

Meta Learning

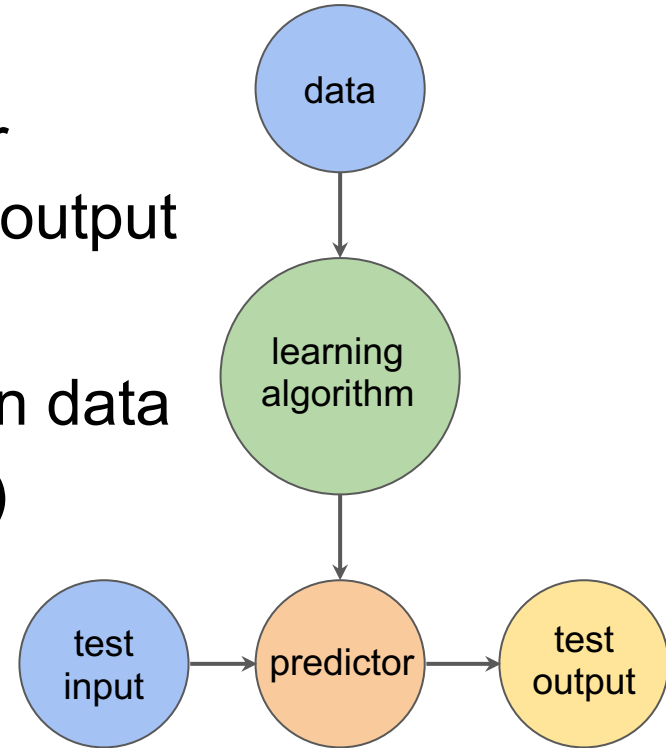
MIT

Iddo Drori, Fall 2020

Supervised Learning

- Data
- Learning algorithm for finding predictor
- For test input, predictor estimates test output

- Inference of predictor parameters given data
 $p(\text{parameters} \mid \text{data})$



Supervised Learning

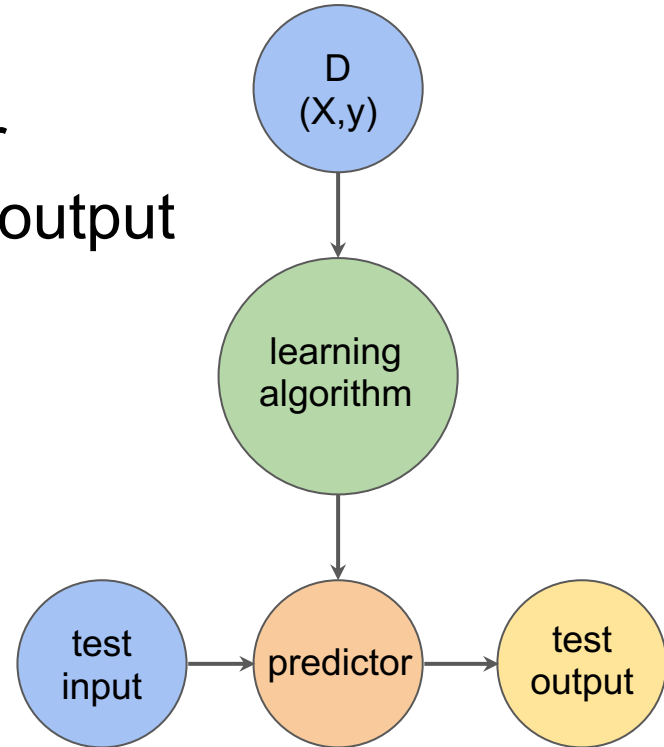
- Data
- Learning algorithm for finding predictor
- For test input, predictor estimates test output

- Inference of ϕ given D : $p(\phi | D)$

$$\phi^* = \operatorname{argmax}_{\phi} \log p(\phi | D)$$

$$= \operatorname{argmax}_{\phi} \log p(D | \phi) + \log p(\phi)$$

$$= \operatorname{argmax}_{\phi} \sum_i \log p(y_i | X_i, \phi) + \log p(\phi)$$



Meta Learning

- Learn meta parameter θ given meta training data D

$$p(\theta | D)$$

$$\theta^* = \operatorname{argmax}_{\theta} \log p(\theta | D)$$

Meta Learning for Few Shot Classification

- Adaptation

$$p(\phi \mid D_t, \theta^*)$$

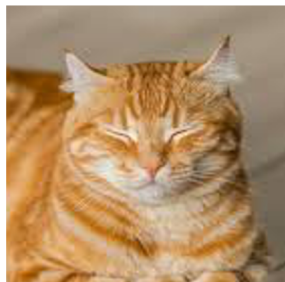
$$\phi^* = \operatorname{argmax}_{\phi} \log p(\phi \mid D_t, \theta^*)$$

Meta Learning for Few Shot Classification

Meta Learning for Few Shot Classification

- Training task 1: cats vs. dogs

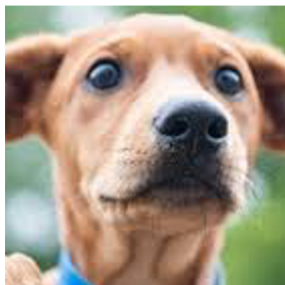
cat



?



dog



1

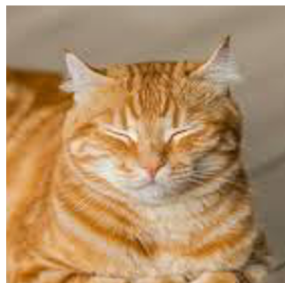
2

3

Meta Learning for Few Shot Classification

- Training task 1: 2 way (classes), 3 shot (samples)

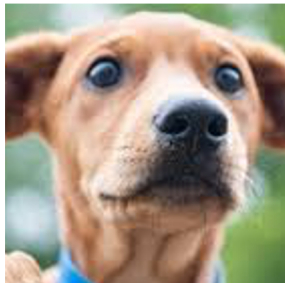
class 1



?



class 2



sample 1

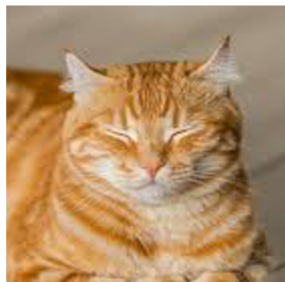
sample 2

sample 3

Meta Learning for Few Shot Classification

- Training task 1: $c = 2$ classes, $k = 3$ samples

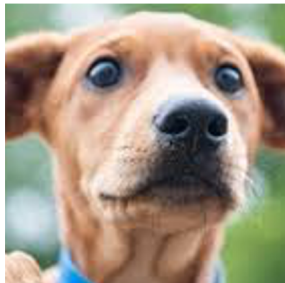
class 1



?



class 2



sample 1

sample 2

sample 3

Meta Learning for Few Shot Classification

- Training task 2: flower vs. bird

flower



?



bird



1

2

3

Meta Learning for Few Shot Classification

- Testing task: lion vs. monkey

lion



?



monkey

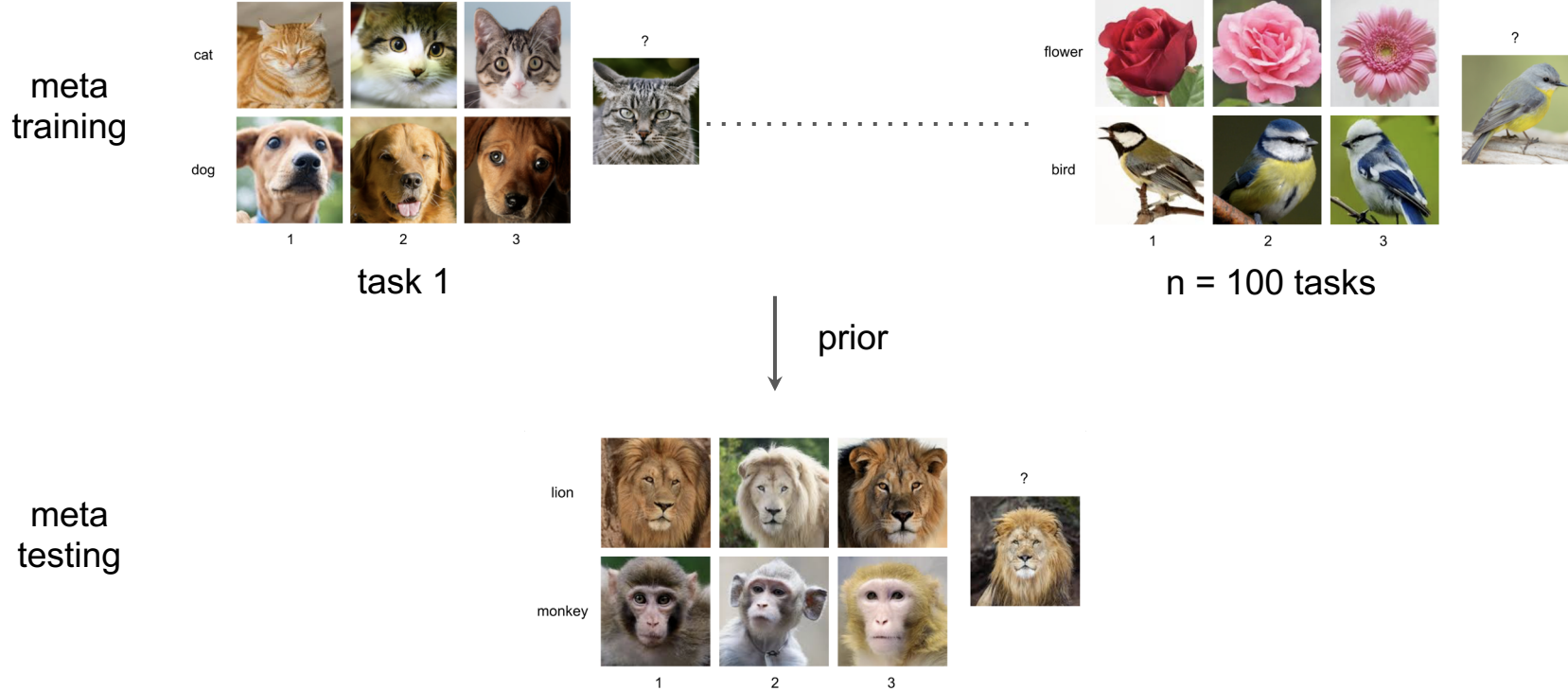


1

2

3

Meta Learning for Few Shot Classification



Meta Learning for Few Shot Classification

- c classes
- k samples per class for training
- n tasks for meta training

Task Support and Query Sets

- For each task i with meta training dataset $D_i = D_{s_i} \cup D_{q_i}$
 - Training set D_{s_i} (support set)
 - Testing set D_{q_i} (query set)

task i

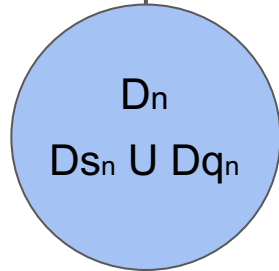
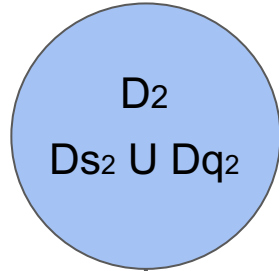
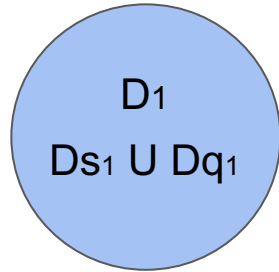


D_{s_i} support set



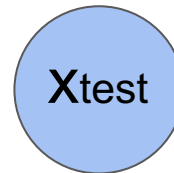
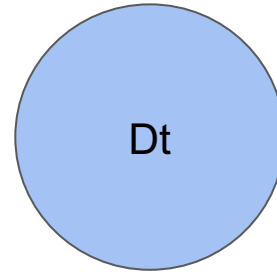
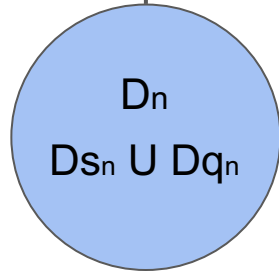
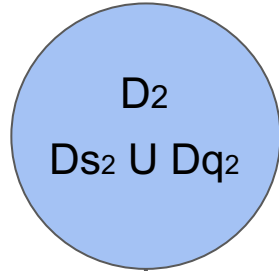
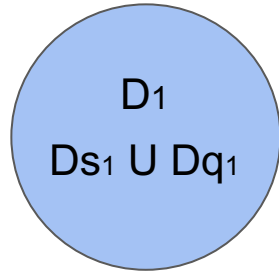
D_{q_i} query set

Meta Data



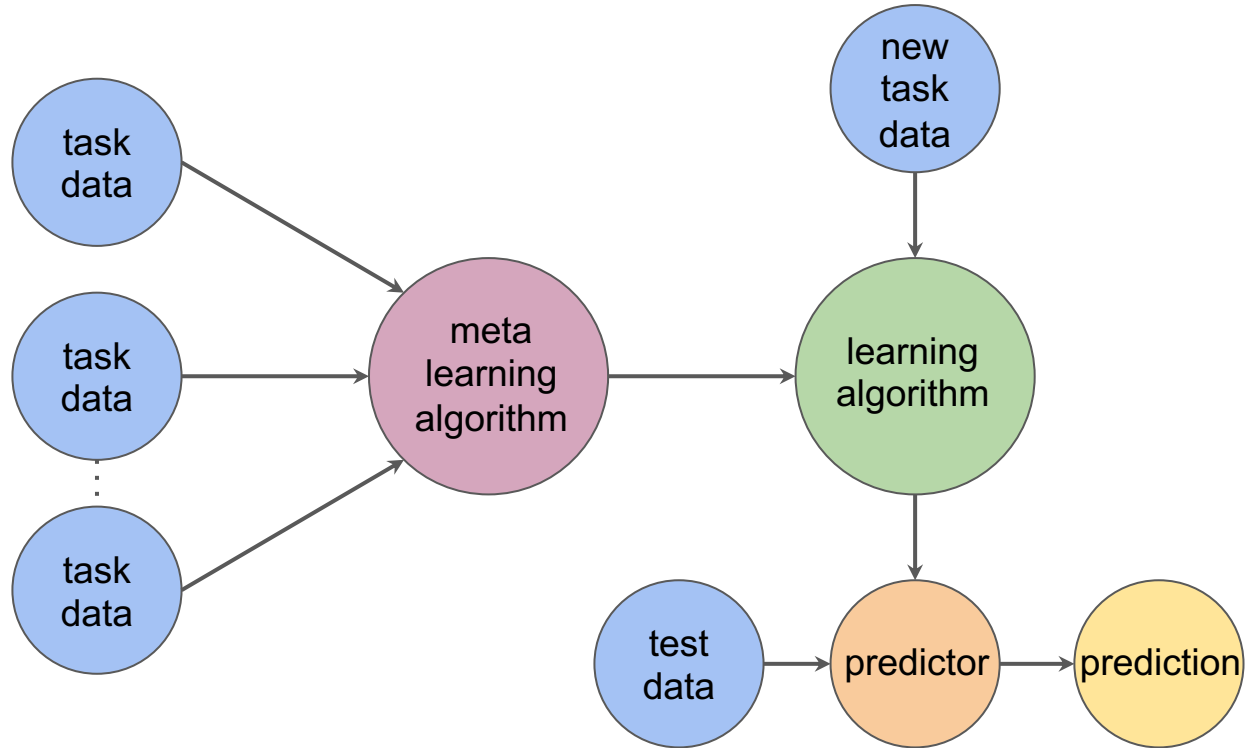
$$D = (D_1, \dots, D_n)$$

Meta Learning



Meta Learning

task = data splits, priors



Black-Box Methods

Black-Box Meta Learning for Few Shot Classification

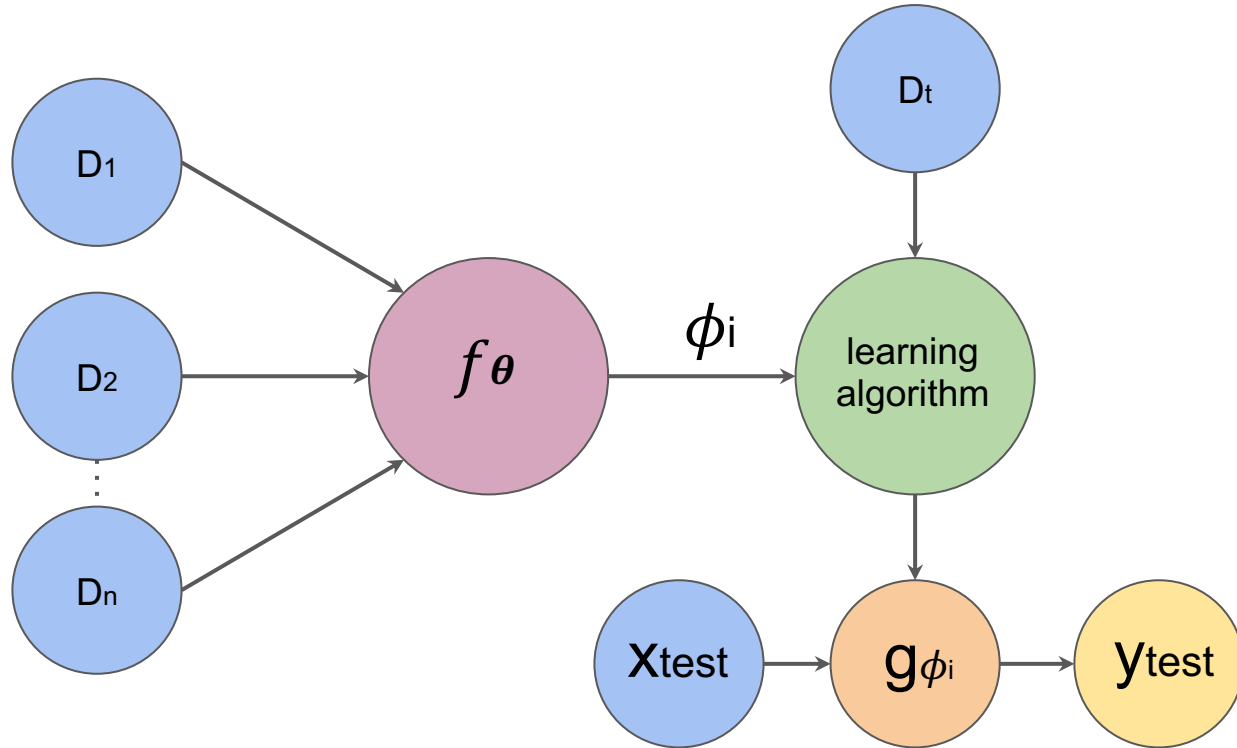
- Meta training data: $D = (D_1, \dots, D_n)$
- Inference over task specific parameters ϕ_i given meta training dataset and meta parameters

$$p(\phi_i | D_{s_i}, \theta)$$

$$\max_{\theta} \sum_i \log(\phi_i | D_{q_i})$$

Black-Box Meta Learning for Few Shot Classification

task = data splits, priors



Meta Learning for Few Shot Classification

- $p(\phi_i | D_{s_i}, \theta)$
- Optimize θ MLE using meta training dataset D
- Model as $p(\phi_i | D_{s_i}, \theta)$ as NN f_θ
- Meta NN f_θ with input D_i and output ϕ_i

$$\phi_i = f_\theta(D_{s_i})$$

- Second task specific NN g with parameters ϕ_i computing

$$y_{\text{test}} = g_{\phi_i}(x_{\text{test}})$$

- $\max_{\theta} \sum_i \sum_{(x,y) \sim D_{q_i}} \log g_{\phi_i}(y|x)$
- $\max_{\theta} \sum_i \mathcal{L}(f_\theta(D_{s_i}), D_{q_i})$

Meta Learning Algorithm for Few Shot Classification

- Sample task i
- Sample task i dataset $D_i = D_{s_i} \cup D_{q_i}$:
 - Training set D_{s_i} (support set)
 - Testing set D_{q_i} (query set)
- Compute $\phi_i = f_{\theta}(D_{s_i})$
- Update θ by $\nabla_{\theta} \mathcal{L}(\phi_i, D_{q_i})$

Gradient-based Methods

Gradient-based Inference

- Meta model parameters θ is a prior, model initialization
- For each task i : task adapted parameter ϕ_i

$$\max_{\theta} \log p(D_{s_i} | \phi_i) + \log p(\phi_i | \theta)$$

Gradient-based Inference

- Meta model parameters θ is a prior, model initialization
- For each task i : task adapted parameter ϕ_i
- Fine tuning
- Initialization with pre-trained parameters θ
 - CNN parameters trained on image dataset
 - Transformer parameters trained on text corpus
- Training data for new task D_t

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_t)$$

Gradient-based Bi-Level Optimization

- Meta model parameters θ is a prior, model initialization
- Optimize θ across many tasks so fine tuning does well
- For each task i : task adapted parameter ϕ_i

$$\min_{\theta} \frac{1}{n} \sum_i \mathcal{L}_i(\phi_i, Dq_i)$$

$$\phi_i = \text{algorithm}(\theta, Ds_i)$$

$$\min_{\theta} \frac{1}{n} \sum_i \mathcal{L}_i(\text{algorithm}(\theta, Ds_i), Dq_i)$$

Model Agnostic Meta Learning (MAML)

- Meta training

$$\min_{\theta} \frac{1}{n} \sum_i \mathcal{L}_i(\phi_i, D_{q_i})$$

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_{s_i})$$

$$\min_{\theta} \frac{1}{n} \sum_i \mathcal{L}_i(\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_{s_i}), D_{q_i})$$

- Meta testing
- D_s : training data of new task
- θ^* : pre-trained parameters

$$\phi = \theta^* - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_s)$$

Meta Algorithm

- Sample task i
- Sample task i dataset $D_i = D_{s_i} \cup D_{q_i}$:
 - Training set D_{s_i} (support set)
 - Testing set D_{q_i} (query set)
- Optimize $\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_{s_i})$
- Update θ by $\nabla_{\theta} \mathcal{L}(\phi_i, D_{q_i}) = \nabla_{\theta} \mathcal{L}(\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_{s_i}), D_{q_i})$

Gradient-based Meta Learning

- Meta training

$$\min_{\theta} \frac{1}{n} \sum_i \mathcal{L}_i(\phi_i, D_{q_i})$$

- Update algorithm

$$\phi_i = \text{algorithm}(\theta, D_{s_i})$$

- Meta testing
- D_s : training data of new task
- θ^* : pre-trained parameters

$$\phi = \theta^* - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_s)$$

Gradient-based Meta Learning

- Meta training

$$\min_{\theta} \frac{1}{n} \sum_i \mathcal{L}_i(\phi_i, D_{q_i})$$

$$\phi_i = \text{algorithm}(\theta, D_{s_i})$$

- MAML $= \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_s)$
- MetaSGD $= \theta - \alpha \text{diag}(w) \nabla_{\theta} \mathcal{L}(\theta, D_s)$
- Tnet $= \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, w, D_s)$
- Meta curvature $= \theta - \alpha B(\theta, w) \nabla_{\theta} \mathcal{L}(\theta, D_s)$
- Wrap-grad $= \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, w, D_s)$

Gradient-based Meta Learning

- Second order derivatives

$$\min_{\theta} \mathcal{L}(\phi, Dq_i)$$

$$\phi = \text{algorithm}(\theta, D_s)$$

$$\min_{\theta} \mathcal{L}(\text{algorithm}(\theta, D_s), Dq_i)$$

$$d_{\theta} \mathcal{L}(\phi, Dq_i) = \nabla_{\theta} \mathcal{L}(a, Dq_i)|_{a=\text{algorithm}(\theta, D_s)} d_{\theta} \text{algorithm}(\theta, D_s)$$

Gradient-based Meta Learning

- Second order derivatives

$$\min_{\theta} \mathcal{L}(\phi, Dq_i)$$

$$\phi = \text{algorithm}(\theta, D_s) = \theta - \alpha d_{\theta} \mathcal{L}(\theta, D_s)$$

$$d_{\theta} \text{algorithm}(\theta, D_s) = I - \alpha d_{\theta}^2 \mathcal{L}(\theta, D_s)$$

$$d_{\theta} \mathcal{L}(\phi, Dq_i) = \nabla_{\theta} \mathcal{L}(a, Dq_i) \Big|_{a=u(\theta, D_s)} d_{\theta} \text{algorithm}(\theta, D_s)$$

Gradient-based Meta Learning

- Second order derivatives

$$\min_{\theta} \frac{1}{n} \sum_i \mathcal{L}_i(\phi_i, Dq_i)$$

$$\phi_i = u(\theta, Ds_i)$$

$$\nabla_{\theta} \mathcal{L}(f_{\phi}, Dq) = (I - \alpha H_s(\theta)) g_q(\phi)$$

- Reptile update for θ : $\theta - \beta \frac{1}{n} (\theta - \phi_i)$

Meta Learning for Few Shot Classification

- Why not take all meta training data together with meta testing data to learn a representation from all of them together?
- This may work better than other meta learning methods.
Rethinking Few-Shot Image Classification: A Good Embedding Is All You Need?, Tian et al, 2020.

Metric-based Methods (non-parametric)

Metric-based Meta Learning

- Matching network
- Prototypical network
- Relation network
- GNN
- MetaOptNet

Naive Approach

- Compare D_{q_i} with each sample in D_{s_i}
- Label by nearest neighbor.
- Other methods?

Siamese Networks

- Are two samples from the same class?
- Training: pairwise comparisons of x_{test} with all D_{s_i}
- Binary classification
- Testing: one vs. many
- $\varphi(x_i, x_j) = \|\phi(x_i) - \phi(x_j)\|$

Matching Network

- Training on multi-class classification
- Nearest neighbors at test time
- Learn an embedding at train time such that nearest neighbors at test time provides accurate predictions
- Meta training: learn g_θ and f_θ

Similarity score $f_\theta(x_{\text{test}}, x_k)$

$$y_{\text{test}} = \sum_{(x_k, y_k) \text{ in } D_s} f_\theta(x_{\text{test}}, x_k) y_k$$

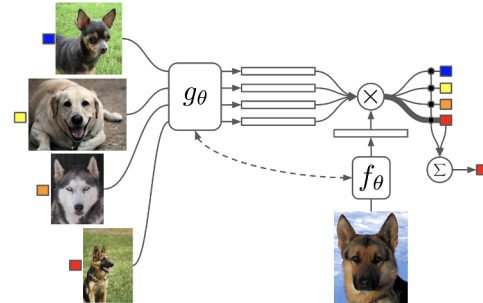


Figure source: Matching networks for one shot learning, Vinyals et al, 2016

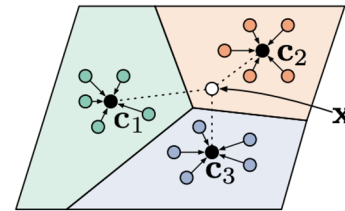
Non-Parametric Meta Learning Algorithm

- Sample task i
- Sample task i dataset $D_i = D_{s_i} \cup D_{q_i}$:
 - Training set D_{s_i} (support set)
 - Testing set D_{q_i} (query set)
- Compute $y_{\text{test}} = \sum_{(x_k, y_k) \in D_s} f_{\theta}(x_{\text{test}}, x_k) y_k$
- Update θ by $\nabla_{\theta} \mathcal{L}(y'_{\text{test}}, y_{\text{test}})$

Non-parametric, independent of ϕ

Prototypical Network

- Aggregate class information, prototypical for each class
- Embed each training image in each class and take mean
- Embed test image
- Embedding of data and nearest neighbors
- $c_k = 1/|D_{Si}| \sum_{(x,y) \text{ in } D_{Si}} f_{\theta}(x)$
- $p_{\theta}(y = k|x) = \exp(-d(f_{\theta}(x), c_k)) / \sum_k \exp(-d(f_{\theta}(x), c_k))$
- Euclidean or cosine distance



(a) Few-shot

Relation Network

- Instead of defining d (Euclidean or cosine), learn d
- Relation module

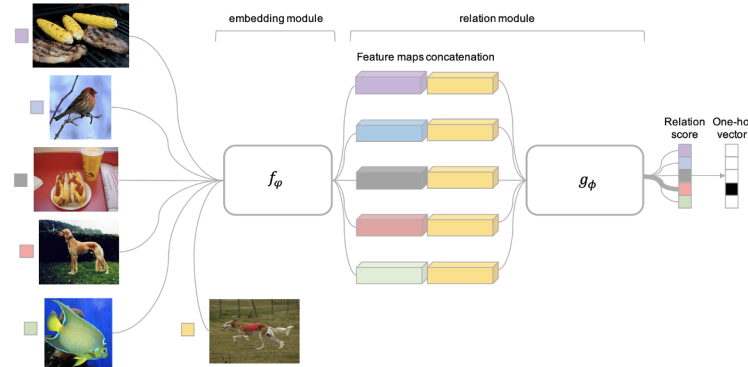


Figure source: Learning to compare: Relation network for few-shot learning, Sung et al, 2018

Graph neural network (GNN)

- Embedding using GNN

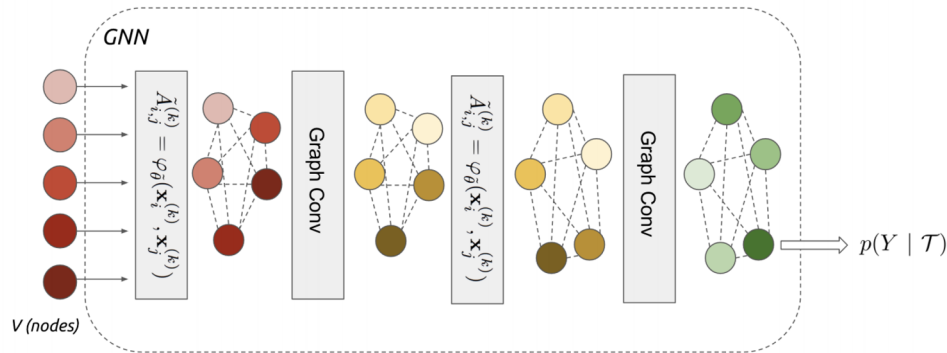


Figure source: Few-shot learning with graph neural networks, Garcia and Bruna, 2018

MetaOptNet

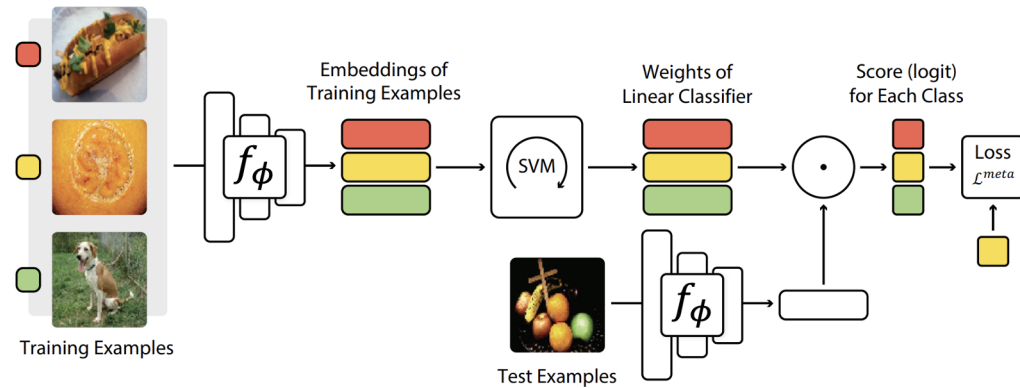


Figure source: Meta-learning with differentiable convex optimization, Lee et al, 2019

Comparison of Approaches

- Black-box: $y_{\text{test}} = f_{\theta}(D_{S_i}, x_{\text{test}})$
- Gradient-based (optimization): $y_{\text{test}} = f(D_{S_i}, x_{\text{test}})$
 $= f\phi_i(x_{\text{test}})$ where $\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_{S_i})$
- Metric-based (non-parametric): $y_{\text{test}} = f(D_{S_i}, x_{\text{test}}) =$
 $\text{softmax}(-d(f_{\theta}(x), c_k)), c_k = 1/|D_{S_i}| \sum_{(x,y) \text{ in } D_{S_i}} f_{\theta}(x)$

Comparison of Approaches

- Black-box: data intensive
- Gradient-based (optimization): classification, regression, reinforcement learning; second order, computation intensive
- Metric-based (non-parametric): classification; simple feed forward; fast; dependent on distance metric

Meta Learning

MIT

Iddo Drori, Fall 2020