

Learning to Few-Shot Learn  
Across Diverse Natural Language Classification Tasks

Trapit Bansal, Rishikesh Jha, Andrew McCallum

# Motivation

- Fine-tuning on a new task still requires large amounts of task-specific labelled data
- Generalize to new tasks with few examples as a meta-learning problem
- Enables optimization-based meta-learning across tasks with **different number of classes**

# Notation

Task: Episodic task (consist of support set and query set)

$$\{T_1, \dots, T_M\}$$

Train: Support set  $\mathcal{D}_i^{tr} \sim T_i$

Validation: Query set  $\mathcal{D}_i^{val} \sim T_i$

# Model-Agnostic Meta Learning (MAML)

Inner loop

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

Outer loop

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{T_i \sim P(\mathcal{T})} \mathcal{L}_i(\theta'_i, \mathcal{D}_i^{val})$$

# Generating Softmax Parameters

Task-dependent softmax parameters

$$w_i^n, b_i^n = \frac{1}{|C_i^n|} \sum_{x_j \in C_i^n} g_\psi(f_\theta(\mathbf{x}_j)) \quad C_i^n = \{x_j | y_j = n\}, \text{ where } n \in [N_i]$$

Prediction for a new instance  $\mathbf{x}^*$

$$p(y|\mathbf{x}^*) = \text{softmax} \{ \mathbf{W}_i h_\phi(f_\theta(\mathbf{x}^*)) + \mathbf{b}_i \}$$

# Learn to Adapt

- Task-agnostic: BERT lower layers
- Task-specific: BERT higher layers, softmax-params
- Inner loop: Update task-specific params

$$\Theta = \theta_{\leq \nu} \cup \{\psi\}$$

$$\Phi_i = \theta_{> \nu} \cup \{\phi, \mathbf{W}_i, \mathbf{b}_i\}$$

$$\Phi_i^{(s+1)} = \Phi_i^{(s)} - \alpha_s \mathbb{E}_{\mathcal{D}_i^{tr} \sim T_i} [\nabla_{\Phi} \mathcal{L}_i(\{\Theta, \Phi_i\}, \mathcal{D}_i^{tr})]$$

- Outer loop: update the task-agnostic params
  - Use first-order approximation for efficient computation

---

**Algorithm 1** LEOPARD

---

**Require:** set of  $M$  training tasks and losses  $\{(T_1, L_1), \dots, (T_M, L_M)\}$ , model parameters  $\Theta =$

$\{\theta, \psi, \alpha\}$ , hyper-parameters  $\nu, G, \beta$

Initialize  $\theta$  with pre-trained BERT-base;

- 1: **while** not converged **do**
  - 2:   *# sample batch of tasks*
  - 3:   **for all**  $T_i \in T$  **do**
  - 4:      $\mathcal{D}_i^{tr} \sim T_i$    *# sample a batch of train data*
  - 5:      $C_i^n \leftarrow \{x_j | y_j = n\}$    *# partition data according to class labels*
  - 6:      $w_i^n, b_i^n \leftarrow \frac{1}{|C_i^n|} \sum_{x_j \in C_i^n} g_\psi(f_\theta(\mathcal{D}_i^{tr}))$    *# generate softmax parameters*
  - 7:      $\mathbf{W}_i \leftarrow [w_i^1; \dots; w_i^{N_i}]; \mathbf{b}_i \leftarrow [b_i^1; \dots; b_i^{N_i}]$
  - 8:      $\Phi_i^{(0)} \leftarrow \theta_{>\nu} \cup \{\phi, \mathbf{W}_i, \mathbf{b}_i\}$    *# task-specific parameters*
  - 9:     **for**  $s := 0 \dots G - 1$  **do**
  - 10:        $\mathcal{D}_i^{tr} \sim T_i$    *# sample a batch of train data*
  - 11:        $\Phi_i^{(s+1)} \leftarrow \Phi_i^{(s)} - \alpha_s \nabla_{\Phi} \mathcal{L}_i(\{\Theta, \Phi_i\}, \mathcal{D}_i^{tr})$    *# adapt task-specific parameters*
  - 12:     **end for**
  - 13:      $\mathcal{D}_i^{val} \sim T_i$    *# sample a batch of validation data*
  - 14:      $g_i \leftarrow \nabla_{\Theta} \mathcal{L}_i(\{\Theta, \Phi_i^{(G)}\}, \mathcal{D}_i^{val})$    *# gradient of task-agnostic parameters on validation*
  - 15:   **end for**
  - 16:    $\Theta \leftarrow \Theta - \beta \cdot \sum_i g_i$    *# optimize task-agnostic parameters*
  - 17: **end while**
-

# Evaluation

- Training and Validation
  - GLUE benchmark
- Testing:
  - Entity typing: CoNLL-2003
  - Rating classification: Amazon Review
  - Text classification:



# Results

## Entity Typing

	$N$	$k$	BERT <sub>base</sub>	MT-BERT <sub>softmax</sub>	MT-BERT	Proto-BERT	LEOPARD
CoNLL		4	50.44 $\pm$ 08.57	52.28 $\pm$ 4.06	<b>55.63</b> $\pm$ 4.99	32.23 $\pm$ 5.10	54.16 $\pm$ 6.32
		8	50.06 $\pm$ 11.30	65.34 $\pm$ 7.12	58.32 $\pm$ 3.77	34.49 $\pm$ 5.15	<b>67.38</b> $\pm$ 4.33
		16	74.47 $\pm$ 03.10	71.67 $\pm$ 3.03	71.29 $\pm$ 3.30	33.75 $\pm$ 6.05	<b>76.37</b> $\pm$ 3.08
MITR		4	49.37 $\pm$ 4.28	45.52 $\pm$ 5.90	<b>50.49</b> $\pm$ 4.40	17.36 $\pm$ 2.75	49.84 $\pm$ 3.31
		8	49.38 $\pm$ 7.76	58.19 $\pm$ 2.65	58.01 $\pm$ 3.54	18.70 $\pm$ 2.38	<b>62.99</b> $\pm$ 3.28
		16	69.24 $\pm$ 3.68	66.09 $\pm$ 2.24	66.16 $\pm$ 3.46	16.41 $\pm$ 1.87	<b>70.44</b> $\pm$ 2.89
<b>Text Classification</b>							
Airline		4	42.76 $\pm$ 13.50	43.73 $\pm$ 7.86	46.29 $\pm$ 12.26	40.27 $\pm$ 8.19	<b>54.95</b> $\pm$ 11.81
		8	38.00 $\pm$ 17.06	52.39 $\pm$ 3.97	49.81 $\pm$ 10.86	51.16 $\pm$ 7.60	<b>61.44</b> $\pm$ 03.90
		16	58.01 $\pm$ 08.23	58.79 $\pm$ 2.97	57.25 $\pm$ 09.90	48.73 $\pm$ 6.79	<b>62.15</b> $\pm$ 05.56

# Results

$k$	Model	Entity Typing	Sentiment Classification	NLI
16	LEOPARD <sub>10</sub>	37.62 $\pm$ 7.37	58.10 $\pm$ 5.40	78.53 $\pm$ 1.55
	LEOPARD <sub>5</sub>	62.49 $\pm$ 4.23	71.50 $\pm$ 5.93	73.27 $\pm$ 2.63
	LEOPARD	69.00 $\pm$ 4.76	76.65 $\pm$ 2.47	76.10 $\pm$ 2.21
	LEOPARD-ZERO	44.79 $\pm$ 9.34	74.45 $\pm$ 3.34	74.36 $\pm$ 6.67

Table 3: Ablations: LEOPARD <sub>$\nu$</sub>  does not adapt layers 0 –  $\nu$  (inclusive) in the inner loop, while LEOPARD adapts all parameters. Note that the outer loop still optimizes all parameters for all models. For new tasks (like entity typing) adapting all parameters is beneficial while for tasks seen at training time (like NLI) adapting fewer parameters is better. LEOPARD-ZERO is the model without the softmax-generator and a zero initialized softmax classifier which shows the importance of the softmax generator in LEOPARD.