

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

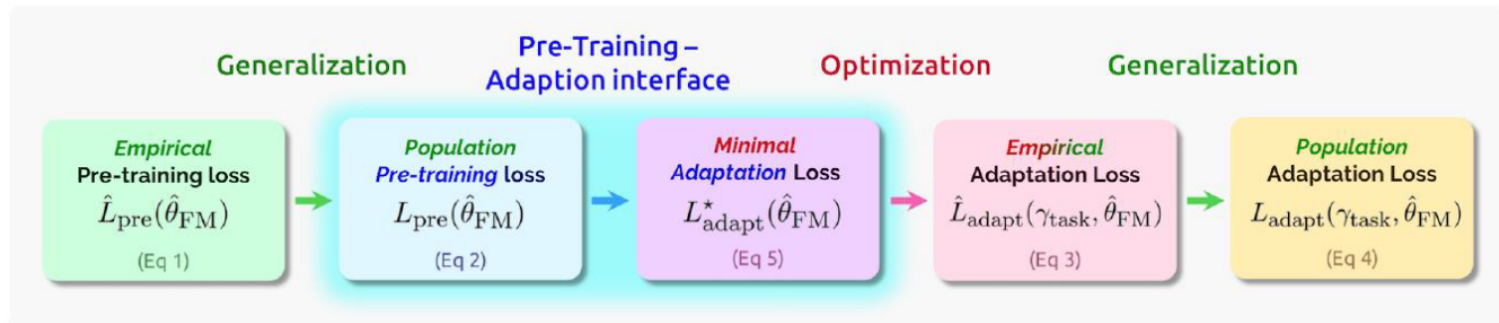


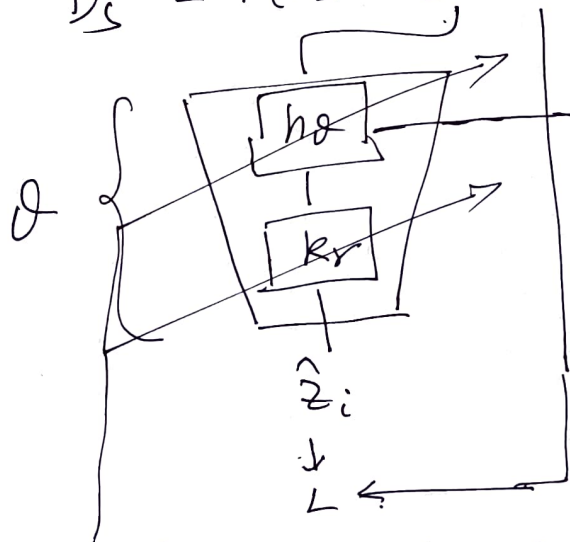
Fig. 22. The analysis of foundation models from pretraining on diverse data to downstream performance on adapted tasks involves capturing the relation between different loss terms as shown above. The main challenge is to analyze the highlighted pretraining-adaptation interface which requires reasoning 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.

PRE-TRAINING

$$D_S = \{x_i^{(k)}\}_{i=1}^M \quad M \gg N$$



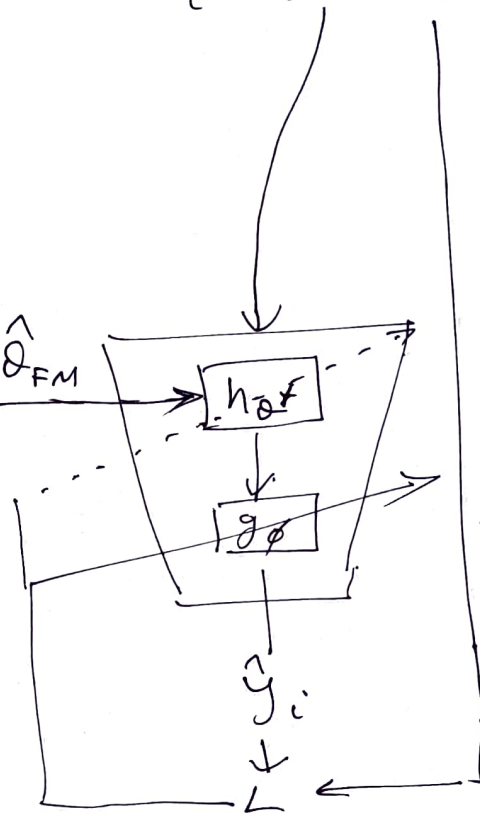
$$\bar{D}_S = P(D_S) = \{x_i, z_i\}_{i=1}^M$$



$$L(\hat{z}_i, z_i) \rightarrow \ell_{\text{pre}}(x; \theta)$$

ADAPTATION
[Downstream task]

$$D_T = \{x_i^{(t)}, y_i^{(t+1)}\}_{i=1}^N$$

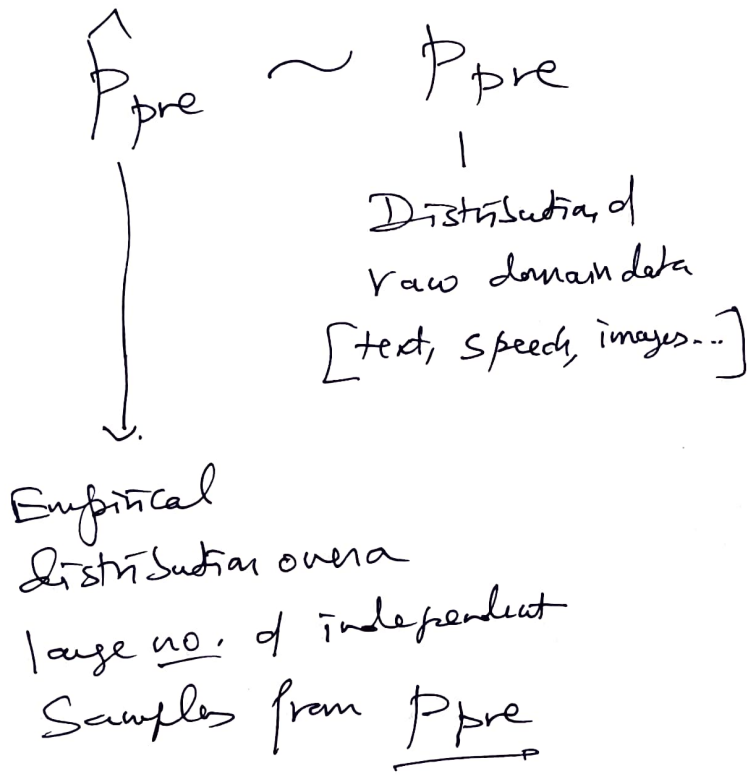


$$L(\hat{y}_i, y_i) \rightarrow \ell_{\text{adapt}}(x; \gamma, \hat{\theta}_{F_M})$$

$$\left. \begin{array}{l} \hat{\theta}_{F_M} \\ \gamma \\ \mathcal{L}_P \end{array} \right\} \gamma \in \Pi$$

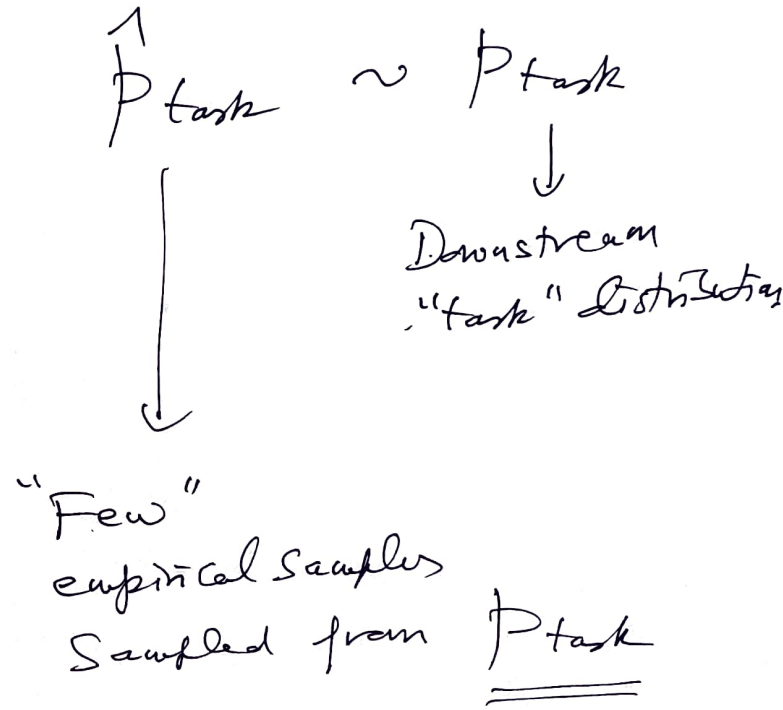
(2)

Pretraining Phase



Optimization-based Adaptation Phase

(3)



$$\ell_{\text{pre}}(x; \theta) \triangleq L(z_i, \hat{z}_i)$$

• ERM [Empirical Risk minimizer]

$$\hat{\theta}_{\text{pre}} = \underset{\theta \sim \hat{\mathcal{P}}_{\text{pre}}}{\text{argmin}} [\ell_{\text{pre}}(x; \theta)]$$

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

↓
Derived FOUNDATION MODEL [ERM]

ERM of Task →
↑
adapted from FM $\hat{\theta}_{\text{FM}}$

Π : Space of adapted model parameters

$$\subseteq \hat{\mathcal{Q}}_{\text{FM}}$$

adaptation
Since different methods
could modify different
subsets of the $\hat{\mathcal{Q}}_{\text{FM}}$

Recall: h₀, g₀ split / calibration.

• ERM [Downstream task.]

$$\hat{\ell}_{\text{adapt}}(x, \hat{\theta}_{\text{FM}})$$

$$= \mathbb{E}_{x \sim \hat{\mathcal{P}}_{\text{task}}} [\ell_{\text{adapt}}(x; \theta, \hat{\theta}_{\text{FM}})]$$

$$\gamma_{\text{task}}(\hat{\theta}_{\text{FM}}) = \underset{\gamma \in \Pi}{\text{argmin}} \hat{\ell}_{\text{adapt}}(\gamma, \hat{\theta}_{\text{FM}})$$

$$\text{Regularization } \|C(\gamma, \hat{\theta}_{\text{FM}}) \leq C_0$$

Expected Risk/Population loss
on the population distribution
 P_{pre}

$$L_{pre}(\theta) = \mathbb{E}_{x \sim P_{pre}} [l_{pre}(x; \theta)]$$

True, unknown distribution

on the population distribution
 P_{task} . (5)

$$L_{adapt}(\theta, \hat{Q}_{FM})$$

$$= \mathbb{E}_{x \sim P_{task}} [l_{adapt}(x; \theta, \hat{Q}_{FM})]$$

$$\theta_{task}^* = \arg \min_{\theta \in \Gamma} L_{adapt}(\theta, \hat{Q}_{FM})$$

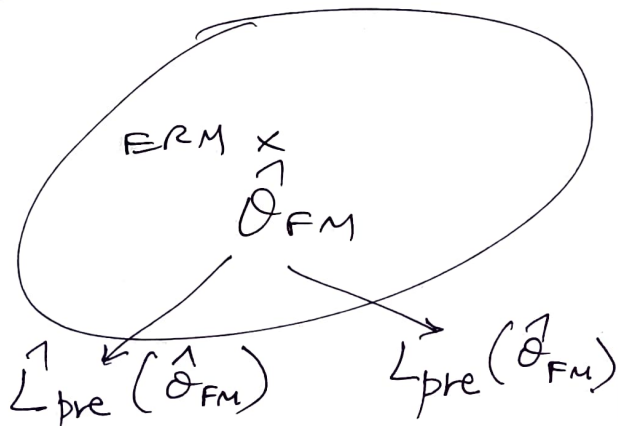
$$C(\theta, \hat{Q}_{FM}) \leq \epsilon_0$$

$$L_{adapt}^*(\hat{Q}_{FM}) = L_{adapt}(\theta_{task}^*, \hat{Q}_{FM})$$

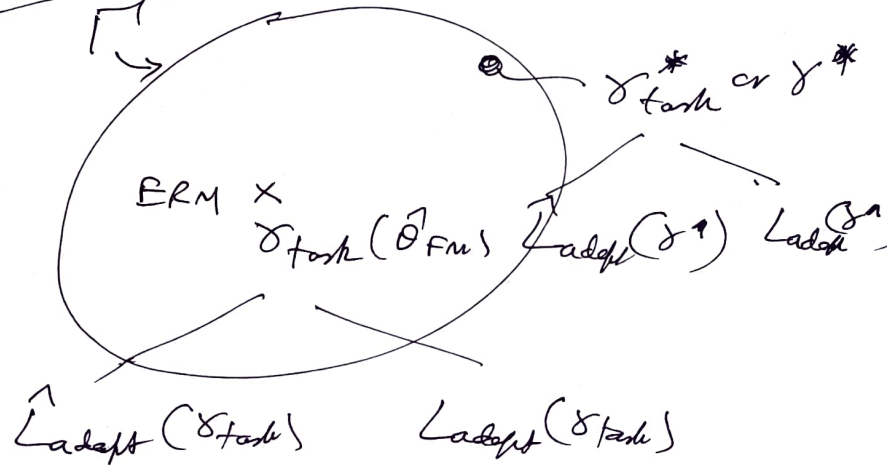
| | Empirical Risk (ERM) | Population Loss | $\hat{\theta}_{FM}$ (6) $LP \leftarrow \hat{\theta}_{FM} \rightarrow FT$ | |
|-------------|---|--|---|---|
| Pretraining | $\hat{L}_{pre} / \hat{\theta}_{FM}$ on \hat{p}_{pre} | L_{pre} on p_{pre} | $\hat{\theta}_{FM}$ ↓ $\hat{\theta}_{FM}$ ↓ $\gamma \in \mathbb{R}^k$ | $\hat{\theta}_{FM}$ ↓ $\hat{\theta}_{FM}$ + $\gamma \in \mathbb{R}^k$ |
| Adaptation | $\hat{L}_{adapt} / \gamma_{task}$ on \hat{p}_{task} | $L_{adapt} / \gamma_{task}^*$ on p_{task} | | |

Pretraining

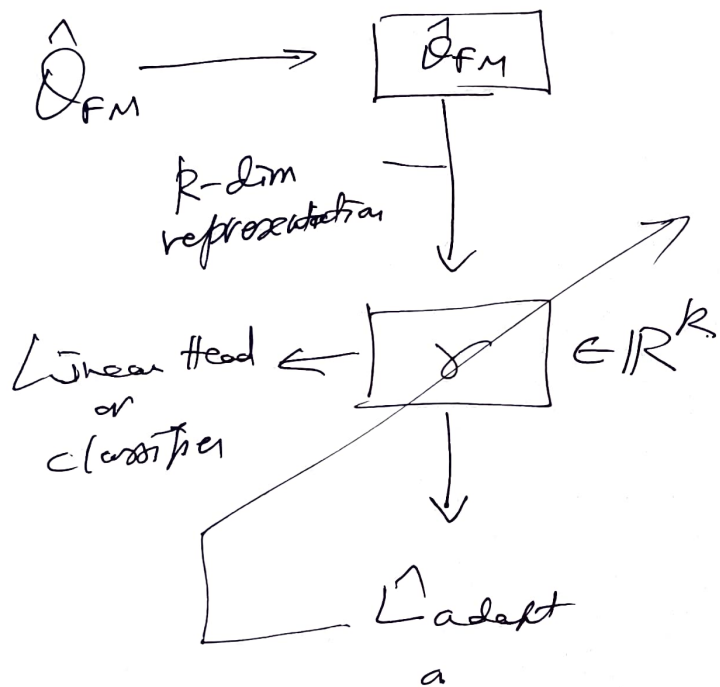
Q-space



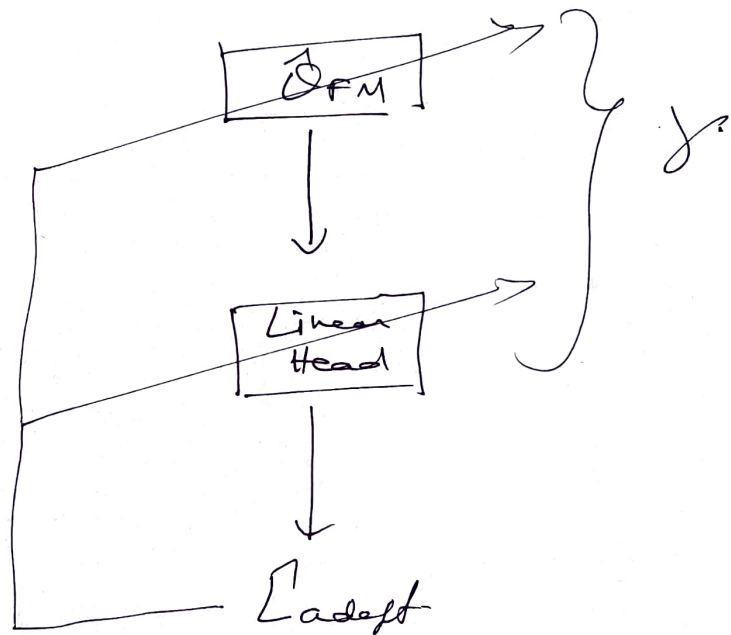
Adaptation



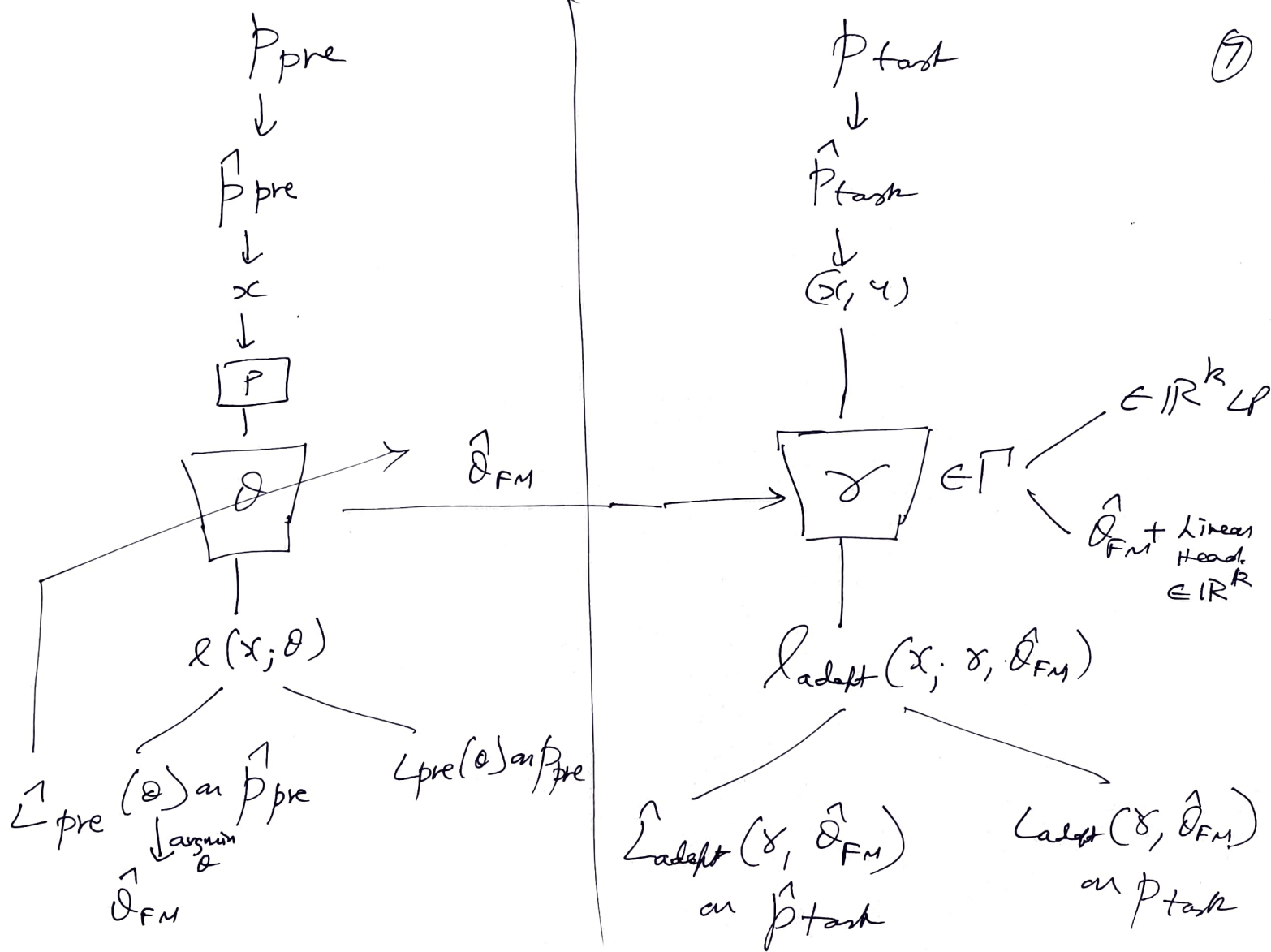
Linear Probing



Fine Tuning



⑦



Modularized Phases - Analysis

8

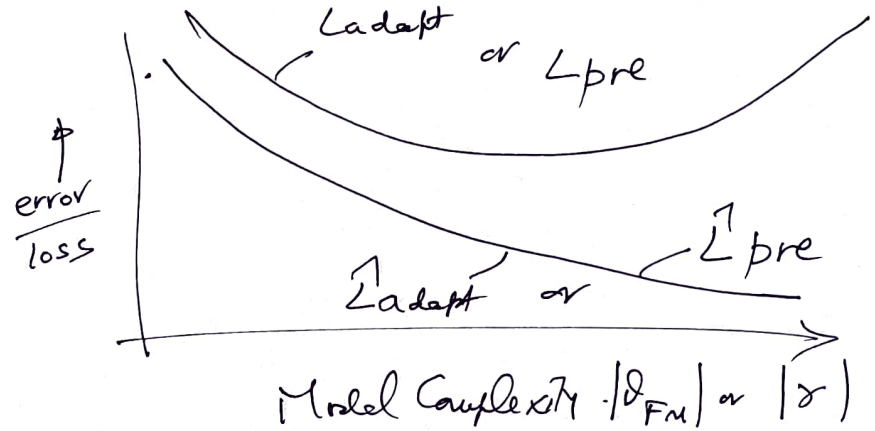
SLT Concerns how

$$\hat{L}_{pre} \approx L_{pre}$$

$$\hat{L}_{adapt} \approx L_{adapt}$$

i.e. Generalization Error

$$= L_{pre} - \hat{L}_{pre}$$



Standard Error Decomposition

⑨

Generalization Error

$$= \left[L_{\text{adapt}}(\gamma_{\text{task}}, \hat{\theta}_{\text{FM}}) - \hat{L}_{\text{adapt}}(\gamma_{\text{task}}, \hat{\theta}_{\text{FM}}) \right] - \textcircled{1}$$

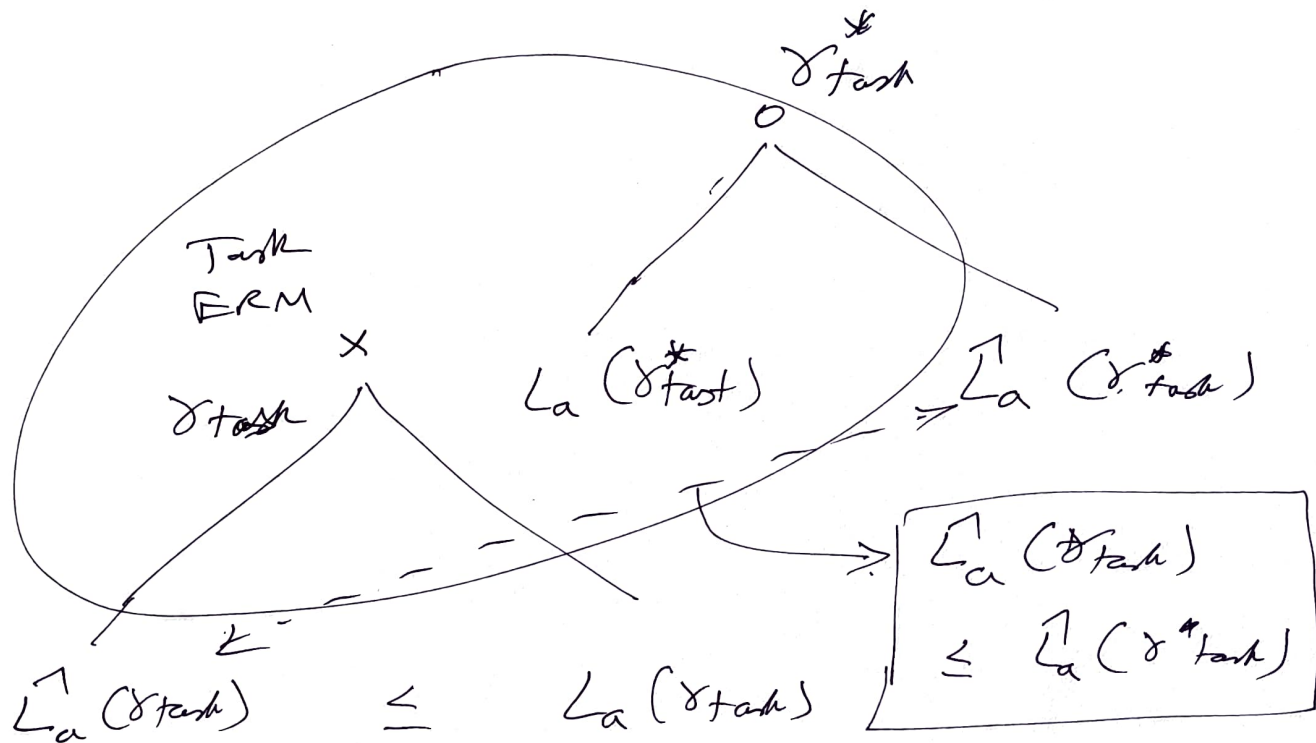
$$+ \left[\hat{L}_{\text{adapt}}(\gamma_{\text{task}}^*, \hat{\theta}_{\text{FM}}) - L_{\text{adapt}}(\gamma_{\text{task}}^*, \hat{\theta}_{\text{FM}}) \right]$$

Setting L_{adapt} as L_a & dropping $\hat{\theta}_{\text{FM}}$ in $\textcircled{1}$ for simplicity

$$GE = \left[L_a(\gamma_{\text{task}}) - \hat{L}_a(\gamma_{\text{task}}) \right] - \textcircled{2}$$

$$+ \left[\hat{L}_a(\gamma_{\text{task}}^*) - L_a(\gamma_{\text{task}}^*) \right]$$

$$GE = [L_a(\delta_{task}) - \hat{L}_a(\delta_{task})] + [\hat{L}_a(\delta_{task}^*) - L_a(\delta_{task}^*)] \quad (10)$$



$$CE = \left[L_a(\gamma_{task}) - \hat{L}_a(\gamma_{task}) \right] + \left[\hat{L}_a(\gamma_{task}^*) - L_a^* \right] \quad \text{⑪}$$

$\underline{L_a^* = L_a(\gamma_{task}^*)}$

Minimal Adaptation Loss

Using $\hat{L}_a(\gamma_{task}) \leq \hat{L}_a(\gamma_{task}^*)$

Best loss model
on Empirical Risk minimization
ie on \hat{p}_{task}

Best loss model from
Expected Risk minimization
ie a p_{task}

— γ_{task}^* optimized on \hat{p}_{task} loss is NOT optimal for loss surface
 $\hat{p}_{task} \sim p_{task}$

$$GE = \left[L_a(\gamma_{task}) - \hat{L}_a(\gamma_{task}) \right] + \left[\hat{L}_a(\gamma_{task}^*) - L_a^* \right]$$

(12)

$$\text{Thus } \hat{L}_a(\gamma_{task}) \leq \hat{L}_a(\gamma_{task}^*)$$

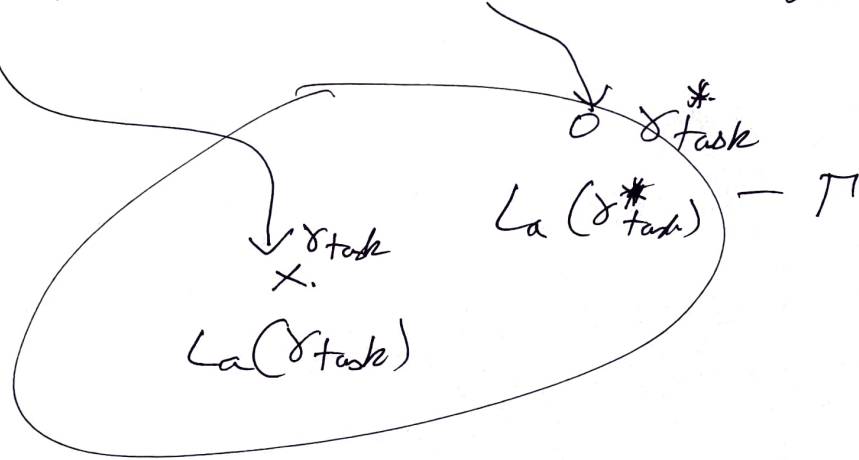
we have

(Excess)

$$L_a(\gamma_{task}, \hat{\theta}_{fM}) \leq L_a^*(\hat{\theta}_{fM}) + \text{Generalization Error}$$

(13)

$$L_a(\gamma_{\text{task}}, \hat{Q}_{FM}) \leq L_a^*(\hat{Q}_{FM}) + \text{Excess Generalization Error}$$



ie

$$\text{Population Loss}(\text{ERM}_{\text{task}} \text{ Model}) \leq \text{Population Loss}(\text{Best Model in } \mathcal{F}) + GE$$

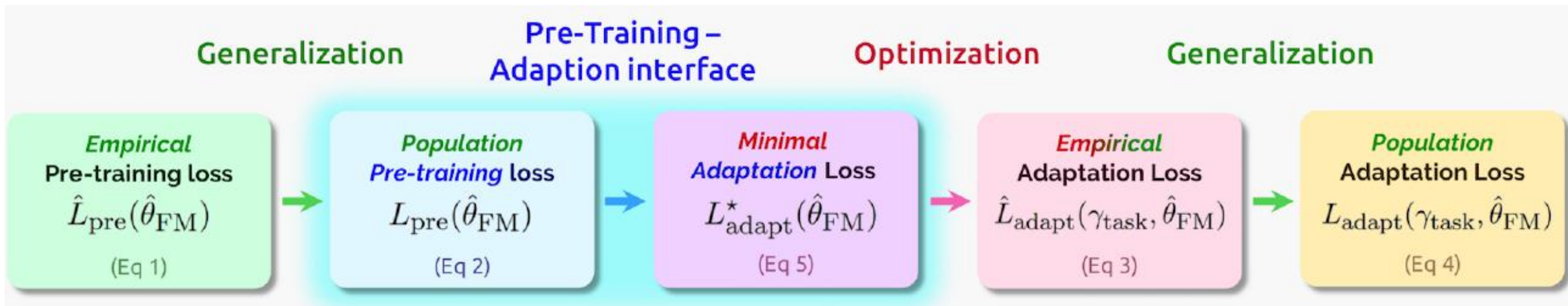
Kind of Main results shown

(14)

$$\text{Population Loss}(\text{ERM Task Model}) \leq \text{Population Loss}(\text{Best Model in } \mathcal{M}) + \text{GE}$$

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

|| i.e. we can Bound the Performance of downstream
adapted model ERM task [Popl. Loss a P-task]
by PRETRAINING Population Loss of FM ||



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?*

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).

Conditions / Factors affecting / influencing the statement / analysis (16)

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.

Thank you !!