

The Power and Limitation of Pretraining-Finetuning for Linear Regression under Covariate Shift

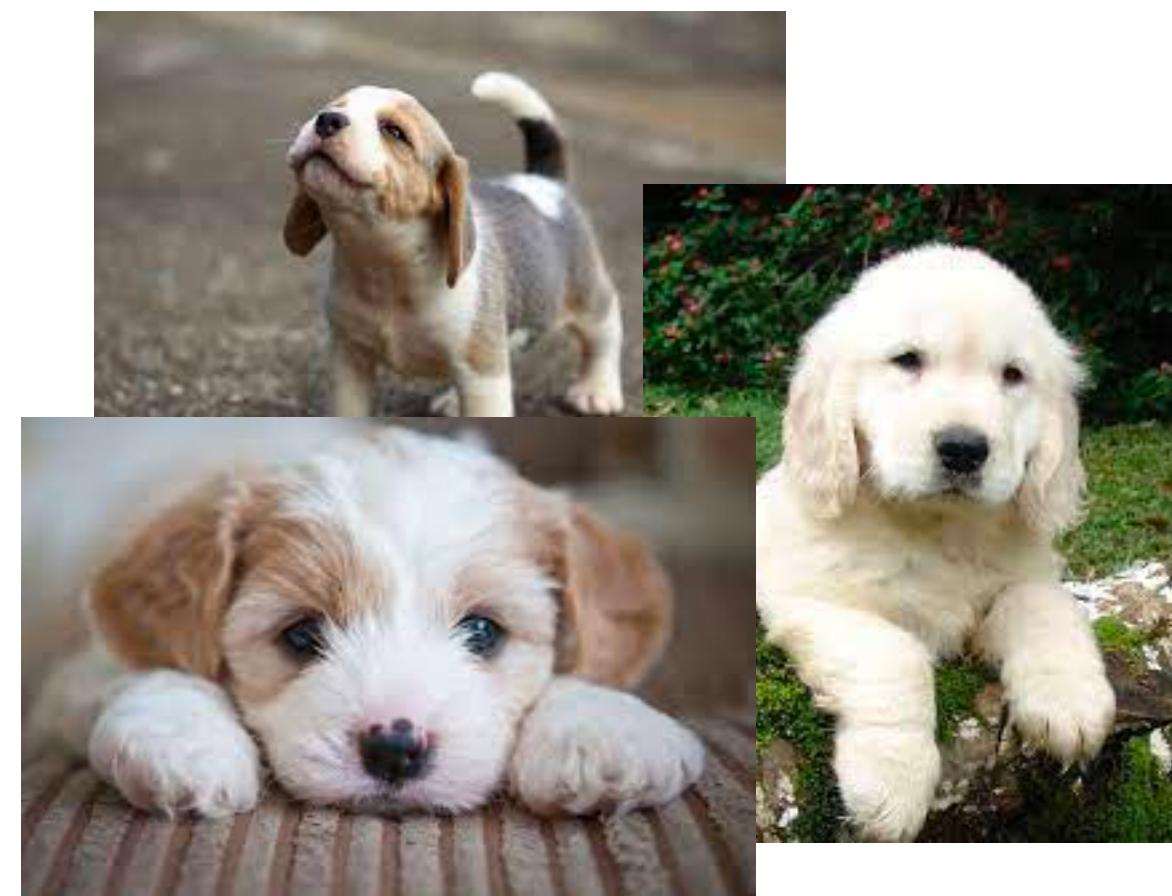
Jingfeng Wu

Joint work with Difan Zou, Vladimir Braverman, Quanquan Gu, Sham M. Kakade

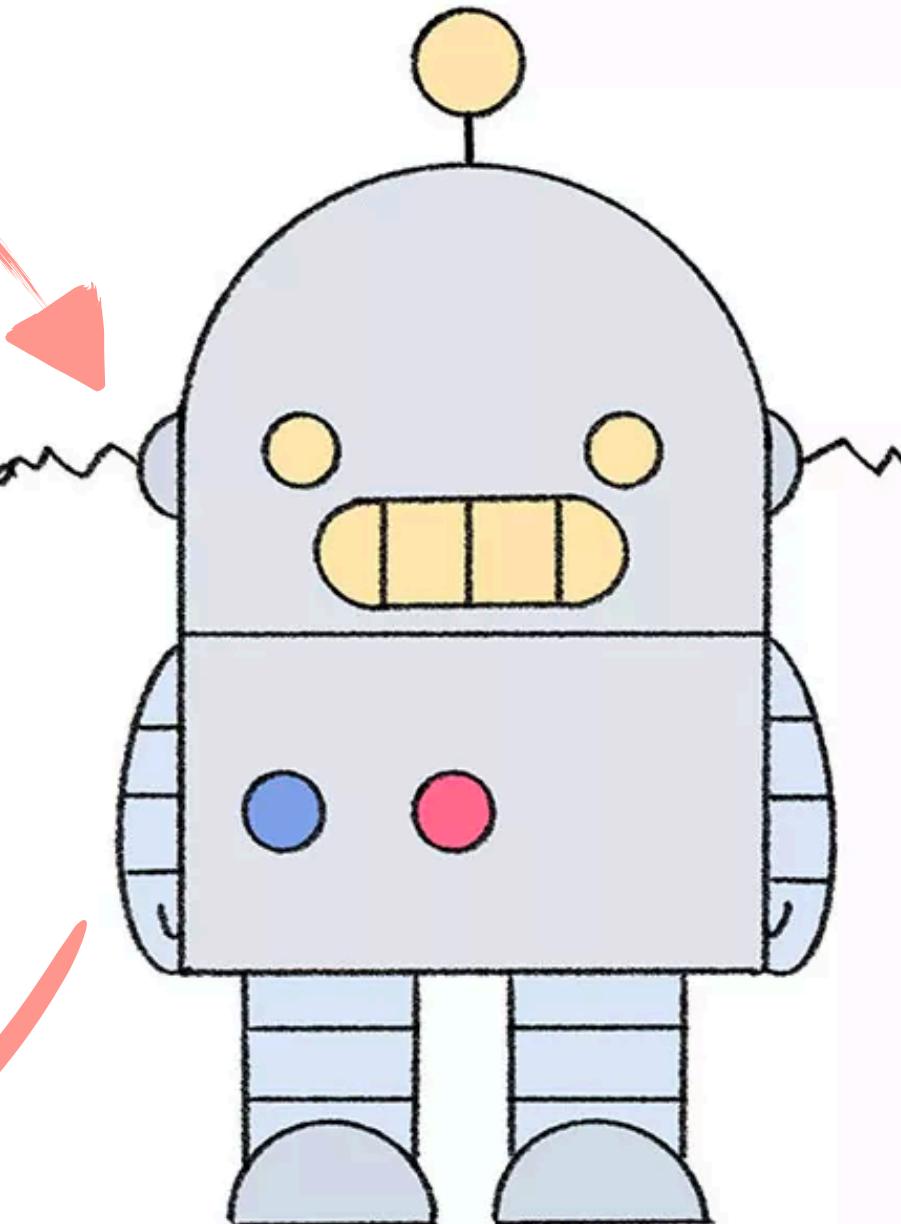
Covariate Shift

$$\mathbb{P}_{\text{source}}(x) \neq \mathbb{P}_{\text{target}}(x) \text{ but } \mathbb{P}_{\text{source}}(y|x) = \mathbb{P}_{\text{target}}(y|x)$$

Source domain



Cat or Dog?



Target domain, example 1



Adult Cat



Adult Dog

Target domain, example 2



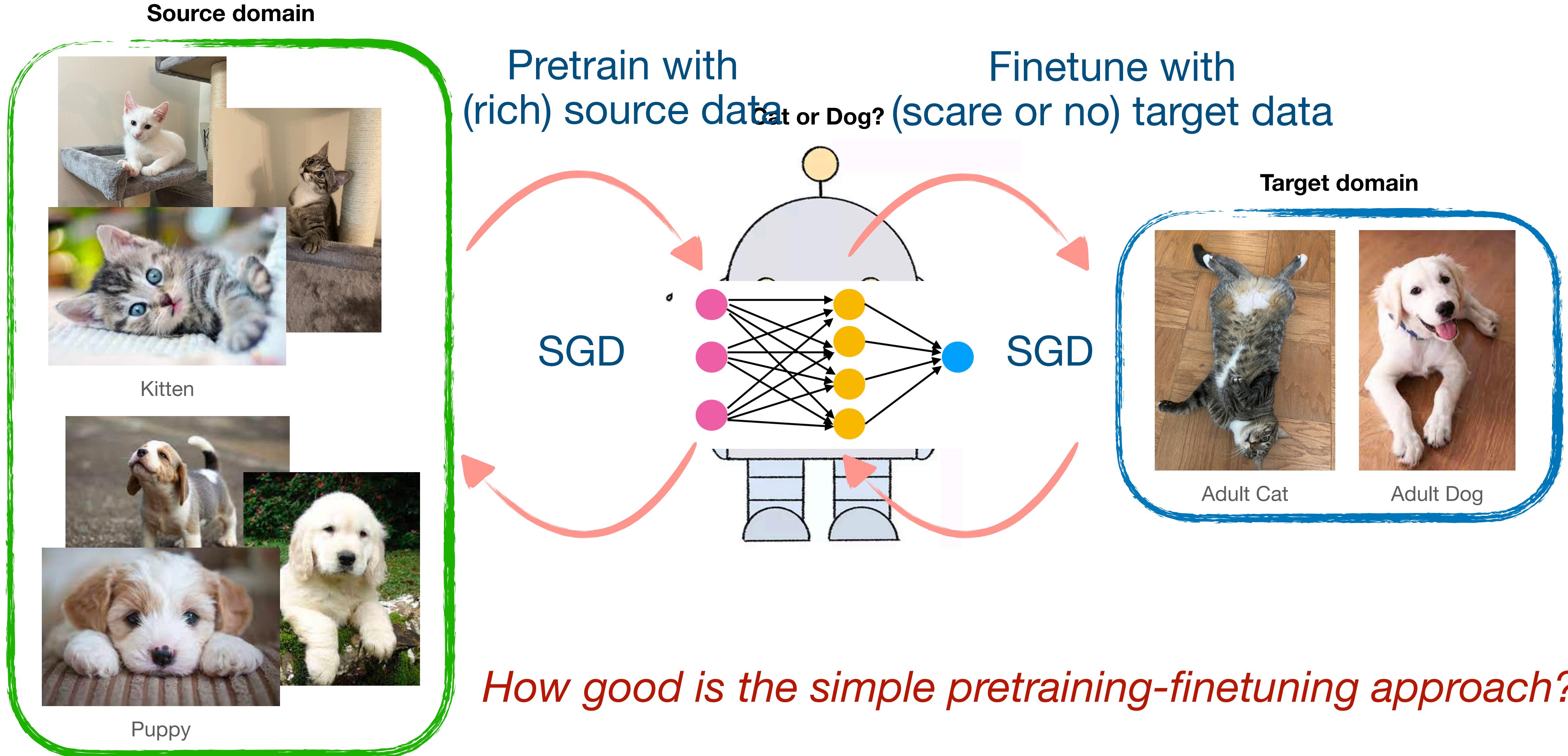
Tiger



Fox

*How to generalize
on the target domain?*

A Basic Yet Fundamental Approach



Problem Formulation

Linear Regression under Covariate Shift

- Shared Labeling Function

$$y = \mathbf{x}^\top \mathbf{w}^* + \mathcal{N}(0, \sigma^2)$$

- Source/Target Covariance Matrix

$$\mathbf{G} := \mathbb{E}_{\text{source}}[\mathbf{x}\mathbf{x}^\top] \quad \mathbf{H} := \mathbb{E}_{\text{target}}[\mathbf{x}\mathbf{x}^\top]$$

- Target Risk

$$\mathcal{L}(\mathbf{w}) := \mathbb{E}_{\text{target}}(y - \mathbf{x}^\top \mathbf{w})^2$$

- Target Excess Risk

$$\Delta(\mathbf{w}) := \mathcal{L}(\mathbf{w}) - \mathcal{L}(\mathbf{w}^*) = (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

Problem Formulation

Pretraining-Finetuning via Online SGD

Input

- M source data $(\mathbf{x}_t, y_t)_{t=1}^M \in \mathbb{R}^{d \times 1}$
- N target data $(\mathbf{x}_{M+t}, y_{M+t})_{t=1}^N \in \mathbb{R}^{d \times 1}$
- Initial stepsize η_0 for pretraining, initial stepsize η_M for finetuning

Output := \mathbf{w}_{M+N} , given by

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \eta_{t-t} \cdot (y_t - \mathbf{x}_t^\top \mathbf{w}_{t-1}) \cdot \mathbf{x}_t$$

$$\eta_t = \begin{cases} \eta_0/2^\ell, & 0 \leq t < M, \ell = \lfloor t/\log(M) \rfloor \\ \eta_M/2^\ell, & M \leq t < N, \ell = \lfloor (t-M)/\log(N) \rfloor \end{cases}$$

Main Result

A Sharp Risk Bound for Pretraining-Finetuning

$$\sigma^2 \cdot \left(\frac{D_{\text{eff}}^{\text{ft}}}{M} + \frac{D_{\text{eff}}}{N} \right) \lesssim \mathbb{E}[\Delta(\mathbf{w}_{M+N})] \lesssim (1 + \text{SNR}) \cdot \sigma^2 \cdot \left(\frac{D_{\text{eff}}^{\text{ft}}}{M} + \frac{D_{\text{eff}}}{N} \right)$$

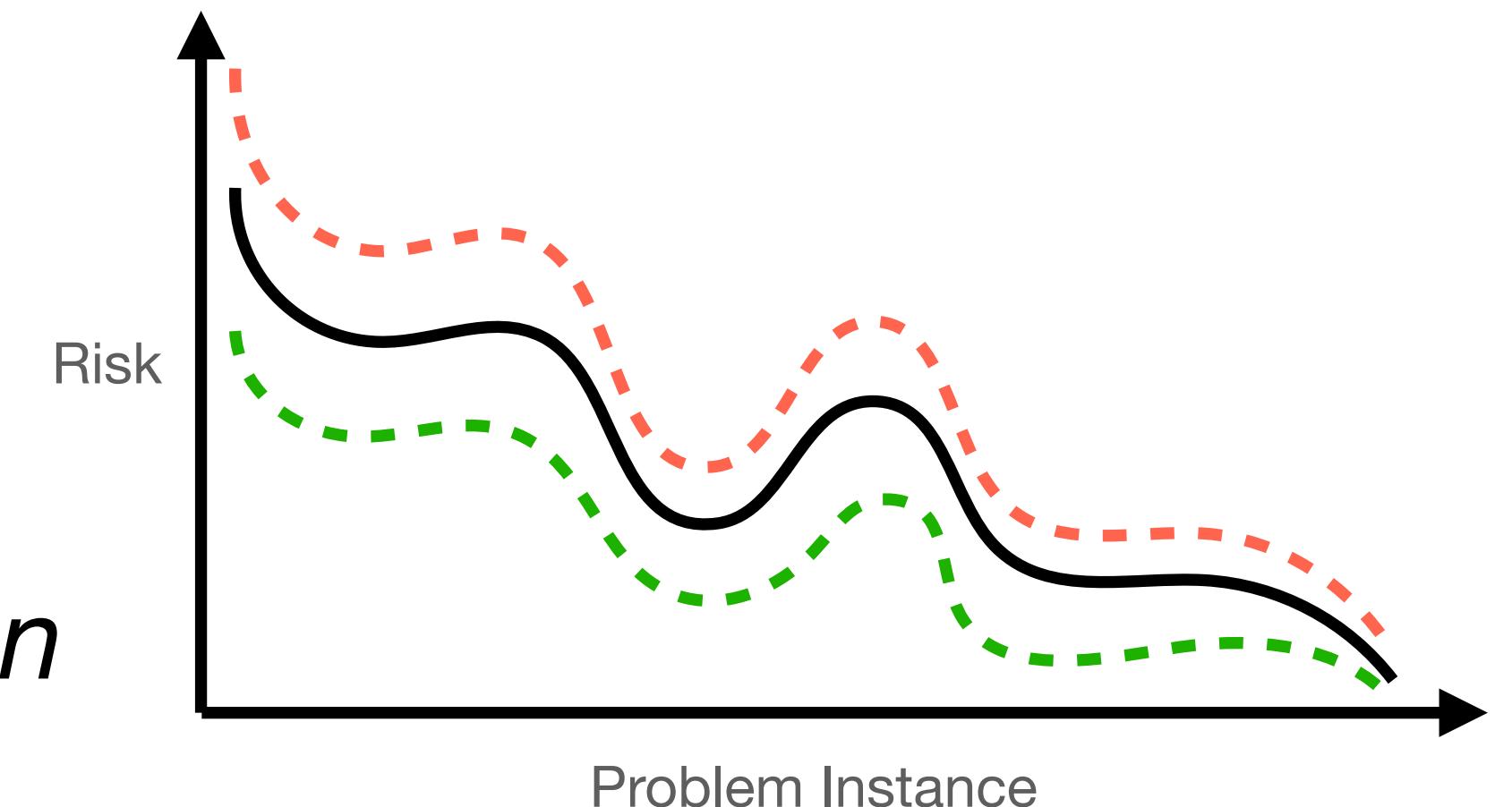
Effective Dimension

- $D_{\text{eff}}^{\text{ft}}$ is a function of $M, N, \eta_0, \eta_M, \mathbf{G}, \mathbf{H}$
- D_{eff} is a function of N, η_M, \mathbf{H}

Assumption

- Source & target data: *hyper-contractivity condition*
- $\text{SNR} := \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2 / \sigma^2 \lesssim 1$

Our bound is point-wisely sharp



Implication 1

Pretraining with M data vs. SL with N^{s_1} data

For **every** problem instance in \mathcal{C} ,

$$\mathcal{C} := \left\{ \mathbf{w}^*, \mathbf{H}, \mathbf{G}, \sigma^2 : \|\mathbf{w}^*\|_{\mathbf{G}}^2 \lesssim \sigma^2, \mathbf{GH} = \mathbf{HG} \right\}$$

we have

$$\mathbb{E}\Delta(\mathbf{w}_{M+0}) \lesssim \mathbb{E}\Delta(\mathbf{w}_{0+N^{s_1}}) \leftarrow M \gtrsim (N^{s_1})^2 \cdot \frac{\|\mathbf{H}_K\|_{\mathbf{G}}}{D_{\text{eff}}^{s_1}}$$

$O(n^2)$ source data $\gtrsim n$ target data

≈ 1 when \mathbf{G} weakly aligns with \mathbf{H}

```
graph LR; A["M \gtrsim (N^{s_1})^2"] --> B["(N^{s_1})^2"]; C["\|\mathbf{H}_K\|_{\mathbf{G}}"] --> D["\approx 1 when \mathbf{G} weakly aligns with \mathbf{H}"]
```

Implication 2

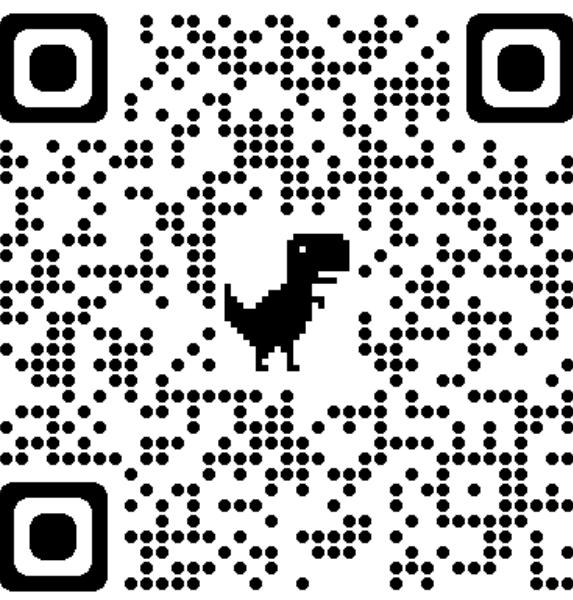
PT with M data vs. SL with N data vs. PT+FT with $M+N$ data

For each $\epsilon > 0$, there exists a problem instance in \mathcal{C} , such that:

to achieve ϵ -excess risk:

- Pretraining $\Rightarrow M \gtrsim \epsilon^{-2}$
- Supervised Learning $\Rightarrow N \gtrsim \epsilon^{-1.5}$
- Pretraining-Finetuning $\Leftarrow M \approx \epsilon^{-1} \log \epsilon^{-1}, N \approx \epsilon^{-1} \log^2 \epsilon^{-1}$

PT+FT could save polynomially than PT or SL alone



Take Home

1. Point-wisely sharp bounds for linear regression under covariate shift
2. $O(n^2)$ source data $\gtrsim n$ target data, when H weakly aligns with G
3. Pretraining-finetuning could save *polynomially* than pretraining (or SL)



Vova Braverman



Quanquan Gu



Sham Kakade



Difan Zou