

PRE-TRAINING – ADAPTATION INTERFACE: (FM)

Analysis of Foundation Models

4.10 Theory

Authors: Aditi Raghunathan, Sang Michael Xie, Ananya Kumar, Niladri Chatterji, Rohan Taori, Tatsunori Hashimoto, Tengyu Ma

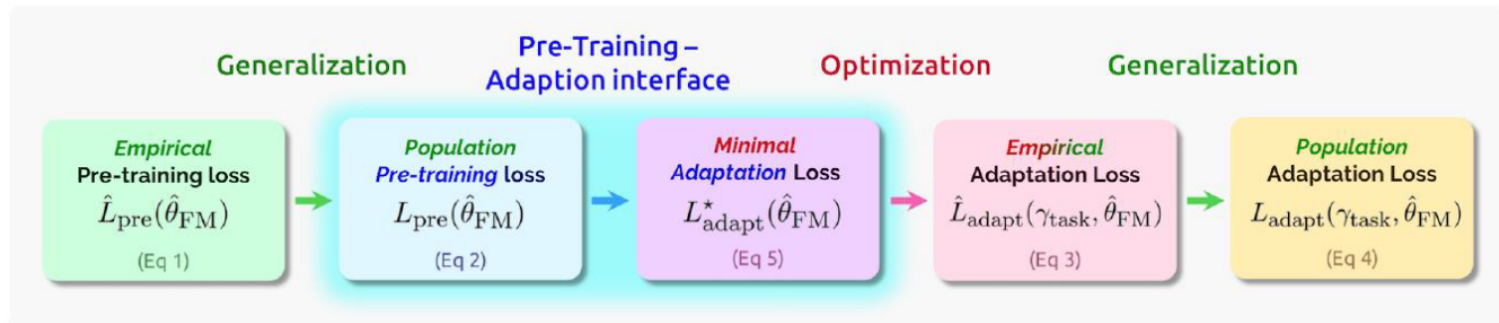
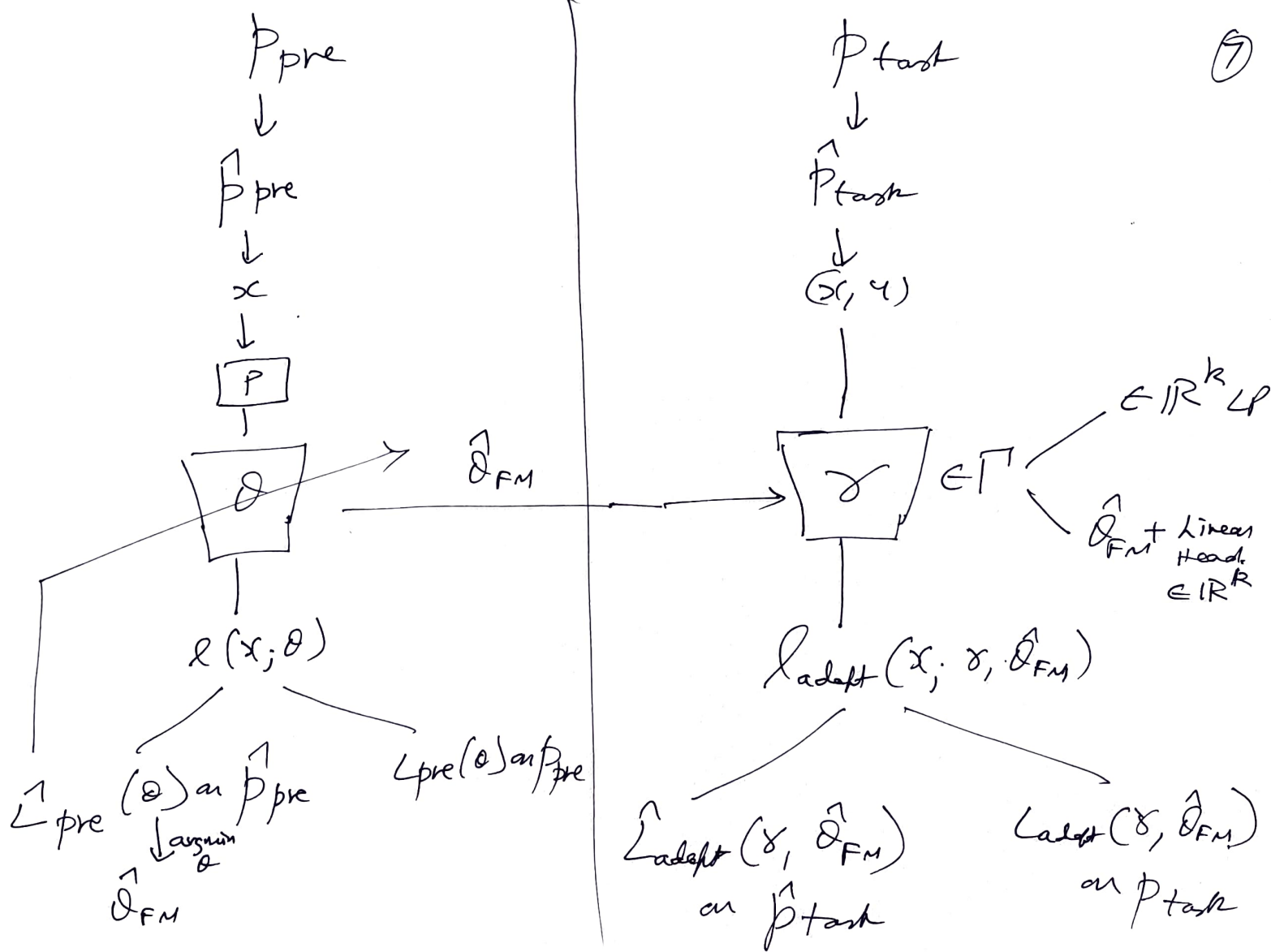


Fig. 22. The analysis of foundation models from pretraining on diverse data to downstream performance on adapted tasks involves capturing the relation between different loss terms as shown above. The main challenge is to analyze the highlighted pretraining-adaptation interface which requires reasoning carefully about the population losses in addition to the model architecture, losses and data distributions of the pretraining and adaptation stages (§4.10.2: THEORY-INTERFACE). Analysis of generalization and optimization largely reduces to their analysis in standard supervised learning.

⑦

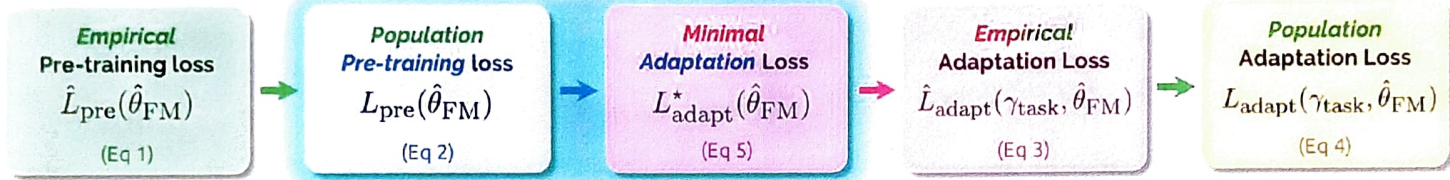


Generalization

Pre-Training -
Adaption interface

Optimization

Generalization



$$\hat{L}_{pre}(\theta) = \mathbb{E}_{x \sim \hat{p}_{pre}} [l_{pre}(x; \theta)] \quad - (1)$$

$$\hat{\theta}_{FM} = \underset{\theta \in \Theta}{\operatorname{argmin}} \hat{L}_{pre}(\theta)$$

$$L_{pre}(\theta) = \mathbb{E}_{x \sim p_{pre}} [l_{pre}(x; \theta)] \quad - (2)$$

$$\hat{L}_{adapt}(\gamma, \hat{\theta}_{FM}) = \mathbb{E}_{x \sim \hat{p}_{task}} [l_{adapt}(x; \gamma, \hat{\theta}_{FM})] \quad - (3)$$

$$\gamma_{task}(\hat{\theta}_{FM}) = \underset{\gamma \in \Gamma}{\operatorname{argmin}} \hat{L}_{adapt}(\gamma, \hat{\theta}_{FM})$$

$$L_{adapt}(\gamma, \hat{\theta}_{FM}) = \mathbb{E}_{x \sim p_{task}} [l_{adapt}(x; \gamma, \hat{\theta}_{FM})] \quad - (4)$$

$$\gamma_{task}^* = \underset{\gamma \in \Gamma, C(\gamma, \hat{\theta}_{FM}) \leq c_0}{\operatorname{argmin}} L_{adapt}(\gamma, \hat{\theta}_{FM})$$

$$L_{adapt}^*(\hat{\theta}_{FM}) = L_{adapt}(\gamma_{task}^*, \hat{\theta}_{FM}) \quad - (5)$$

PRETRAINING

ADAPTATION

Generalization

Pre-Training -
Adaption interface

Optimization

Generalization

Empirical
Pre-training loss

$$\hat{L}_{\text{pre}}(\hat{\theta}_{\text{FM}})$$

(Eq 1)

Population
Pre-training loss

$$L_{\text{pre}}(\hat{\theta}_{\text{FM}})$$

(Eq 2)

Minimal
Adaptation Loss

$$L_{\text{adapt}}^*(\hat{\theta}_{\text{FM}})$$

(Eq 5)

Empirical
Adaptation Loss

$$\hat{L}_{\text{adapt}}(\gamma_{\text{task}}, \hat{\theta}_{\text{FM}})$$

(Eq 3)

Population
Adaptation Loss

$$L_{\text{adapt}}(\gamma_{\text{task}}, \hat{\theta}_{\text{FM}})$$

(Eq 4)

\mathcal{Q} -space
 $\mathcal{Q} \in \mathcal{Q}$

γ -space
 $\gamma \in \Gamma$

$\gamma_{\text{task}}^* \approx \gamma^*$

Transfer & Adapt

ERM \times $\hat{\theta}_{\text{FM}}$

ERM \times $\gamma_{\text{task}}(\hat{\theta}_{\text{FM}})$

PRE-TRAINING



ADAPTATION

Generalization

Pre-Training -
Adaption interface

Optimization

Generalization

Empirical
Pre-training loss

$$\hat{L}_{\text{pre}}(\hat{\theta}_{\text{FM}})$$

(Eq 1)



Population
Pre-training loss

$$L_{\text{pre}}(\hat{\theta}_{\text{FM}})$$

(Eq 2)



Minimal
Adaptation Loss

$$L_{\text{adapt}}^*(\hat{\theta}_{\text{FM}})$$

(Eq 5)



Empirical
Adaptation Loss

$$\hat{L}_{\text{adapt}}(\gamma_{\text{task}}, \hat{\theta}_{\text{FM}})$$

(Eq 3)



Population
Adaptation Loss

$$L_{\text{adapt}}(\gamma_{\text{task}}, \hat{\theta}_{\text{FM}})$$

(Eq 4)

$$L_a(\gamma_{\text{task}}) \approx L_a^* + \text{Generalization Error.}$$

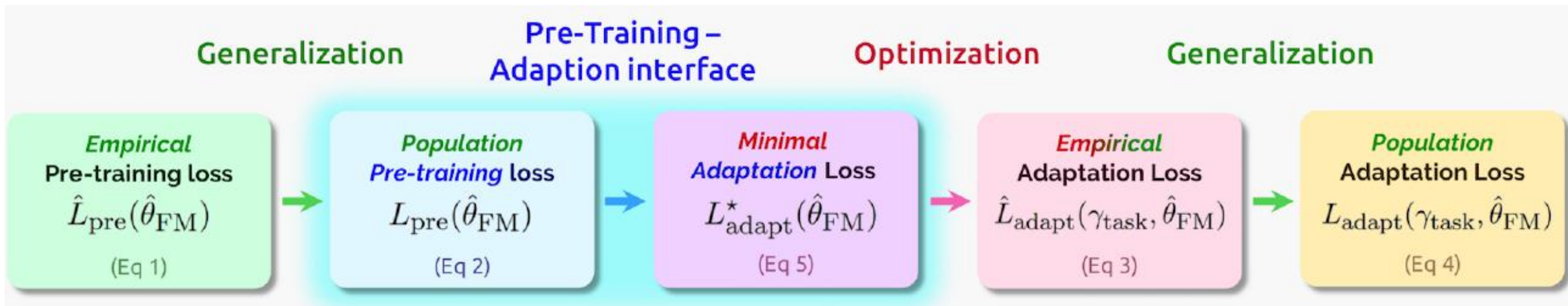
$$\downarrow$$

$$L_a(\gamma_{\text{task}}^*) \quad \text{Minimal Adaptation Loss}$$

To be replaced by

$$\rightarrow L_{\text{pre}}(\hat{\theta}_{\text{FM}})$$

Smaller $L_{\text{pre}}(\hat{\theta}_{\text{FM}}) \Rightarrow$ Smaller $L_a(\gamma_{\text{task}})$



As shown in Figure 22, the main missing link beyond standard supervised theory is:

Under what conditions does a small population pretraining loss $L_{\text{pre}}(\hat{\theta}_{\text{FM}})$ imply a small minimal adaptation loss $L_{\text{adapt}}^(\hat{\theta}_{\text{FM}})$ and why?*

Conditions / Factors affecting / influencing the statement / analysis (15)

1. Pretraining - Adaptation Interface

- Two different population quantities / distributions

Pretraining

Task.

- How do these two distributions relate to each other.
- Effect of distribution shifts \Rightarrow structural shifts

2. Model Architecture.

- Pretraining Distribution $\xrightarrow{\text{impact on}}$ Intermediate / Representations in DFN
(e.g. No, K_y split.)

3. Few Shot Learning in Downstream Supervised Task. (17)

Small "Population Pretraining Loss"

Low Complexity
Task (e.g. LP)

Sample efficiency
[Low, small size Task]

4. IMPORTANT: choice of $L_{pre} \approx L_{adapt}$.

(18)

Minimization of L_{pre}
using a particular L_{pre}
[Best suited for a given
PRETEXT task]

Need not be
optimal
for

Downstream
 L_{adapt} on a
particular L_{adapt}



Different
"Surrogate" objectives



Good adaptation
on a
wide range of
Downstream
tasks.

A Theoretical Analysis of Contrastive Unsupervised Representation Learning

Sanjeev Arora^{1 2} Hrishikesh Khandeparkar¹ Mikhail Khodak³ Orestis Plevrakis¹ Nikunj Saunshi¹

¹Princeton University, Princeton, New Jersey, USA. ²Institute for Advanced Study, Princeton, New Jersey, USA. ³Carnegie Mellon University, Pittsburgh, Pennsylvania, USA. Correspondence to: Orestis Plevrakis <orestisp@cs.princeton.edu>, Nikunj Saunshi <nsaunshi@cs.princeton.edu>.

Proceedings of the 36th International Conference on Machine Learning, Long Beach, California, PMLR 97, 2019. Copyright 2019 by the author(s).

Thank you !!