meta-Learning Problem

Transfer Learning with Many Source Tasks

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, solve new task $\mathcal{T}_{\text{test}}$ more quickly / proficiently / stably

Example Meta-Learning Problem

5-way, 1-shot image classification (Minilmagenet)



Given 1 example of 5 classes:

meta-test



Classify new examples



Can replace image classification with: regression, language generation, skill learning,

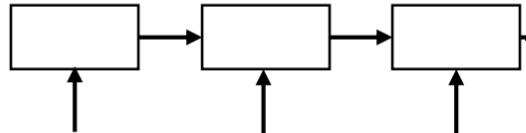
any ML problem

Black-Box Adaptation

Something like
an RNN or LSTM

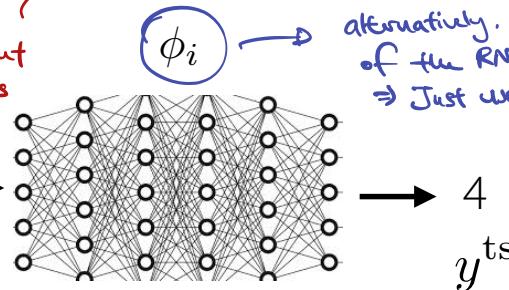
f_θ

Have that NN output
a set of params
for that task



$\mathcal{D}_i^{\text{tr}}$

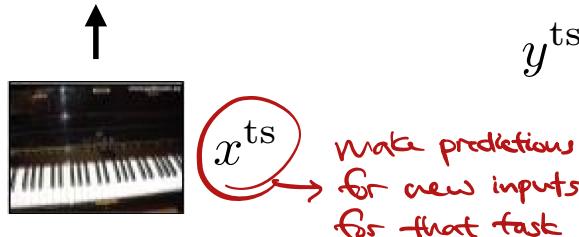
so it is a hypernetwork!



alternatively, it could just give us a hidden state
of the RNN (instead of giving us an entire NN params)
⇒ Just use the RNN from before.

general form:

$$y^{\text{ts}} = f_{\text{black-box}}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}})$$



+ expressive

↳ i.e. can represent lots of diff. learning algorithms

How else can we represent $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$?

- challenging optimization problem

going from training datasets
to a set of fast parameters
↳ NN parameters optimized
to execute a certain task.

What if we treat it as an **optimization** procedure?

Plan for Today

Recap

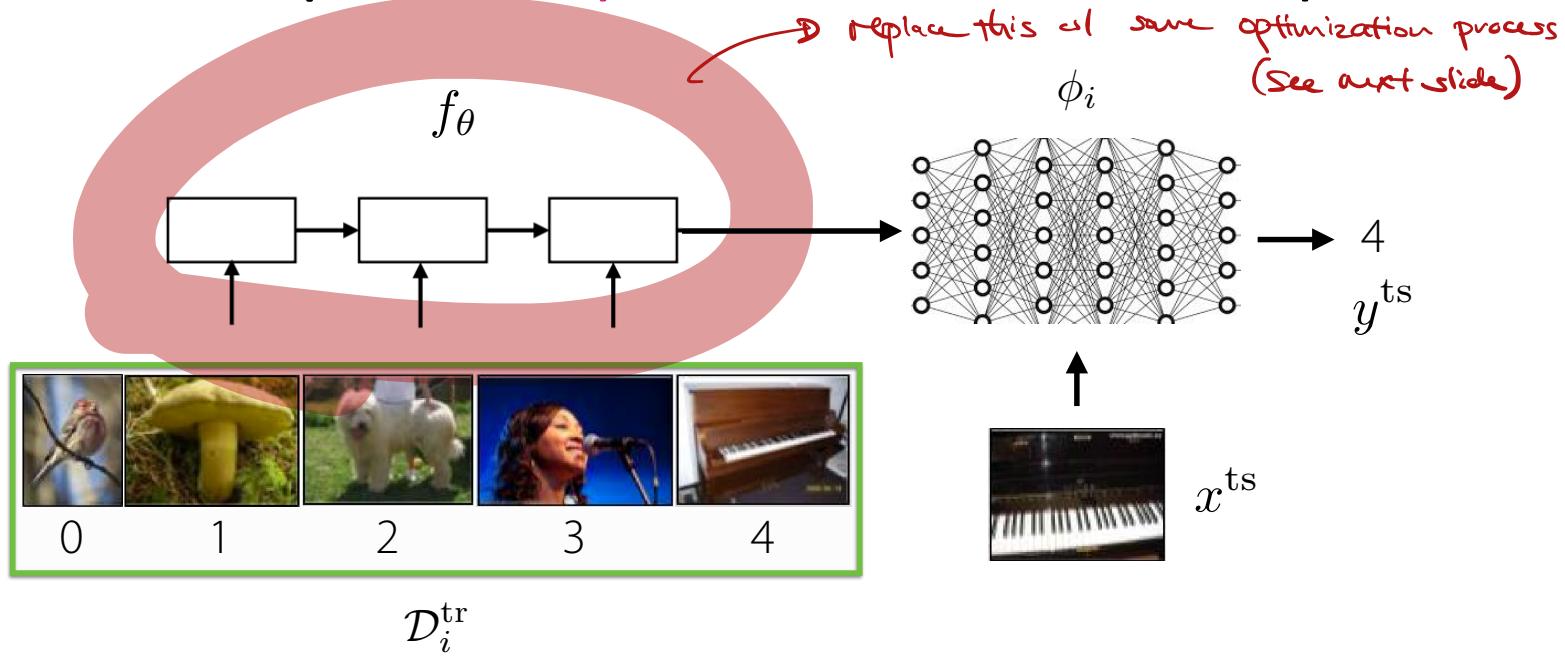
- Meta-learning problem & black-box meta-learning

Optimization Meta-Learning

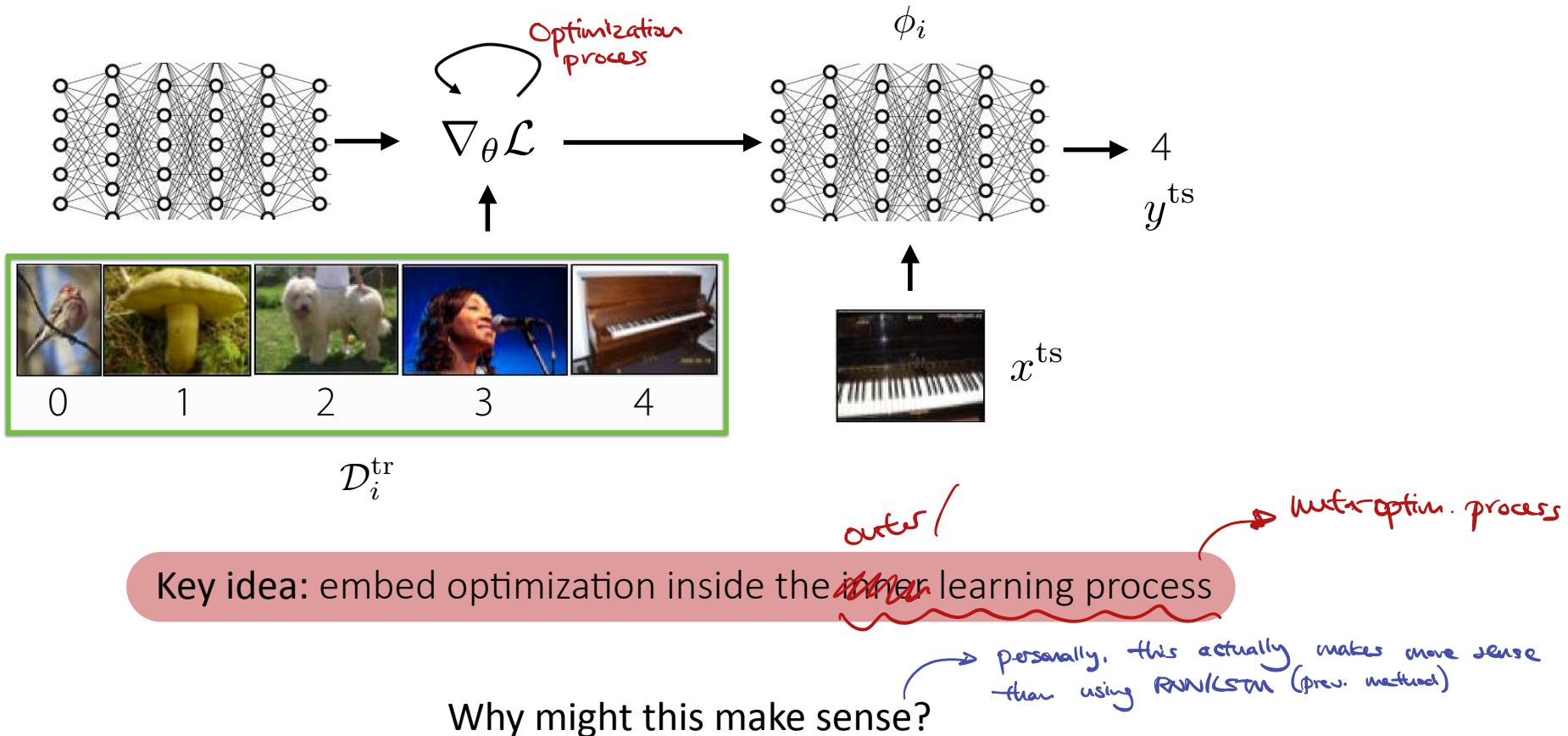
} Part of Homework 2!

- Overall approach
- Compare: optimization-based vs. black-box
- Challenges & solutions
- Case study of land cover classification (time-permitting)

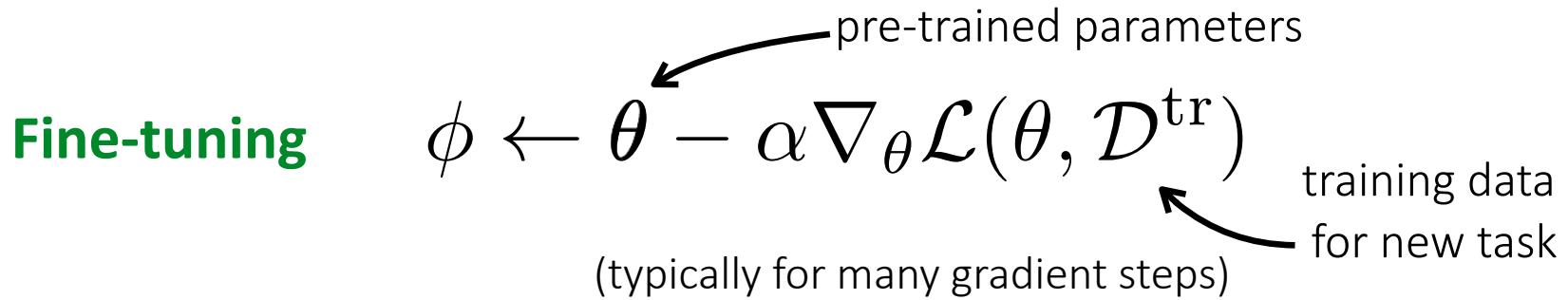
~~Black-Box~~ Adaptation Optimization-Based Adaptation



~~Black-Box~~ Adaptation Optimization-Based Adaptation



Recall: Fine-tuning



Universal Language Model Fine-Tuning for Text Classification. Howard, Ruder. '18

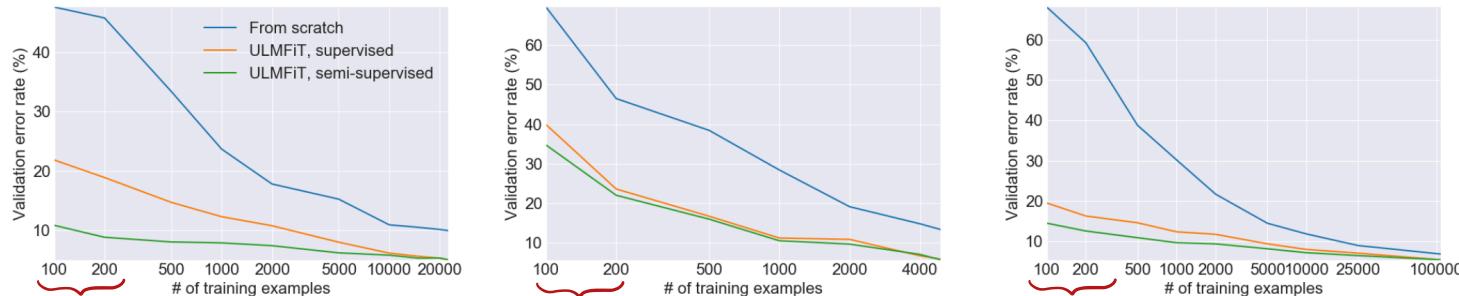


Figure 3: Validation error rates for supervised and semi-supervised ULMFiT vs. training from scratch with different numbers of training examples on IMDb, TREC-6, and AG (from left to right).

Fine-tuning less effective with very small datasets.

why people came up
with meta-learning

→ Could look @ it this way : [Fine Tuning] → [Meta Learning]
 (plugging in FT in ML)

Optimization-Based Adaptation

Fine-tuning [test-time]

$$\phi \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}^{\text{tr}})$$

pre-trained parameters

training data for new task

Meta-learning

$$\min_{\theta} \sum_{\text{task } i} \mathcal{L}(\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}}), \mathcal{D}_i^{\text{ts}})$$

evaluate how good those parameters (for task i) are for new examples in test dataset for the task.

↪ One step of gradient descent to get parameters for task i

↪ initial parameter vector

Key idea: Over many tasks, learn parameter vector θ that transfers via fine-tuning

→ ⚡ Then optimize for this pre-trained parameters
 s.t. one step of gradient descent is generalizing well to new examples. → do this for all tasks

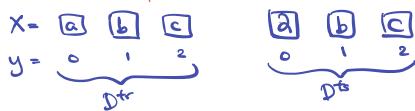
↪ has two layers of loss functions!

3-Way 1-Shot

0. Randomly initialize θ . \rightarrow "meta parameters"

1. Sample task (3 chars)

2. Sample 2 images/char \rightarrow output: meta params ϕ_i



3. Take one gradient descent step on the training dataset

$$\phi_i = \theta - 2 \nabla_{\theta} L(\theta, D_i^{\text{tr}}) \quad \text{⊗ key difference from the BlackBox.}$$

4. Update θ

Take our test examples, \rightarrow run it through our NN w/ params ϕ_i to get a corresponding prediction, \rightarrow compare the prediction w/ the corresponding label. See how well this new network is generalizing to new examples, \rightarrow back-propagate to our initial parameter vector θ to update our meta-parameters.

Q: How are θ & ϕ_i related to each other?

$\Rightarrow \theta$ is a meta starting point, ϕ_i is a set of params optimized found for a specific task.

\Rightarrow Through the above process, we're trying to determine what could have been a better initialization that would allow us to better generalize w/ the given model.

Q: Can we use different L for different tasks?

\Rightarrow Yes. Change L to L_i in the above equations.

Meta-Test

1. Given task j , D_j^{tr}

2. $\phi_j \leftarrow \theta - 2 \nabla_{\theta} L(\theta, D_j^{\text{tr}})$

$$x \rightarrow f_{\phi_j}(x) \rightarrow \hat{y}$$

\hookrightarrow "literally just running fine-tuning"

\Rightarrow What does she mean by this?

\rightarrow on D_j^{ts} \hat{y}_j ??

cf) For Black-Box approach, this was

3. $\phi_i \leftarrow f_{\phi_i}(D_i^{\text{tr}}) \rightarrow$ get the training data for a task, output the model parameters (or the hidden state-equivalent)

\rightarrow cf regression

$$\textcircled{2} \quad \textcircled{3} \quad f_{\phi_i} \rightarrow L_i(\theta_i, D_i^{\text{tr}}) = \sum_{x, y_i \in D_i^{\text{tr}}} \| \hat{y}_i - y_i \|^2$$

Q: Last time we were taking gradients on the query set, but this time we're doing it for the support set. Correct?

\Rightarrow In both Black-Box approach & the Optimization-based approach, we're updating the gradients on both the support set and the query set of the test dataset.

S1 ch1
S2 Sotokale

Q: Can we add a term for regularization?

\Rightarrow Yes, but oftentimes for Step 3 ($\phi_i = \theta - 2 \nabla_{\theta} L(\theta, D_i^{\text{tr}})$), we don't need a regularization (despite the tiny dataset used in the inner loop), because we're optimizing explicitly for generalization, it's optimizing for part of the optimization landscape that's well-behaved.

$\textcircled{2}$ In step #3, we're taking a gradient wrt. θ , not ϕ_i

$$\phi_i = \theta - 2 \nabla_{\theta} L(\theta, D_i^{\text{tr}})$$

\Rightarrow Treating ϕ_i more as activations than weights

Optimization-Based Adaptation

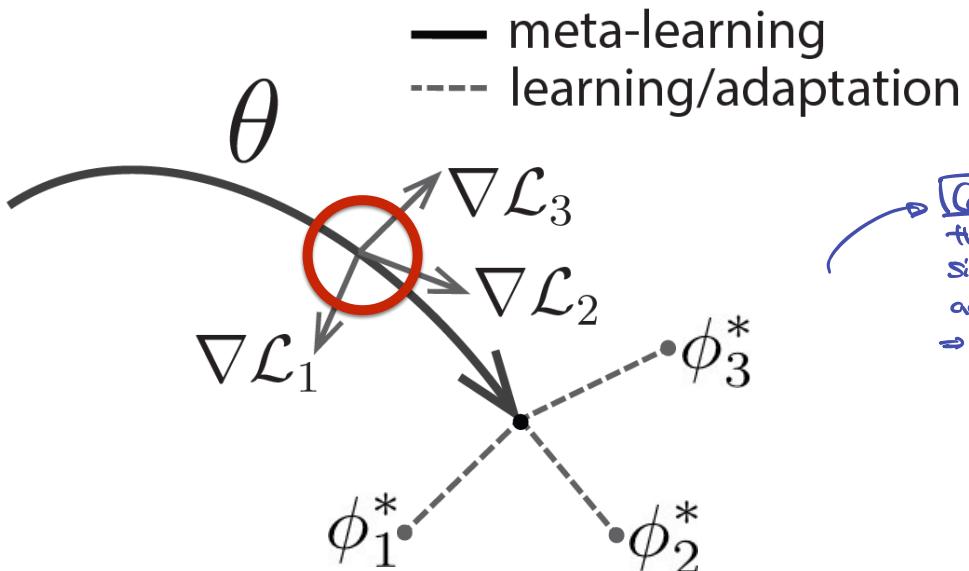
$$\min_{\theta} \sum_{\text{task } i} \mathcal{L}(\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}}), \mathcal{D}_i^{\text{ts}})$$

In practice,
this LR would be
much larger than
your typical LR
 \therefore you want to traverse
w. in the param space

- Q: Can you use hyperparam optimization to tune α ?
A: Yes, but a lot of them are designed to optimize a w. low dimensional set of hyperparams. \rightarrow won't scale well for millions of model params.
 \Rightarrow That's why we use end-to-end differentiation/gradient-based updating.

θ parameter vector
being meta-learned

ϕ_i^* optimal parameter
vector for task i



Caret
there isn't just a single optimum for any given task
 \Rightarrow More of a region.

Key Visualization for MAML!

Model-Agnostic Meta-Learning

Optimization-Based Adaptation

Key idea: Acquire ϕ_i through optimization.

General Algorithm:

~~Black box approach~~ Optimization-based approach

1. Sample task \mathcal{T}_i (*or mini batch of tasks*)
2. Sample disjoint datasets $\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{test}}$ from \mathcal{D}_i
3. ~~Compute $\phi_i \leftarrow f_\theta(\mathcal{D}_i^{\text{tr}})$~~ Optimize $\phi_i \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}})$
4. Update θ using $\nabla_\theta \mathcal{L}(\phi_i, \mathcal{D}_i^{\text{test}})$

→ brings up **second-order** derivatives

Do we need to compute the full Hessian? 😱

Both w.r.t. θ ,
not ϕ_i .

Bad bcs. N -dim params \rightarrow Hessian N^2
curse of dimensionality.

-> whiteboard

Do we get higher-order derivatives with more inner gradient steps?

Hessian-Vector Products

→ can compute much more cheaply than computing Hessian itself.
→ Goal: Compute $H(x)v$

↪ PyTorch uses this

Function: f

fs gradient: g

$$g(x+\alpha x) \approx g(x) + H(x)\alpha x \rightarrow \text{1st order Taylor Expansion approximation}$$

replace

$$g(x+rv) \approx g(x) + r H(x)v$$

$H(x)v \approx \frac{g(x+rv) - g(x)}{r}$ → can approximately compute Hessian-vector product

w/ just two gradient evals

→ Just a few more backward passes on the NN

Notation

$$\frac{d}{dx} : \text{full deriv.} \quad \nabla_x \Big|_{x=x} : \text{partial deriv.}$$

$$\phi_i = \theta - \alpha \frac{d}{d\theta} \mathcal{L}(\theta, D_i^{tr})$$

MAML objective

$$\min_{\theta} \mathcal{L}(\phi_i, D_i^{tr})$$

$$\begin{aligned} & \text{chain rule} \\ &= \frac{d}{d\theta} \mathcal{L}(\phi_i, D_i^{tr}) \\ &= \nabla_{\phi} \mathcal{L}(\phi_i, D_i^{tr}) \Big|_{\phi=\phi_i} \frac{d\phi_i}{d\theta} \\ & \quad \text{row vector } v \quad \text{matrix} \\ &= vI - \alpha V H \end{aligned}$$

from here

$$\begin{aligned} & \text{(2-step case example)} \\ & \text{For multiple steps, } \underbrace{\phi_i = \theta - \alpha \frac{d}{d\theta} \mathcal{L}(\theta, D_i^{tr}) - \alpha \frac{d}{d\theta} \mathcal{L}(\phi_i, D_i^{tr})}_{\theta'} \\ & \quad \text{this is w.r.t. } \theta' \text{, not } \theta \end{aligned}$$

$$\begin{aligned} \frac{d\phi_i}{d\theta} &= I - \alpha \frac{d^2}{d\theta^2} \mathcal{L}(\theta, D_i^{tr}) - \alpha \nabla_{\phi}^2 \mathcal{L}(\phi_i, D_i^{tr}) \Big|_{\phi=\phi_i} \frac{d\phi_i}{d\theta} \\ & \quad \text{added term for 2nd chain} \\ & \Rightarrow \text{Products, no 3rd order deriv terms.} \end{aligned}$$



$$\begin{aligned} & \frac{d}{d\theta} \mathcal{L}(\phi_i, \mathcal{D}_i^{\text{ts}}) \\ &= \nabla_{\bar{\phi}} \mathcal{L}(\bar{\phi}, \mathcal{D}_i^{\text{ts}})|_{\bar{\phi}=\phi_i} \frac{d\phi_i}{d\theta} \\ &= \nabla_{\bar{\phi}} \mathcal{L}(\bar{\phi}, \mathcal{D}_i^{\text{ts}})|_{\bar{\phi}=\phi_i} \left(I - \alpha \frac{d^2}{d\theta^2} \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}}) \right) \end{aligned}$$

PyTorch

Deep learning libraries handle the math for you.

Optimization-Based Adaptation

Key idea: Acquire ϕ_i through optimization.

Meta-Test Time:

Optimization-based approach

1. Given task \mathcal{T}_j
2. Given training data $\mathcal{D}_j^{\text{tr}}$
3. Fine-tune $\phi_j \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_j^{\text{tr}})$
4. Make predictions on new datapoints $f_{\phi_j}(x)$

Plan for Today

Recap

- Meta-learning problem & black-box meta-learning

Optimization Meta-Learning

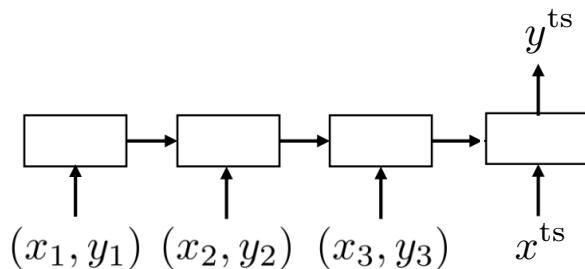
} Part of Homework 2!

- Overall approach
- **Compare: optimization-based vs. black-box**
- Challenges & solutions
- Case study of land cover classification (time-permitting)

Optimization vs. Black-Box Adaptation

Black-box adaptation

general form: $y^{\text{ts}} = f_{\text{black-box}}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}})$



Model-agnostic meta-learning

$$\begin{aligned} y^{\text{ts}} &= f_{\text{MAML}}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}}) \\ &= f_{\phi_i}(x^{\text{ts}}) \end{aligned}$$

where $\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}})$

MAML can be viewed as computation graph, \otimes
with embedded gradient operator interesting...

Note: Can mix & match components of computation graph

Learn initialization but replace gradient update with learned network

$$\begin{aligned} \text{where } \phi_i &= \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}}) \\ &f(\theta, \mathcal{D}_i^{\text{tr}}, \nabla_{\theta} \mathcal{L}) \end{aligned} \quad \text{→ Interesting, but in practice, MAML is preferred.}$$

Ravi & Larochelle ICLR '17

(actually precedes MAML)

This **computation graph view** of meta-learning will come back again!

All trained on the original Omniglot dataset,
but tested on "modified" datasets \rightarrow warped.

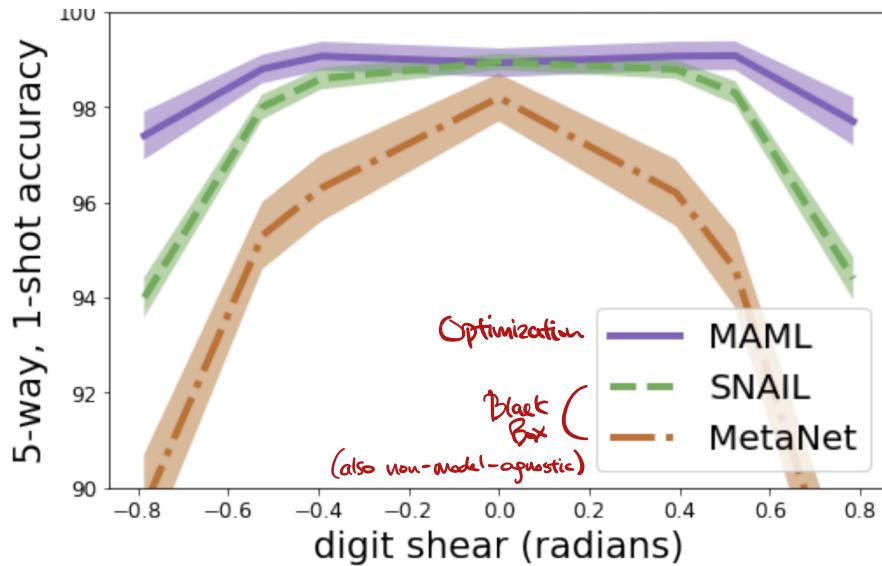
\Rightarrow MAML generalizes better than other Black-Box meta learners

Optimization vs. Black-Box Adaptation

How well can learning procedures generalize to similar, but extrapolated tasks?

Omniglot image classification

performance



Optimization
(also non-model-agnostic)

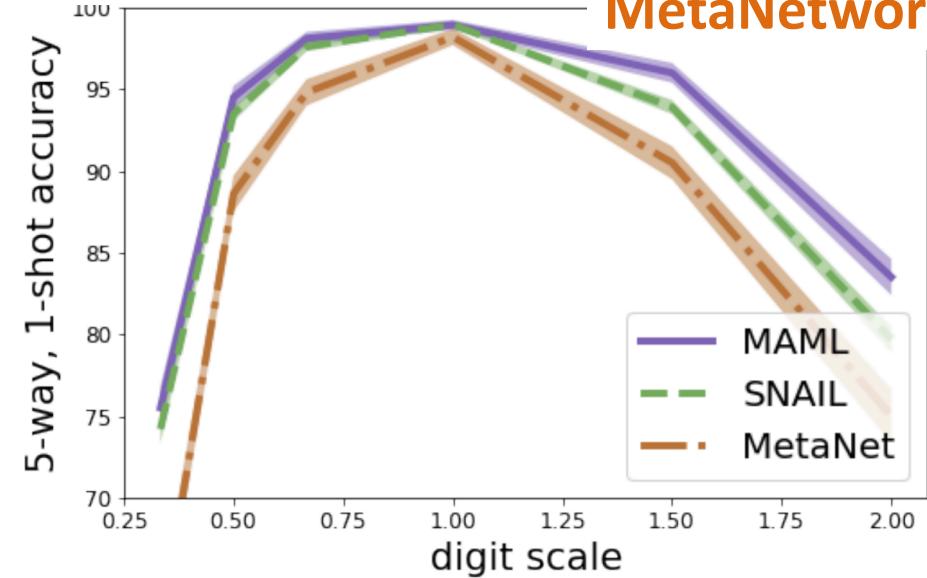
Black Box



$\xrightarrow{\text{of using gradient descent}}$

Does this structure come at a cost?

MAML SNAIL, MetaNetworks



Black-box adaptation

$$y^{\text{ts}} = f_{\text{black-box}}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}})$$

Optimization-based (MAML)

$$y^{\text{ts}} = f_{\text{MAML}}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}})$$

Does this structure come at a cost?

(For a sufficiently deep network,) *Why is this the condition?*

MAML function can approximate any function of $\mathcal{D}_i^{\text{tr}}, x^{\text{ts}}$



Finn & Levine, ICLR 2018

Assumptions:

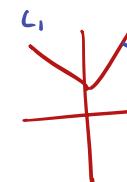
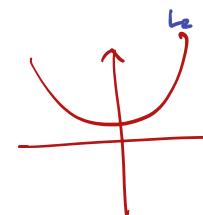
- diff. gradients?* $\cancel{\text{X}}$ - nonzero α

L_2 , CrossEntropy \vdots
 L_1 \vdots

If you have duplicates, - datapoints in $\mathcal{D}_i^{\text{tr}}$ are unique
or you want it to give you different f_i 's for diff. numbers of duplicates in the dataset, it's hard to understand how many duplicates there are in the dataset.

- $\cancel{\text{X}}$ - loss function gradient does not lose information about the label

\Rightarrow two examples of loss functions



If you just look @ the gradient of this function, you can't tell where you are \Rightarrow no info

Why is this interesting?

MAML has benefit of inductive bias without losing expressive power.

Set of assumptions a learning algorithm uses to predict outputs of given inputs that it has not encountered.

Plan for Today

Recap

- Meta-learning problem & black-box meta-learning

Optimization Meta-Learning

} Part of Homework 2!

- Overall approach
- Compare: optimization-based vs. black-box
- **Challenges & solutions**
- Case study of land cover classification (time-permitting)

Optimization-Based Adaptation Good set of references!

Challenges. Bi-level optimization can exhibit instabilities.

Ways to
stabilize
the
process

Idea: Automatically learn inner vector learning rate, tune outer learning rate
(Li et al. Meta-SGD, Behl et al. AlphaMAML)

Idea: Optimize only a subset of the parameters in the inner loop
(Zhou et al. DEML, Zintgraf et al. CAVIA)

Idea: Decouple inner learning rate, BN statistics per-step (Antoniou et al. MAML++)

Idea: Introduce context variables for increased expressive power.
(Finn et al. bias transformation, Zintgraf et al. CAVIA)

 **Takeaway:** a range of simple tricks that can help optimization significantly

Optimization-Based Adaptation

$$\textcircled{X} \frac{d\phi_i}{d\theta} = I - \alpha \frac{d^2}{d\theta^2} \mathcal{L}(\theta, D_i) \\ \text{--- this is set to zero.}$$

Challenges. Backpropagating through many inner gradient steps is

compute- & memory-intensive.

Assuming the loss landscape of ϕ_i is v. similar to that of θ → unclear why this works. Maybe because they're both well-conditioned?

Idea: [Crudely] approximate $\frac{d\phi_i}{d\theta}$ as identity \textcircled{X}

Started off as
a Tensorflow bug :-)

(Finn et al. first-order MAML '17, Nichol et al. Reptile '18)

Surprisingly works for simple few-shot problems, but (anecdotally) not for more complex meta-learning problems.

Idea: Only optimize the *last layer* of weights.

ridge regression, logistic regression

(Bertinetto et al. R2-D2 '19)

→ much cheaper computationally.

→ There's another paper that says optimizing everything but the last layer yields good results. By Boyle(!) et al.

support vector machine

(Lee et al. MetaOptNet '19)

→ leads to a *closed form* or *convex* optimization on top of meta-learned features

Idea: Derive meta-gradient using the implicit function theorem

(Rajeswaran, Finn, Kakade, Levine. Implicit MAML '19)

→ compute full meta-gradient *without differentiating through optimization path*

Optimization-Based Adaptation

Challenges. How to choose architecture that is effective for inner gradient step?

Idea: Progressive neural architecture search + MAML
(Kim et al. Auto-Meta)

- finds highly non-standard architecture (deep & narrow)
 - different from architectures that work well for standard supervised learning

Optimization-Based Adaptation

Key idea: Acquire ϕ_i through optimization.

Takeaways: Construct *bi-level optimization* problem.

- ∴ ↗ + positive inductive bias at the start of meta-learning
- ↗ + tends to extrapolate better via structure of optimization
- ↗ + maximally expressive with sufficiently deep network
- ↗ + model-agnostic (easy to combine with your favorite architecture)
- typically requires second-order optimization
- usually compute and/or memory intensive

↗ - expensive ~ challenging

Before you do any meta learning,
you're already running gradient descent
on your training dataset on the
inner loop → already have a pretty
good starting point.
cf) Black-Box starts off w/ a
randomly initialized network.

Plan for Today

Recap

- Meta-learning problem & black-box meta-learning

Optimization Meta-Learning

} Part of Homework 2!

- Overall approach
- Compare: optimization-based vs. black-box
- Challenges & solutions
- **Case study of land cover classification** (time-permitting)

Case Study



이거 구현해보면 재미있을듯.

Q: What was the core meta learner?
MAWL?

Meta-Learning for Few-Shot Land Cover Classification

Marc Rußwurm^{1,*†}, Sherrie Wang^{2,3,*}, Marco Körner¹, and David Lobell²

¹Technical University of Munich, Chair of Remote Sensing Technology

²Stanford University, Center on Food Security and the Environment

³Stanford University, Institute for Computational and Mathematical Engineering

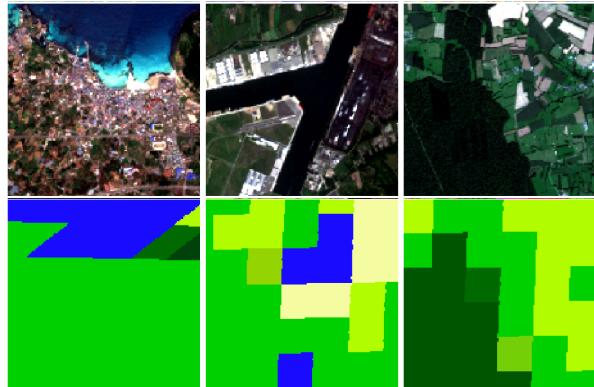
CVPR 2020 EarthVision Workshop

Link: <https://arxiv.org/abs/2004.13390>

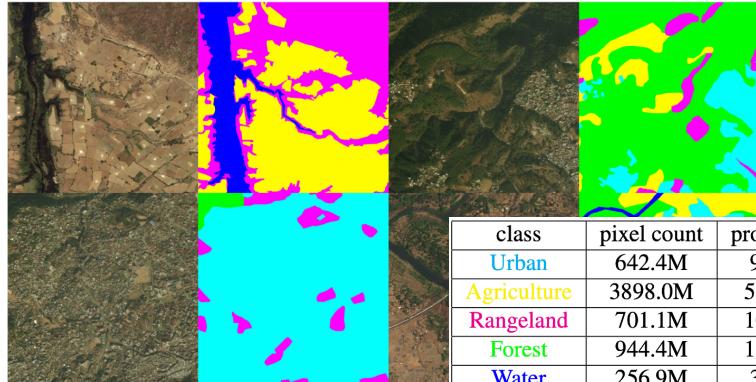
Predict how the land is being used.

Problem: Map land covering from satellite images

SEN12MS dataset
(Schmitt et al. 2019)



DeepGlobe dataset
(Demir et al. 2018)



Applications in global urban planning, climate change research

Challenges:

- Labeling data is expensive. :-
- Different regions look different & have different land use proportions
- Transfer learning may not generalize well.*

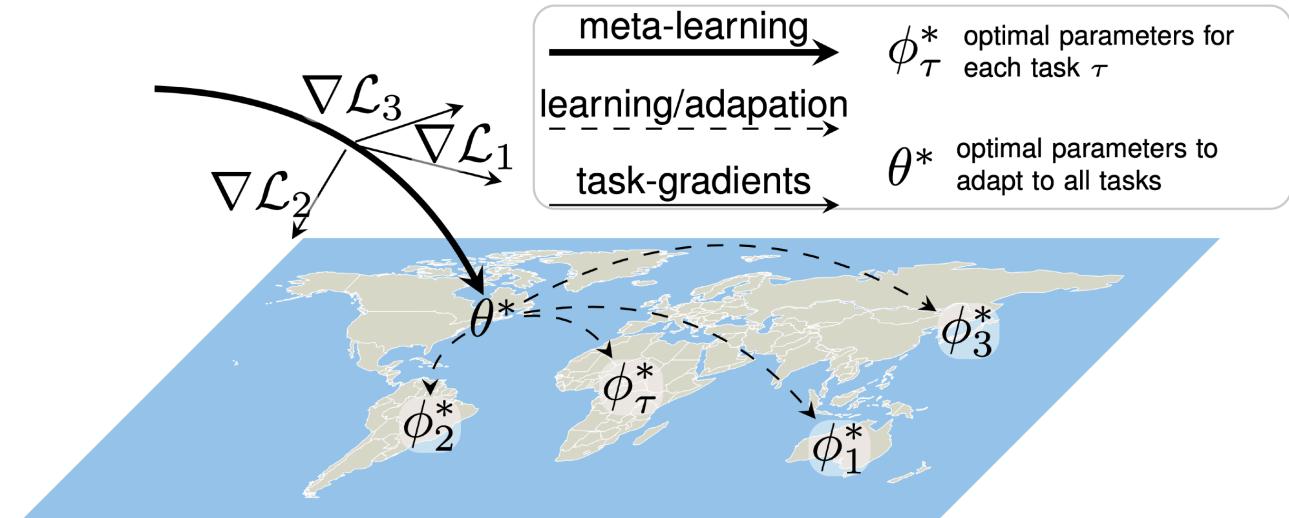
Framing land cover mapping as a meta-learning problem



Different tasks: different regions of the world

Goal: Segment/classify images from a new region with a small amount of data

The beauty of meta-learning



Croplands from four countries.

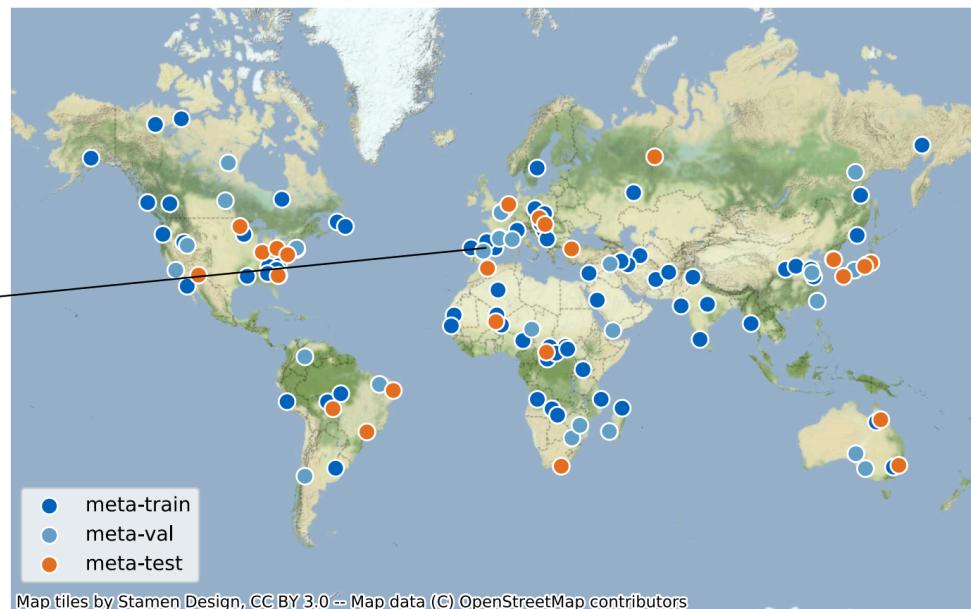
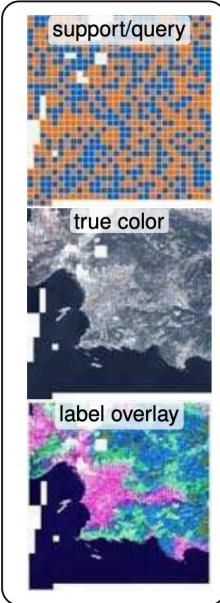
Framing land cover mapping as a meta-learning problem

Goal: Segment/classify images from a new region with a small amount of data

SEN12MS dataset (Schmitt et al. 2019)

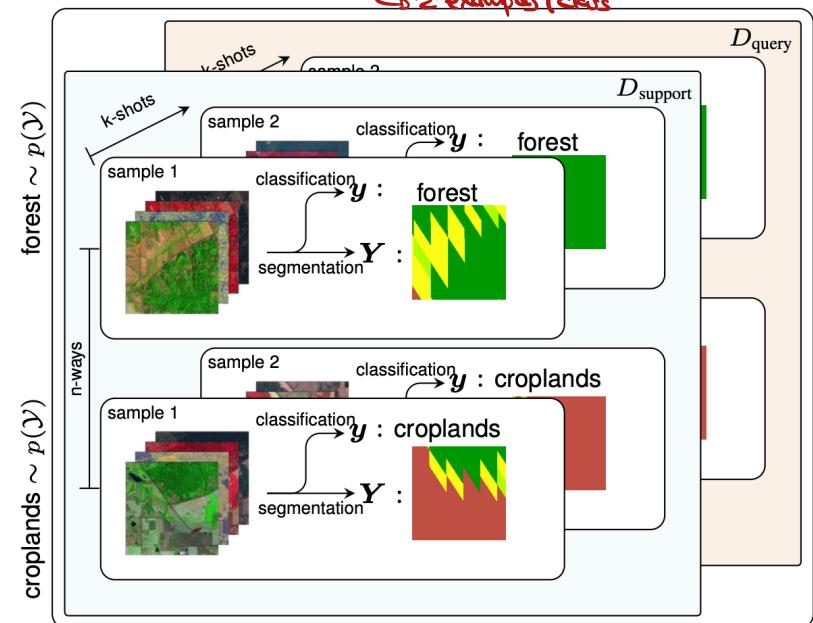
Geographic meta-data provided

Marseille (summer)



Example 2-way 2-shot classification task

Forest or crop land?
2 examples (clues)

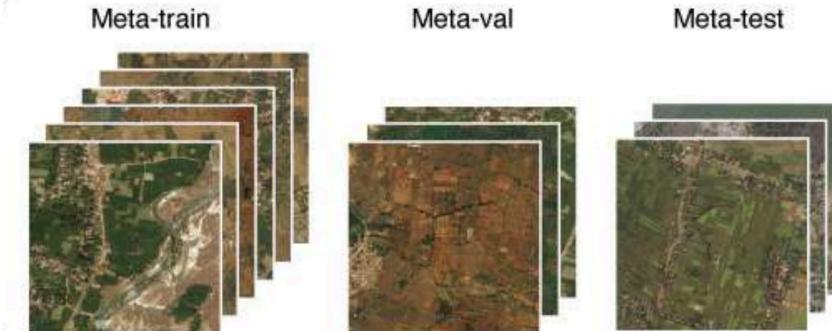


Framing land cover mapping as a meta-learning problem

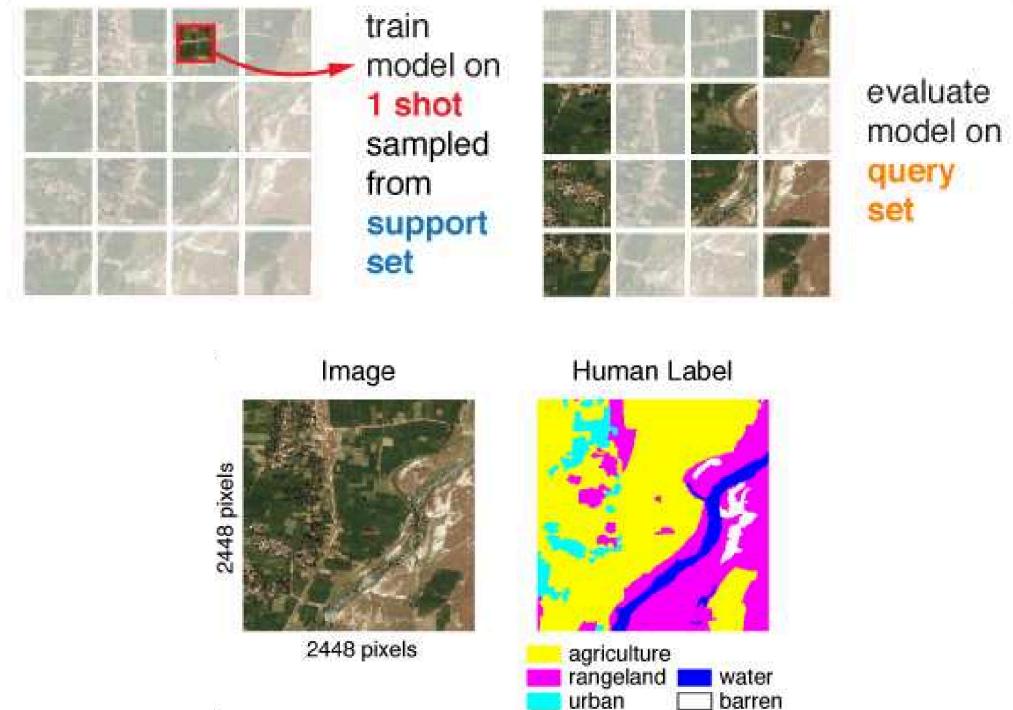
Goal: Segment/classify images from a new region with a small amount of data

DeepGlobe dataset (Demir et al. 2018)

No geographic metadata, used clustering to guess region



Example 1-shot learning segmentation task.



Evaluation

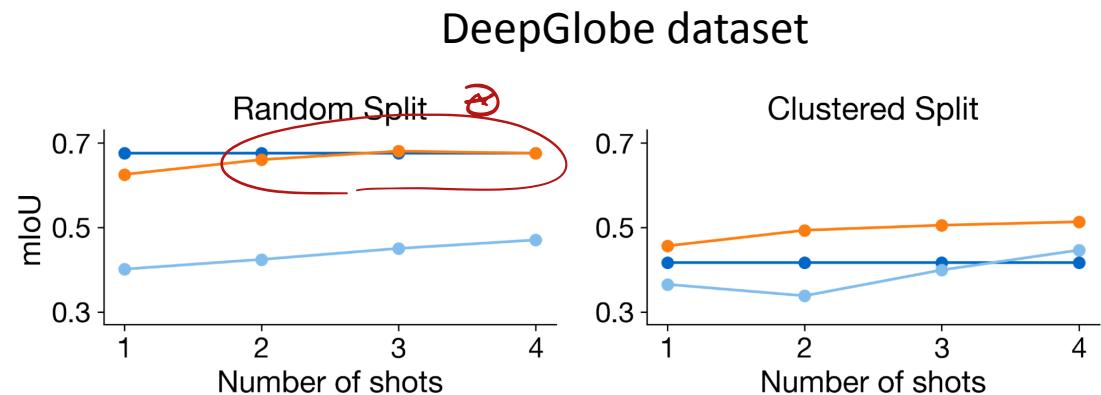
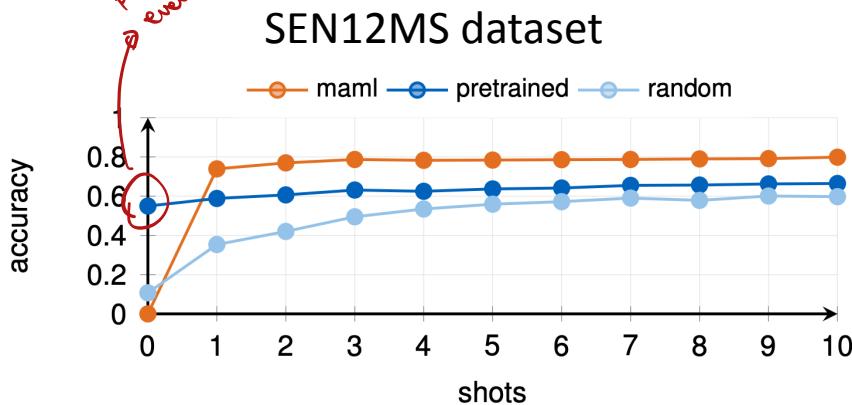
Meta-training data: $\{\mathcal{D}_1, \dots, \mathcal{D}_T\}$

Meta-test time: small amount of data from new region: $\mathcal{D}_j^{\text{tr}}$
(meta-test training set / meta-test support set)

Random init: Train from scratch on $\mathcal{D}_j^{\text{tr}}$

Compare: **Pre-train** on meta-training data $\mathcal{D}_1 \cup \dots \cup \mathcal{D}_T$, fine-tune on $\mathcal{D}_j^{\text{tr}}$

MAML on meta-training data $\{\mathcal{D}_1, \dots, \mathcal{D}_T\}$, adapt with $\mathcal{D}_j^{\text{tr}}$



More visualizations and analysis in the paper!

④ Makes sense: pre-training ~ target data
are more similar

Plan for Today

Recap

- Meta-learning problem & black-box meta-learning

Optimization Meta-Learning

} Part of Homework 2!

- Overall approach
- Compare: optimization-based vs. black-box
- Challenges & solutions
- Case study of land cover classification (time-permitting)

Goals for by the end of lecture:

- Basics of optimization-based meta-learning techniques (& how to implement)
- Trade-offs between black-box and optimization-based meta-learning

Roadmap for upcoming lectures

Wednesday: Non-parametric few-shot learners, comparison of approaches

Next week: Unsupervised pre-training for few-shot learning

Following week: Advanced meta-learning topics (e.g. memorization, large-scale meta-optimization)

Course Reminders

Project group form due **tonight.**
(for assigning project mentors)

Homework 1 due **Wednesday**