

A Notation Glossary

| Notation | Definition |
|--|--|
| $[n]$ | $\{1, \dots, n\}$ |
| $\ \cdot\ _2$ | Euclidean vector norm |
| $\ \cdot\ _F$ | Frobenius matrix norm |
| $\langle \cdot, \cdot \rangle$ | Euclidean inner product OR Frobenius inner product |
| $\sigma_i(A)$ | i^{th} largest singular value of A |
| $\lambda_i(A)$ | i^{th} largest eigenvalue of A |
| $\lambda_{\max}(A)$ | Largest eigenvalue of A |
| $\lambda_{\min}(A)$ | Smallest eigenvalue of A |
| $A^{1/2}$ | Principal square root of PSD A |
| $\text{Col } A$ | Column span of A |
| P_A | Projection onto $\text{Col } A$ |
| P_A^\perp | Projection onto $(\text{Col } A)^\perp$ |
| ℓ | Loss function |
| $\mathcal{L}_\infty(g, g^*)$ | Population risk of g when the true predictor is g^* |
| $\mathcal{L}(g, g^*)$ | Empirical risk of g when the true predictor is g^* (samples suppressed) |
| $\mathcal{L}_\infty^{\text{ex}}(g, g^*)$ | Excess population risk of g when the true predictor is g^* . |
| $\mathcal{A}_\mathcal{C}(\theta)$ | Predictors whose parameters lie in $\theta + \mathcal{C}$ |
| $\mathcal{F}^{\otimes T}$ | $\{(x_1, \dots, x_T) \mapsto (f_1(x_1), \dots, f_T(x_T)) \mid f_1, \dots, f_T \in \mathcal{F}\}$ |
| $\mathcal{R}_n(\mathcal{H})$ | Rademacher complexity of function class \mathcal{H} on n samples |

B Proof of Theorem 3.1

In this section, we will prove the performance guarantee in the linear representation setting presented in Theorem 3.1. We first compute a bound on the difference in the spans of the true underlying representation B^* and the representation B_0 obtained from training on the source tasks. Having done so, we then analyze the performance of the best predictor found by projected gradient descent.

For clarity of presentation, we will write $\theta_t^* = B_t^* w_t^* + \delta_t^*$ and $\hat{\theta}_t = (B + \Delta_t) w_t$ throughout this section. Furthermore, let $\hat{\delta}_t = \Delta_t w_t$. Finally, we will be making use of the following covariance concentration results throughout this section, allowing us to connect empirical averages to population averages and vice versa:

Lemma B.1 (Source covariance concentration, Du et al. (2020), Claim A.1). *If $n_S \gg \rho^4[d + \log(T/\delta)]$, then with probability at least $1 - \delta/9$ over the random draw of $n_S T$ source inputs,*

$$0.9\Sigma \preceq \frac{1}{n_S} X_t^\top X_t \preceq 1.1\Sigma$$

for any $t \in [T]$.

Lemma B.2 (Target covariance concentration). *If $n_T \gg \rho^4[d + \log(T/\delta)]$, then with probability at least $1 - \delta/9$ over the random draw of n_T target inputs,*

$$0.9\Sigma \preceq \frac{1}{n_T} X^\top X \preceq 1.1\Sigma$$

for any $t \in [T]$.

Proof. The proof is similar to that of Lemma B.1, and is omitted for brevity. \square

B.1 Source Guarantees for the Linear Setting

We proceed to analyze the representation B_0 obtained from the source training procedure outlined in (4). Key to the analysis is a bound on the average population loss over the source tasks that the global minimizer of (4) can achieve as a function of n_S :

Lemma B.3 (Source training bound). *Let B_0 be a minimizer of (4), and $\{(\Delta_t, w_t)\}_{t \in [T]}$ be minimizers for the inner optimization problem given B_0 . If the regularizer coefficients are chosen such that*

$$\lambda \asymp \frac{1}{T\delta_0} \left[\frac{\sigma \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} + \frac{\sigma}{\sqrt{n_S}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{T}{\delta}\right)} \right],$$

and $\gamma = \delta_0^2 \lambda$, then with probability at least $1 - \delta/3$,

$$\begin{aligned} & \frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 \\ & \leq \frac{\sigma^2}{n_S T} \left(kT + kd \log \kappa n_S + \log \frac{1}{\delta} \right) + \frac{\sigma \delta_0 \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} \\ & \quad + \frac{\sigma \delta_0}{\sqrt{n_S}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{T}{\delta}\right)}. \end{aligned}$$

Proof. Throughout this proof, we instantiate the high-probability event in Lemma B.1, which occurs with probability at least $1 - \delta/9$.

Note that we can express δ_t^* as $[\delta_t^* (w_t^*)^\top / \|w_t^*\|_2^2] w_t^* = \Delta_t^* w_t^*$. Thus, via the optimality of B_0 and $\{(\Delta_t, w_t)\}_{t \in [T]}$, we can form the basic inequality

$$\begin{aligned} & \sum_{t \in [T]} \frac{1}{2n_S T} \left\| y_t - X_t \hat{\theta}_t \right\|_2^2 + \frac{\lambda}{2} \|\Delta_t\|_F^2 + \frac{\gamma}{2} \|w_t\|_2^2 \\ & \leq \sum_{t \in [T]} \frac{1}{2n_S T} \|z_t\|_2^2 + \frac{\lambda}{2} \|\Delta_t^*\|_F^2 + \frac{\gamma}{2} \|w_t^*\|_2^2 \\ & \leq \sum_{t \in [T]} \frac{1}{2n_S T} \|z_t\|_2^2 + \frac{\lambda}{2} \sum_{t \in [T]} \frac{\|\delta_t^*\|_2^2}{\|w_t^*\|_2^2} + \frac{\gamma}{2} \sum_{t \in [T]} \|w_t^*\|_2^2 \\ & \leq \sum_{t \in [T]} \frac{1}{2n_S T} \|z_t\|_2^2 + \left(\frac{\lambda \delta_0^2 + \gamma}{2} \right) T \\ & \leq \underbrace{\frac{\sigma \delta_0 \|\Sigma\|_2^{1/2}}{\sqrt{n T}} \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} + \frac{\sigma \delta_0}{\sqrt{n}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{T}{\delta}\right)}}_{=: C_{n_S, T}} \\ & \quad + \sum_{t \in [T]} \frac{1}{2n_S T} \|z_t\|_2^2. \end{aligned}$$

Note that the simplification of the regularizer on the optimum holds since $\|w_t^*\|_2 = \Theta(1)$ by Assumption 3.2. Equivalently, by rearranging,

$$\begin{aligned} & \sum_{t \in [T]} \frac{1}{2n_S T} \left\| X_t (\theta_t^* - \hat{\theta}_t) \right\|_2^2 + \frac{\lambda}{2} \|\Delta_t\|_F^2 + \frac{\gamma}{2} \|w_t\|_2^2 \\ & \leq \sum_{t \in [T]} \frac{1}{2n_S T} \left\langle z_t, X_t (\theta_t^* - \hat{\theta}_t) \right\rangle + C_{n_S, T}. \end{aligned}$$

Finally, by Proposition I.4, the regularizer on Δ_t and w_t can be rewritten as a regularizer on $\hat{\delta}_t$, i.e.

$$\sum_{t \in [T]} \frac{1}{2n_S T} \left\| X_t (\theta_t^* - \hat{\theta}_t) \right\|_2^2 + \sqrt{\lambda \gamma} \left\| \hat{\delta}_t \right\|_2 \leq \sum_{t \in [T]} \frac{1}{2n_S T} \left\langle z_t, X_t (\theta_t^* - \hat{\theta}_t) \right\rangle + C_{n_S, T}. \quad (8)$$

To proceed, define the set S as

$$\left\{ [\alpha_t + \beta_t]_{t \in [T]} \mid \text{rank}[\alpha_t] \leq 2k, \|\beta_t\|_2 \leq \delta_0 + \|\hat{\delta}_t\|_2, \|X_t(\alpha_t + \beta_t)\|_2 \leq \|X_t(\theta_t^* - \hat{\theta}_t)\|_2 \right\},$$

and observe that $[\theta_t^* - \hat{\theta}_t]_{t \in [T]} \in S$, by letting $\alpha_t = B^* w_t^* - B_0 w_t$ and $\beta_t = \delta_t^* - \hat{\delta}_t$. We bound the right-hand side of (8) via bounding the supremum of the inner product over S , i.e.

$$\sum_{t \in [T]} \left\langle z_t, X_t(\theta_t^* - \hat{\theta}_t) \right\rangle \leq \sup_{[\alpha_t + \beta_t] \in S} \sum_{t \in [T]} \langle z_t, X_t(\alpha_t + \beta_t) \rangle.$$

Now, we decompose the Gaussian width as

$$\sup_{[\alpha_t + \beta_t] \in S} \sum_{t \in [T]} \langle z_t, X_t(\alpha_t + \beta_t) \rangle \leq \underbrace{\sup_{[\alpha_t] \in S} \sum_{t \in [T]} \langle z_t, X_t \alpha_t \rangle}_{\text{(I)}} + \underbrace{\sup_{[\beta_t] \in S_2} \sum_{t \in [T]} \langle z_t, X_t \beta_t \rangle}_{\text{(II)}},$$

where $S_2 = \left\{ [\beta_t]_{t \in [T]} \mid \|\beta_t\|_2 \leq \delta_0 + \|\hat{\delta}_t\|_2 \right\}$. Note the abuse of notation in (I), where we say $[\alpha_t]_{t \in [T]} \in S$ if there exists $[\beta_t]_{t \in [T]}$ so that $[\alpha_t + \beta_t]_{t \in [T]} \in S$. This decomposes the Gaussian width into the sum of the Gaussian widths of a low-rank set (I) and a small norm set (II). We proceed to bound both terms accordingly to these two properties.

Bounding the Gaussian width of the low-rank set (I).

To bound the Gaussian width, we first enlarge S to remove β_t from the definition of the feasible set. Fix any (α_t, β_t) pair satisfying the conditions in S , and note that

$$\|X_t \beta_t\|_2 \lesssim \sqrt{n_S} \|\Sigma\|_2^{1/2} \left(\delta_0 + \|\hat{\delta}_t\|_2 \right).$$

Therefore, by the reverse triangle inequality,

$$\begin{aligned} \left| \|X_t \alpha_t\|_2 - \|X_t \beta_t\|_2 \right| &\leq \|X_t(\alpha_t + \beta_t)\|_2 \\ \implies \|X_t \alpha_t\|_2 &\leq \underbrace{\|X_t(\theta_t^* - \hat{\theta}_t)\|_2 + \sqrt{n_S} \|\Sigma\|_2^{1/2} \left(\delta_0 + \|\hat{\delta}_t\|_2 \right)}_{=:\rho_t}. \end{aligned}$$

Consequently, we can enlarge the feasible set to

$$S_1 := \{[\alpha_t] \mid \text{rank}[\alpha_t] \leq 2k, \|X_t \alpha_t\|_2 \leq \rho_t\}.$$

Having relaxed the constraints, we now proceed to the main argument.

Since $\text{rank}[\alpha_t] \leq 2k$, there exists an orthogonal matrix $V \in \mathbb{R}^{d \times 2k}$ (dependent on α_t) and vectors $r_t \in \mathbb{R}^{2k}$ such that $\alpha_t = V r_t$. Therefore, the inner product in the Gaussian width would be unchanged if we project z_t onto $X_t V$, i.e.

$$\sup_{[\alpha_t] \in S_1} \sum_{t \in [T]} \langle z_t, X_t \alpha_t \rangle = \sup_{[\alpha_t] \in S_1} \sum_{t \in [T]} \langle P_{X_t V} z_t, X_t V r_t \rangle.$$

The key idea we will leverage is that if V were chosen independently of z_t , then $P_{X_t V} z_t$ is Gaussian in a $2k$ -dimensional space, and thus norm bounded by $\tilde{O}(\sqrt{2k})$ with high probability. However, due to the supremum over α_t , this independence assumption is not satisfied. Nevertheless, we can obtain a *fixed* finite covering of the set of all rank- $2k$ matrices, and ensure that the aforementioned norm bound on $P_{X_t V} z_t$ holds for every V in the covering via a union-bound. By choosing the discretization level of the covering appropriately, we can control the error resulting from approximating the supremum by some element of the covering.

Formally, let $\mathcal{O}^{d \times 2k}$ be the set of orthogonal matrices in $\mathbb{R}^{d \times 2k}$. By Proposition I.1, there exists an ε -covering of $\mathcal{O}^{d \times 2k}$ in the Frobenius norm with at most $(6\sqrt{2k}/\varepsilon)^{2kd}$ elements. Let \bar{V} be an

element of the covering such that $\|V - \bar{V}\|_F \leq \varepsilon$. Then,

$$\begin{aligned}
\langle z_t, X_t \alpha_t \rangle &= \langle z_t, X_t \bar{V} r_t \rangle + \langle z_t, X_t (V - \bar{V}) r_t \rangle \\
&\leq \|P_{X_t \bar{V}} z_t\|_2 \|X_t \bar{V} r_t\|_2 + \|z_t\|_2 \|X_t (V - \bar{V}) r_t\|_2 \\
&\leq \|P_{X_t \bar{V}} z_t\|_2 \|X_t V r_t\|_2 + (\|P_{X_t \bar{V}} z_t\|_2 + \|z_t\|_2) \|X_t (V - \bar{V}) r_t\|_2 \\
&\leq \|P_{X_t \bar{V}} z_t\|_2 \|X_t V r_t\|_2 + (\|P_{X_t \bar{V}} z_t\|_2 + \|z_t\|_2) \|X_t (V - \bar{V}) r_t\|_2 \\
&\lesssim \|P_{X_t \bar{V}} z_t\|_2 \|X_t \alpha_t\|_2 + \|z_t\|_2 \|X_t (V - \bar{V}) r_t\|_2,
\end{aligned}$$

and thus

$$\begin{aligned}
\sum_{t \in [T]} \langle z_t, X_t V r_t \rangle &\leq \underbrace{\sqrt{\left(\sum_{t \in [T]} \|P_{X_t \bar{V}} z_t\|_2^2 \right)}}_{(A)} \underbrace{\left(\sum_{t \in [T]} \|X_t \alpha_t\|_2^2 \right)}_{(C)} \\
&\quad + \underbrace{\sqrt{\left(\sum_{t \in [T]} \|z_t\|_2^2 \right)}}_{(B)} \underbrace{\left(\sum_{t \in [T]} \|X_t (V - \bar{V}) r_t\|_2^2 \right)}_{(C)}.
\end{aligned}$$

We will bound (A), (B), and (C) individually.

(A) For a fixed \bar{V} , (A) is a chi-squared random variable with kT degrees of freedom scaled by σ^2 . However, since \bar{V} depends on V , we need to have a high probability bound for any element of the covering. By using known concentration bounds for chi-squared random variables together with the union-bound, we find that uniformly over the covering,

$$(A) \leq \sigma^2 \left(kT + \log \frac{1}{\delta'} \right) \quad \text{w.p.} \geq 1 - \delta' \left(\frac{6\sqrt{2k}}{\varepsilon} \right)^{2kd}.$$

(B) Note that (B) is a chi-squared random variable with $n_S T$ degrees of freedom scaled by σ^2 , and thus

$$(B) \lesssim \sigma^2 \left(n_S T + \log \frac{1}{\delta} \right) \quad \text{w.p.} \geq 1 - \frac{\delta}{18}.$$

(C) Since $(1/n_S) X^\top X$ is concentrated about Σ via Lemma B.1,

$$\|X_t (V - \bar{V}) r_t\|_2^2 \lesssim n_S \|\Sigma\|_2 \varepsilon^2 \|r_t\|_2^2 = n_S \|\Sigma\|_2 \varepsilon^2 \|V r_t\|_2^2 \lesssim \kappa \varepsilon^2 \|X_t V r_t\|_2^2.$$

Putting these bounds together, and setting $\varepsilon = \sqrt{k/\kappa n_S}$ and $\delta' = \delta/[18(6\sqrt{2k}/\varepsilon)^{2kd}]$, we obtain

$$\begin{aligned}
\sum_{t \in [T]} \langle z_t, X_t \alpha_t \rangle &\lesssim \left(\sigma \sqrt{kT + \log \frac{1}{\delta'}} + \varepsilon \sigma \sqrt{\kappa n_S T + \log \frac{1}{\delta}} \right) \sqrt{\left(\sum_{t \in [T]} \|X_t V r_t\|_2^2 \right)} \\
&\lesssim \left(\sigma \sqrt{kT + \log \frac{1}{\delta'}} \right) \sqrt{\left(\sum_{t \in [T]} \|X_t V r_t\|_2^2 \right)} \quad (k < n_S, \delta' < \delta) \\
&\lesssim \sigma \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} \sqrt{\left(\sum_{t \in [T]} \|X_t \alpha_t\|_2^2 \right)}.
\end{aligned}$$

Taking the supremum over S_1 , we thus obtain the following bound on the Gaussian width:

$$\begin{aligned} & \sup_{[\alpha_t] \in S_1} \sum_{t \in [T]} \langle z_t, X_t \alpha_t \rangle \\ & \lesssim \sigma \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} \left[\sqrt{\sum_{t \in [T]} \|X_t(\theta_t^* - \hat{\theta}_t)\|_2^2} + \sqrt{\|\Sigma\|_2 n_S T} (\delta_0 + \|\hat{\delta}_t\|_2) \right]. \end{aligned}$$

Note that the events for this sub-argument occur with probability at least $1 - \delta/9$.

Bounding the Gaussian width of the low-norm set (II).

Recall that we want to bound

$$\sup_{[\beta_t] \in S_2} \sum_{t \in [T]} \langle z_t, X_t \beta_t \rangle, \quad S_2 := \left\{ [\beta_t] \mid \|\beta_t\|_2 \leq \delta_0 + \|\hat{\delta}_t\|_2 \right\}.$$

For any $t \in [T]$,

$$\langle z_t, X_t \beta_t \rangle = \langle X_t^\top z_t, \beta_t \rangle \leq (\delta_0 + \|\hat{\delta}_t\|_2) \|X_t^\top z_t\|_2.$$

Furthermore, by the Hanson-Wright inequality, we have that with probability at least $1 - \delta/9T$,

$$\|X_t^\top z_t\|_2 \lesssim \sigma \sqrt{n_S \text{tr} \Sigma \left(1 + \log \frac{T}{\delta}\right)}.$$

Putting everything together, we thus find that with probability at least $1 - \delta/9$,

$$\sup_{[\beta_t] \in S_2} \sum_{t \in [T]} \langle z_t, X_t \beta_t \rangle \leq \sigma \left(T \delta_0 + \sum_{t \in [T]} \|\hat{\delta}_t\|_2 \right) \sqrt{n_S \text{tr} \Sigma \left(1 + \log \frac{T}{\delta}\right)}.$$

Combining the bounds and concluding.

Having bounded both Gaussian widths, we can thus bound the right-hand side of (8) as

$$\begin{aligned} & \sum_{t \in [T]} \frac{1}{2n_S T} \|X_t(\theta_t^* - \hat{\theta}_t)\|_2^2 + \left(\sqrt{\lambda \gamma} - \frac{1}{\delta_0 T} C_{n_S, T} \right) \|\hat{\delta}_t\|_2 \\ & \lesssim \frac{\sigma}{\sqrt{n_S T}} \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} \left(\sum_{t \in [T]} \frac{1}{2n_S T} \|X_t(\theta_t^* - \hat{\theta}_t)\|_2^2 \right)^{1/2} + C_{n_S, T}. \end{aligned}$$

Therefore, as long as $\sqrt{\lambda \gamma} \geq (1/\delta_0 T) C_{n_S, T}$,

$$\begin{aligned} & \sum_{t \in [T]} \frac{1}{2n_S T} \|X_t(\theta_t^* - \hat{\theta}_t)\|_2^2 \\ & \lesssim \frac{\sigma}{\sqrt{n_S T}} \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} \left(\sum_{t \in [T]} \frac{1}{2n_S T} \|X_t(\theta_t^* - \hat{\theta}_t)\|_2^2 \right)^{1/2} + C_{n_S, T}. \end{aligned}$$

Finally, by solving the quadratic inequality using Proposition I.2, we find that

$$\begin{aligned} & \sum_{t \in [T]} \frac{1}{2n_S T} \|X_t(\theta_t^* - \hat{\theta}_t)\|_2^2 \\ & \lesssim \frac{\sigma^2}{n_S T} \left(kT + kd \log \kappa n_S + \log \frac{1}{\delta} \right) + \frac{\sigma \delta_0 \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} \\ & \quad + \frac{\sigma \delta_0}{\sqrt{n_S}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{T}{\delta}\right)}, \end{aligned}$$

from which the desired performance bound follows due to the concentration of the empirical covariance from Lemma B.1. Since the concentration of the source covariance matrices and each of the sub-arguments all hold with probability at least $1 - \delta/9$, it follows that the main claim holds with probability at least $1 - \delta/3$. \square

The prior bound is central to the analysis, as the performance of the learner can be tied to how well B_0 spans the correct space. To see why this is the case, note that for large n_S , the effect of the noise on the optimization in (4) is negligible. In this regime, no matter which representation B_0 the learner has chosen, the optimal predictor would satisfy $P_{X_t B_0} X_t \hat{\theta}_t \approx P_{X_t B_0} X_t \theta_t^*$ and $P_{X_t B_0}^\perp X_t \hat{\theta}_t \approx P_{X_t B_0}^\perp X_t \theta_t^*$. Consequently, the performance of the predictors chosen by the learner can be tied to the chosen representation B_0 . We formalize this intuition in the following result:

Lemma B.4 (Transfer Lemma). *Under Assumption 3.2, we have that*

$$\left\| P_{\Sigma^{1/2} B_0} \Sigma^{1/2} B^* \right\|_F^2 \lesssim k \left[\frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 \right].$$

Proof. Throughout this proof, we will write $P := P_{\Sigma^{1/2} B_0}$ and $P^\perp := P_{\Sigma^{1/2} B_0}^\perp$ for readability. To proceed, note that we intuitively expect that for a learner that has learned the correct spaces, $P \Sigma^{1/2} \hat{\theta}_t \approx \Sigma^{1/2} B^* w_t^*$ as it is the low-rank component of the estimator, and consequently, $P^\perp \Sigma^{1/2} \hat{\theta}_t \approx \Sigma^{1/2} \delta_t^*$. Then, decomposing into the corresponding errors, we have that for any $t \in [T]$,

$$\begin{aligned} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 &\gtrsim \left\| \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \hat{\theta}_t \right\|_2^2 + \left\| \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\|_2^2 \\ &\quad + 2 \left\langle \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \hat{\theta}_t, \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\rangle \\ &\geq \left\| P^\perp \Sigma^{1/2} B^* w_t^* \right\|_2^2 + \left\| \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\|_2^2 \\ &\quad - 2 \left| \left\langle \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \hat{\theta}_t, \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\rangle \right|. \end{aligned}$$

We proceed to bound the inner product above. We do so by observing that if we were to replace $\hat{\theta}_t$ by θ_t^* , then

$$\begin{aligned} &\left\langle \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \theta_t^*, \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \delta_t^* \right\rangle \\ &= \left\langle P^\perp \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \delta_t^*, P^\perp \Sigma^{1/2} \delta_t^* - P \Sigma^{1/2} B^* w_t^* \right\rangle \\ &= (w_t^*)^\top (B^*)^\top \Sigma \delta_t^* \\ &= 0. \end{aligned}$$

To translate this result into a bound on the original inner product, we note that as $n_S \rightarrow \infty$, and the impact of the noise and the regularizer on the optimization is diminished, we expect $\hat{\theta}_t$ to learn the projections of $X_t \theta_t^*$ onto $X_t B_0$ and its complement. With these two insights in mind, we decompose the inner product as

$$\begin{aligned} &\left| \left\langle \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \hat{\theta}_t, \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\rangle \right| \\ &\leq \left| \left\langle \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \theta_t^*, P^\perp \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\rangle \right| \\ &\quad + \left| \left\langle P \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t), \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\rangle \right| \\ &\quad + \underbrace{\left| \left\langle \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \theta_t^*, \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \theta_t^* \right\rangle \right|}_{=0} \\ &= \left| \left\langle P^\perp \Sigma^{1/2} B^* w_t^*, P^\perp \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\rangle \right| + \left| \left\langle P \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t), \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\rangle \right|. \end{aligned}$$

Putting everything together, we find that

$$\begin{aligned}
& \frac{1}{T} \sum_{t \in [T]} \left\| P^\perp \Sigma^{1/2} B^* w_t^* \right\|_2^2 + \frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\|_2^2 \\
& \lesssim \frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 + \frac{1}{T} \sum_{t \in [T]} \left| \left\langle \Sigma^{1/2} B^* w_t^* - P \Sigma^{1/2} \hat{\theta}_t, \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\rangle \right| \\
& \leq \frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 + \frac{1}{T} \sum_{t \in [T]} \left| \left\langle P^\perp \Sigma^{1/2} B^* w_t^*, P^\perp \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\rangle \right| \\
& \quad + \frac{1}{T} \sum_{t \in [T]} \left| \left\langle P \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t), \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\rangle \right| \\
& \leq \frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 + \sqrt{\frac{1}{T} \sum_{t \in [T]} \left\| P^\perp \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2} \sqrt{\frac{1}{T} \sum_{t \in [T]} \left\| P^\perp \Sigma^{1/2} B^* w_t^* \right\|_2^2} \\
& \quad + \sqrt{\frac{1}{T} \sum_{t \in [T]} \left\| P \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2} \sqrt{\frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} \delta_t^* - P^\perp \Sigma^{1/2} \hat{\theta}_t \right\|_2^2}.
\end{aligned}$$

This is a quadratic inequality in the two terms on the left-hand side, and so by applying Proposition I.2,

$$\begin{aligned}
\frac{1}{T} \sum_{t \in [T]} \left\| P^\perp \Sigma^{1/2} B^* w_t^* \right\|_2^2 & \lesssim \frac{1}{T} \sum_{t \in [T]} \left\| P^\perp \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 + \frac{1}{T} \sum_{t \in [T]} \left\| P \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 \\
& \quad + \frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 \\
& \lesssim \frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2,
\end{aligned}$$

where the last line follows by orthogonality. Finally, by Proposition I.3,

$$\frac{1}{T} \sum_{t \in [T]} \left\| P^\perp \Sigma^{1/2} B^* w_t^* \right\|_2^2 = \frac{1}{T} \left\| P^\perp \Sigma^{1/2} B^* W^* \right\|_F^2 \geq \frac{\sigma_k(W^*)^2}{T} \left\| P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* \right\|_F^2,$$

which together with the diversity assumption in Assumption 3.2 yields the final bound

$$\left\| P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* \right\|_F^2 \lesssim k \left[\frac{1}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (\theta_t^* - \hat{\theta}_t) \right\|_2^2 \right]. \quad \square$$

B.2 Target Guarantees for the Linear Setting

Having established a connection between the performance on the source tasks and the difference in the spans of B^* and B_0 , we can now analyze the performance of the target training procedure. First, we bound the performance of nearly optimal points in \mathcal{C}_β for several possible choices of c_1, c_2 .

Lemma B.5 (Statistical Rates for \mathcal{L}_β). *Assume that (Δ, w) is ζ -suboptimal for \mathcal{L}_β with the constraint set $\mathcal{C}_\beta = \{(\Delta, w) \mid \|\Delta\|_F \leq c_1/\beta, \|w\|_2 \leq c_2/\beta\}$, i.e.*

$$\mathcal{L}_\beta(\Delta, w) \leq \min_{(\Delta', w') \in \mathcal{C}_\beta} \mathcal{L}_\beta(\Delta', w') + \zeta.$$

We write $\hat{\theta}$ for the predictor corresponding to (Δ, w) , i.e. $\hat{\theta} = \beta(A_{B_0} + \Delta)(w_0 + w)$. Now, let

$$\begin{aligned}
\bar{w} &= (B_0 X^\top X B_0)^\dagger B_0^\top X^\top X \theta^* \\
\bar{\delta} &= \theta^* - B_0 \bar{w} \\
\tilde{w} &= (B_0^\top X^\top X B_0)^\dagger B_0^\top X^\top X B^* w^*.
\end{aligned}$$

Note that $XB_0\bar{w} = P_{XB_0}X\theta^*$, $X\bar{\delta} = P_{XB_0}^\perp X\theta^*$, and $XB\tilde{w} = P_{XB_0}XB^*w^*$. Then, assuming $\beta > \max(\|\bar{w}\|_2, \|\tilde{w}\|_2)$,

$$\frac{1}{2n_T} \|X(\theta^* - \hat{\theta})\|_2^2 \lesssim \zeta + \begin{cases} \left\| P_{\Sigma^{1/2}B_0}^\perp \Sigma^{1/2} B^* w^* \right\|_2^2 + \frac{\sigma^2}{n_T} \left(k + \log \frac{1}{\delta} \right) + \frac{\sigma\delta_0}{\sqrt{n_T}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{1}{\delta} \right)} & c_1 = \delta_0, c_2 = \|\tilde{w}\|_2 \\ \frac{\sigma^2}{n_T} \left(k + \log \frac{1}{\delta} \right) + \frac{\sigma \|\bar{\delta}\|_2}{\sqrt{n_T}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{1}{\delta} \right)} & c_1 = \|\bar{\delta}\|_2, c_2 = \|\bar{w}\|_2 \\ \frac{\sigma^2}{n_T} \left(d + \log \frac{1}{\delta} \right) & c_1 = \|\theta^*\|_2, c_2 = 0 \end{cases}$$

with probability at least $1 - \delta/3$.

Proof. We proceed by proving the three cases separately. Throughout the proof, we instantiate the high-probability event in Lemma B.2, which guarantees that

$$0.9\Sigma \preceq \frac{1}{n_T} X^\top X \preceq 1.1\Sigma.$$

Recall that this event occurs with probability at least $1 - \delta/9$.

$$c_1 = \delta_0, c_2 = \|\tilde{w}\|_2$$

Due to the choice of c_1 and c_2 , there exists a parameter in \mathcal{C}_β corresponding to the prediction vector $P_{XB_0}X\theta^* + X\delta^*$. Writing the corresponding basic inequality, we thus have that

$$\begin{aligned} \frac{1}{2n_T} \|y - X\hat{\theta}\|_2^2 - \zeta &\leq \frac{1}{2n_T} \|P_{XB_0}^\perp XB^*w^*\|_2^2 + \frac{1}{n_T} \langle z, P_{XB_0}^\perp XB^*w^* \rangle + \frac{1}{2n_T} \|z\|_2^2 \\ \implies \frac{1}{2n_T} \|X(\theta^* - \hat{\theta})\|_2^2 &\leq \zeta + \frac{1}{2n_T} \|P_{XB_0}^\perp XB^*w^*\|_2^2 + \frac{1}{n_T} \langle z, P_{XB_0}^\perp XB^*w^* \rangle \\ &\quad - \frac{1}{n_T} \langle z, X(\theta^* - \hat{\theta}) \rangle. \end{aligned}$$

Simplifying further,

$$\begin{aligned} &\frac{1}{n_T} \langle z, P_{XB_0}^\perp XB^*w^* \rangle - \frac{1}{n_T} \langle z, X(\theta^* - \hat{\theta}) \rangle \\ &= \frac{1}{n_T} \langle z, P_{XB_0}X(\hat{\theta} - \theta^*) \rangle - \frac{1}{n_T} \langle z, P_{XB_0}^\perp X(\delta^* - \beta\Delta(w_0 + w)) \rangle. \end{aligned}$$

Now, by the Hanson-Wright inequality, we can bound the last term as

$$\begin{aligned} -\frac{1}{n_T} \langle X^\top P_{XB_0}^\perp z, \delta^* - \beta\Delta(w_0 + w) \rangle &\leq \frac{1}{n_T} \|X^\top P_{XB_0}^\perp z\|_2 \|\delta^* + \beta\Delta(w_0 + w)\|_2 \\ &\lesssim \frac{\sigma\delta_0}{\sqrt{n_T}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{1}{\delta} \right)} \end{aligned}$$

with probability at least $1 - \delta/9$. Therefore, we can rewrite the prior basic inequality as

$$\begin{aligned} \frac{1}{2n_T} \|X(\theta^* - \hat{\theta})\|_2^2 &\lesssim \underbrace{\frac{1}{2n_T} \|P_{XB_0}^\perp XB^*w^*\|_2^2 + \frac{\sigma\delta_0}{\sqrt{n_T}} \sqrt{\text{tr} \Sigma \left(1 + \log \frac{1}{\delta} \right)}}_{=: \mathcal{C}} + \zeta \\ &\quad + \frac{1}{n_T} \langle z, P_{XB_0}X(\theta^* - \hat{\theta}) \rangle. \end{aligned}$$

Now, note that we can form the quadratic inequality

$$\frac{1}{n_T} \left\| P_{XB_0} X(\theta^* - \hat{\theta}) \right\|_2^2 \leq \frac{1}{n_T} \|P_{XB_0} z\|_2 \left\| P_{XB_0} X(\theta^* - \hat{\theta}) \right\|_2 + C,$$

and thus by applying Proposition I.2 to solve the inequality and noting that $P_{XB_0} z/\sigma$ is distributed as a chi-squared random variable with k degrees of freedom,

$$\begin{aligned} \frac{1}{n_T} \left\| X(\theta^* - \hat{\theta}) \right\|_2^2 &\lesssim \zeta + \left\| P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* w^* \right\|_2^2 \\ &\quad + \frac{\sigma^2}{n_T} \left(k + \log \frac{1}{\delta} \right) + \frac{\sigma \delta_0}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma \left(1 + \log \frac{1}{\delta} \right)} \end{aligned}$$

with probability at least $1 - \delta/9$, which is the bound that we wanted to show.

$$c_1 = \|\bar{\delta}\|_2, c_2 = \|\bar{w}\|_2$$

Due to the choice of c_1 and c_2 , there exists a parameter in \mathcal{C}_β corresponding to a predictor that agrees with θ^* on the target samples. Therefore,

$$\begin{aligned} \frac{1}{2n_T} \left\| X(\theta^* - \hat{\theta}) \right\|_2^2 &\leq \zeta + \frac{1}{2n_T} \|z\|_2^2 \\ \implies \frac{1}{2n_T} \left\| X(\theta^* - \hat{\theta}) \right\|_2^2 &\leq \zeta + \frac{1}{n_T} \left\langle z, P_{XB_0} X(\theta^* - \hat{\theta}) \right\rangle + \frac{1}{n_T} \left\langle z, P_{X B_0}^\perp X(\theta^* - \hat{\theta}) \right\rangle. \end{aligned} \tag{9}$$

We proceed with an argument similar to that used in the source guarantee, albeit simpler since the representation B_0 is fixed (and thus no covering argument is required). Along these lines, we bound the first term using the low-rank of B_0 . Via projections,

$$\begin{aligned} \frac{1}{n_T} \left\| P_{XB_0} X(\theta^* - \hat{\theta}) \right\|_2^2 &\lesssim \zeta + \frac{1}{n_T} \|P_{XB_0} z\|_2 \left\| P_{XB_0} X(\theta^* - \hat{\theta}) \right\|_2 + \frac{1}{n_T} \left\langle z, P_{X B_0}^\perp X(\theta^* - \hat{\theta}) \right\rangle. \end{aligned}$$

Therefore, by applying Proposition I.2 to solve the quadratic inequality, we have that with probability at least $1 - \delta/9$,

$$\frac{1}{2n_T} \left\| X(\theta^* - \hat{\theta}) \right\|_2^2 \leq \zeta + \frac{\sigma^2}{n_T} \left(k + \log \frac{1}{\delta} \right) + \frac{1}{n_T} \left\langle z, P_{X B_0}^\perp X(\theta^* - \hat{\theta}) \right\rangle.$$

To bound the second term, we simply make use of the norm constraints defining the feasible set \mathcal{C}_β , which we note is analogous to the low-norm sub-argument of the source guarantee. Formally, the Hanson-Wright inequality implies that with probability at least $1 - \delta/9$,

$$\begin{aligned} \frac{1}{n_T} \left\langle z, P_{X B_0}^\perp X(\theta^* - \hat{\theta}) \right\rangle &= \frac{1}{\sqrt{n_T}} \left\| \frac{1}{\sqrt{n_T}} X^\top P_{X B_0}^\perp z \right\|_2 \|\bar{\delta} + \beta \Delta w_0 + \beta \Delta w\|_2 \\ &\lesssim \frac{\sigma \|\bar{\delta}\|_2}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma \left(1 + \log \frac{1}{\delta} \right)}. \end{aligned}$$

Putting everything together, we thus have that

$$\frac{1}{2n_T} \left\| X(\theta^* - \hat{\theta}) \right\|_2^2 \lesssim \zeta + \frac{\sigma^2}{n_T} \left(k + \log \frac{1}{\delta} \right) + \frac{\sigma \|\bar{\delta}\|_2}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma \left(1 + \log \frac{1}{\delta} \right)}.$$

$$c_1 = \|\theta^*\|_2, c_2 = 0$$

Due to the choice of c_1 and c_2 , there exists a parameter in \mathcal{C}_β corresponding to a predictor that agrees with θ^* on the target samples. Therefore, we can write the basic inequality

$$\frac{1}{n_T} \left\| y - X \hat{\theta} \right\|_2^2 \leq \zeta + \frac{1}{n_T} \|z\|_2^2 \implies \frac{1}{n_T} \left\| X(\theta^* - \hat{\theta}) \right\|_2^2 \leq \zeta + \frac{1}{n_T} \left\langle z, X(\theta^* - \theta) \right\rangle.$$

Now, noting that $P_X z / \sigma$ is a chi-squared random variable with d degrees of freedom, we have that with probability at least $1 - 2\delta/9$,

$$\frac{1}{n_T} \left\langle z, X(\theta^* - \hat{\theta}) \right\rangle \leq \frac{1}{n_T} \|P_X z\|_2 \|X(\theta^* - \hat{\theta})\|_2 \leq \frac{1}{n_T} \sqrt{d + \log \frac{1}{\delta}} \|X(\theta^* - \hat{\theta})\|_2.$$

Therefore, by solving the resulting quadratic inequality via Proposition I.2, we obtain the bound

$$\frac{1}{n_T} \|X(\theta^* - \hat{\theta})\|_2^2 \leq \zeta + \frac{\sigma^2}{n_T} \left(d + \log \frac{1}{\delta} \right).$$

Observe that all relevant high-probability events for each case occur simultaneously with probability at least $1 - \delta/3$, as desired. \square

B.3 Optimization Landscape during Target Time Training

Having derived statistical rates on nearly-optimal points for several choices of \mathcal{C}_β in the prior section, all that remains to be shown is that projected gradient descent can indeed find such points. In particular, we will demonstrate that for large enough β , the optimization landscape induced by \mathcal{L}_β is approximately convex. We do so by demonstrating that the objective satisfies the assumptions outlined in Section H, and thus the accompanying guarantees for projected gradient descent hold.

Lemma B.6 (Approximate linearity of function class). *Let $g_\theta(x) = \beta x^\top (A_{B_0} + \Delta)(w_0 + w)$ for $\theta = (\Delta, w) \in \mathcal{C}_\beta$, where \mathcal{C}_β is considered as a subset of \mathbb{R}^{dk+k} . Then, assuming the high-probability event in Lemma B.2 holds, then*

$$\sup_{\theta \in \mathcal{C}_\beta} \frac{1}{n_T} \sum_{i \in [n_T]} \|\nabla_{\theta}^2 g_\theta(x_i)\|_2^2 \lesssim \beta^2 \text{tr } \Sigma \quad \text{and} \quad \frac{1}{n_T} \sum_{i \in [n_T]} \|\nabla_{\theta} g_\theta(x_i)\|_2^2 \lesssim \beta^2 \text{tr } \Sigma.$$

Proof. To bound the average squared Hessian operator norm, note that $\nabla_{\theta}^2 g_\theta(x_i)[\Delta, w] = \beta x_i^\top \Delta w$, which is independent of θ . Then, by the variational characterization of the operator norm,

$$\begin{aligned} \|\nabla_{\theta}^2 g_\theta(x_i)\|_2 &= \sup_{\|\Delta\|_F^2 + \|w\|_2^2 \leq 1} |\nabla_{\theta}^2 g_\theta(x_i)[\Delta, w]| = \sup_{\|\Delta\|_F^2 + \|w\|_2^2 \leq 1} \beta |x_i^\top \Delta w| \\ &\leq \sup_{\|\Delta\|_F^2 + \|w\|_2^2 \leq 1} \beta \|x_i\|_2 \|\Delta\|_F \|w\|_2 \\ &\leq \beta \|x_i\|_2. \end{aligned}$$

and thus $\|\nabla_{\theta}^2 g_\theta(x_i)\|_2^2 \lesssim \beta^2 \|x_i\|_2^2$. Consequently,

$$\sup_{\theta \in \mathcal{C}_\beta} \frac{1}{n_T} \sum_{i \in [n_T]} \|\nabla_{\theta}^2 g_\theta(x_i)\|_2^2 \leq \frac{\beta^2}{n_T} \text{tr } X X^\top \lesssim \beta^2 \text{tr } \Sigma.$$

We now proceed to bound the squared norm of the gradient. Observe that the gradient is given by $\nabla_{\theta} g_\theta(x_i) = \beta [A_{B_0}^\top x_i, \text{vec}(x_i w_0^\top)]$, and therefore,

$$\begin{aligned} \frac{1}{n_T} \sum_{i \in [n_T]} \|\nabla_{\theta} g_\theta(x_i)\|_2^2 &= \frac{\beta^2}{n_T} \sum_{i \in [n_T]} \|A_{B_0}^\top x_i\|_2^2 + \|x_i w_0^\top\|_F^2 \lesssim \frac{\beta^2}{n_T} \sum_{i \in [n_T]} \|x_i\|_2^2 = \frac{\beta^2}{n_T} \text{tr } X X^\top \\ &\lesssim \beta^2 \text{tr } \Sigma, \end{aligned}$$

where the first inequality uses the fact that $\|A_{B_0}^\top x_i\|_2^2 = 2 \|B_0^\top x_i\|_2^2 = 2 \|P_{B_0} x_i\|_2^2 \lesssim \|x_i\|_2^2$, by the definition of A_{B_0} and assumed orthogonality of B_0 . \square

Lemma B.7 (\mathcal{L}_β is Lipschitz). *Let $g_\theta(x) = \beta x^\top (A_{B_0} + \Delta)(w_0 + w)$ for $\theta = (\Delta, w) \in \mathcal{C}_\beta$, where \mathcal{C}_β is considered as a subset of \mathbb{R}^{dk+k} . Furthermore, assume that the high-probability event in Lemma B.2 holds. Then, with probability at least $1 - \delta/3$ over the draw of n_T labels, we have that for any $c_1, c_2 > 0$ and $\beta^2 > c_1^2 + c_2^2$,*

$$\sup_{\theta \in \mathcal{C}_\beta} \|\nabla_g \mathcal{L}_\beta(g_\theta(X))\|_2^2 \lesssim \frac{1}{n_T} \left[\|\Sigma\|_2 \left(\|\theta^*\|_2^2 + c_1^2 + c_2^2 \right) + \sigma^2 \left(1 + \log \frac{1}{\delta} \right) \right].$$

Proof. We assume that $\|z\|_2^2/n_T \lesssim \sigma^2[1 + \log(1/\delta)]$, which occurs with probability at least $1 - \delta/3$ via standard tail bounds on chi-squared random variables. Then, the gradient of the loss with respect to the predictions is given by

$$\nabla_g \mathcal{L}_\beta(g_\theta(X)) = \frac{1}{n_T}(y - g_\theta(X)) = \frac{1}{n_T}[X(\theta^* - \beta A_{B_0} w - \beta \Delta w_0 - \beta \Delta w) + z].$$

Therefore,

$$\begin{aligned} \sup_{\theta \in \mathcal{C}_\beta} \|\nabla_g \mathcal{L}_\beta(g_\theta(X))\|_2^2 &\lesssim \sup_{\theta \in \mathcal{C}_\beta} \frac{1}{n_T^2} \|X(\theta^* - \beta A_{B_0} w - \beta \Delta w_0 - \beta \Delta w)\|_2^2 + \frac{1}{n_T^2} \|z\|_2^2 \\ &\lesssim \frac{1}{n_T} \left[\|\Sigma\|_2 \left(\|\theta^*\|_2^2 + c_1^2 + c_2^2 \right) + \sigma^2 \left(1 + \log \frac{1}{\delta} \right) \right]. \quad \square \end{aligned}$$

Lemma B.8. *Let \tilde{w} , \bar{w} , and $\bar{\delta}$ be defined as in Lemma B.5. Then, assuming that the high-probability event in Lemma B.2 holds, the three quantities above are norm-bounded by $\kappa^{1/2}(\|B^* w^*\|_2 + \delta_0)$, up to constant factors. Furthermore, we can also bound $\|\bar{\delta}\|_2$ as*

$$\|\bar{\delta}\|_2 \leq \delta_0 + \left(\frac{\|w^*\|_2 + \delta_0 \kappa^{1/2}}{\lambda_{\min}^{1/2}} \right) \|P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^*\|_F.$$

Proof. Throughout the proof, we will write λ_{\max} and λ_{\min} as shorthand for $\lambda_{\max}(\Sigma)$ and $\lambda_{\min}(\Sigma)$, respectively. By definition, since B_0 is assumed to be orthonormal, the concentration of the sample covariance in Lemma B.2 implies that

$$\begin{aligned} \|\bar{w}\|_2^2 &\lesssim \frac{1}{n_T \lambda_{\min}} \|X B_0 \bar{w}\|_2^2 = \frac{1}{n_T \lambda_{\min}} \|P_{X B_0} X \theta^*\|_2^2 \lesssim \frac{1}{\lambda_{\min}} \|\Sigma^{1/2} \theta^*\|_2^2 \leq \kappa \|\theta^*\|_2^2 \\ &\leq \kappa (\|B^* w^*\|_2 + \delta_0)^2. \end{aligned}$$

Following similar arguments for \tilde{w} and $\bar{\delta}$, we obtain the same bounds. Thus, we have demonstrated that all three quantities are indeed norm-bounded by $\kappa^{1/2}(\|B^* w^*\|_2 + \delta_0)$, up to constant factors.

Finally, we proceed to derive the final bound on $\bar{\delta}$. Note that $X \bar{\delta} = X \delta^* - P_{X B_0} X \delta^* + P_{X B_0}^\perp X B^* w^*$. Therefore,

$$\begin{aligned} \|\bar{\delta}\|_2 &\lesssim \delta_0 + \frac{1}{\sqrt{n_T \lambda_{\min}}} [\|P_{X B_0} X \delta^*\|_2 + \|P_{X B_0}^\perp X B^* w^*\|_2] \\ &\lesssim \delta_0 + \frac{1}{\lambda_{\min}^{1/2}} \left[\|P_{\Sigma^{1/2} B_0} \Sigma^{1/2} \delta^*\|_2 + \|P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* w^*\|_2 \right] \\ &\lesssim \delta_0 + \frac{1}{\lambda_{\min}^{1/2}} \left[\|\Sigma^{1/2} \delta^*\|_2 \|P_{\Sigma^{1/2} B_0} P_{\Sigma^{1/2} B^*}^\perp\|_F + \|w^*\|_2 \|P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^*\|_F \right] \end{aligned}$$

Now, by applying the properties of the trace operator and Proposition I.3,

$$\begin{aligned} \|P_{\Sigma^{1/2} B_0} P_{\Sigma^{1/2} B^*}^\perp\|_F^2 &= \text{tr } P_{\Sigma^{1/2} B_0} P_{\Sigma^{1/2} B^*}^\perp P_{\Sigma^{1/2} B_0} = \text{tr } (I - P_{\Sigma^{1/2} B^*}) P_{\Sigma^{1/2} B_0} \\ &= \text{tr } P_{\Sigma^{1/2} B_0} - P_{\Sigma^{1/2} B^*} P_{\Sigma^{1/2} B_0} \\ &\leq \text{tr } P_{\Sigma^{1/2} B^*} - P_{\Sigma^{1/2} B^*} P_{\Sigma^{1/2} B_0} = \text{tr } P_{\Sigma^{1/2} B^*} (I - P_{\Sigma^{1/2} B_0}) \\ &= \text{tr } P_{\Sigma^{1/2} B^*} P_{\Sigma^{1/2} B_0}^\perp P_{\Sigma^{1/2} B^*} = \|P_{\Sigma^{1/2} B_0}^\perp P_{\Sigma^{1/2} B^*}\|_F^2 \\ &\leq \frac{1}{\lambda_{\min}} \|P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^*\|_F^2, \end{aligned}$$

and thus

$$\begin{aligned} \|\bar{\delta}\|_2 &\lesssim \delta_0 + \frac{1}{\lambda_{\min}^{1/2}} \left[\frac{1}{\lambda_{\min}^{1/2}} \|\Sigma^{1/2} \delta^*\|_2 + \|w^*\|_2 \right] \|P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^*\|_F \\ &\leq \delta_0 + \left(\frac{\|w^*\|_2 + \delta_0 \kappa^{1/2}}{\lambda_{\min}^{1/2}} \right) \|P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^*\|_F. \quad \square \end{aligned}$$

B.4 Deducing Theorem 3.1

Having proven all the previous results, we can now assemble the main claim in Theorem 3.1. Recall that we have defined the rates

$$\begin{aligned}
r_S(n_S, T) &:= \frac{\sigma^2}{n_S T} \left(kT + kd \log \kappa n_S + \log \frac{1}{\delta} \right) + \frac{\sigma \delta_0 \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log \kappa n_S + \log \frac{1}{\delta}} \\
&\quad + \frac{\sigma \delta_0}{\sqrt{n_S}} \sqrt{\text{tr } \Sigma \left(1 + \log \frac{T}{\delta} \right)} \\
r_T^{(1)}(n_T) &:= kr_S(n_S, T) + \frac{\sigma^2}{n_T} \left(k + \log \frac{1}{\delta} \right) + \frac{\sigma \delta_0}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma \left(1 + \log \frac{1}{\delta} \right)} \\
r_T^{(2)}(n_T) &:= \frac{\sigma^2}{n_T} \left(k + \log \frac{1}{\delta} \right) \\
&\quad + \frac{\sigma}{\sqrt{n_T}} \left[\delta_0 + \left(\frac{\|w^*\|_2 + \delta_0 \kappa^{1/2}}{\lambda_{\min}^{1/2}(\Sigma)} \right) \sqrt{kr_S(n_S, T)} \right] \sqrt{\text{tr } \Sigma \left(1 + \log \frac{1}{\delta} \right)} \\
r_T^{(3)}(n_T) &:= \frac{\sigma^2}{n_T} \left(d + \log \frac{1}{\delta} \right).
\end{aligned}$$

Theorem 3.1 (Performance guarantee, linear representations). *Assume that Assumptions 3.1 and 3.2 hold, $n_S \gg \rho^4(d + \log(T/\delta))$, and $n_T \gg \rho^4(d + \log(1/\delta))$. Then there are $(\lambda, \gamma, \beta, T_{\text{PGD}}, \eta, c_1, c_2)$ such that the training procedure in Section 3.2, with high probability, finds θ achieving excess risk bounded as*

$$\mathbb{E} [(x^\top \theta^* - x^\top \theta)^2] \lesssim \min(r_T^{(1)}(n_T), r_T^{(2)}(n_T), r_T^{(3)}(n_T)).$$

Proof. Assume that during source training, the regularization parameters (λ, γ) are chosen according to Lemma B.3. We instantiate the high-probability events in Lemmas B.3, B.5 and B.7. These events occur altogether with probability at least $1 - \delta$. Throughout the proof, we define $r(n) = \min(r_T^{(1)}(n), r_T^{(2)}(n), r_T^{(3)}(n))$.

Given all the events above, we know that the optimization landscape induced by \mathcal{L}_β is well-behaved. That is, the function class is approximately linear in the parameters by Lemma B.6 and the loss is Lipschitz by Lemma B.7. Furthermore, for any of the proposed feasible sets in Lemma B.5, the bounds on relevant quantities provided by Lemma B.8 imply that

$$c_1^2 + c_2^2 \lesssim \underbrace{\kappa(\|B^* w^*\|_2^2 + \delta_0^2)}_{=: R^2} \implies \sup_{\theta \in \mathcal{C}_\beta} \|\nabla_g \mathcal{L}_\beta(g_\theta(X))\|_2^2 \lesssim \frac{1}{n_T} \underbrace{\left[\|\Sigma\|_2 R^2 + \sigma^2 \left(1 + \log \frac{1}{\delta} \right) \right]}_{=: \alpha^2}.$$

Therefore, we have demonstrated that \mathcal{L}_β satisfies the assumptions in Section H for any $\beta > R^2$. Consequently, by running projected gradient descent with parameters

$$\beta = \max \left(R^2, \frac{\alpha R^2 \sqrt{\text{tr } \Sigma}}{r(n_T)} \right) \quad \text{and} \quad T_{\text{PGD}} = \frac{\alpha^2 R^2 \text{tr } \Sigma}{r(n_T)^2},$$

and setting η as in Theorem H.1, we can guarantee that projected gradient descent finds an $r(n_T)$ -suboptimal point. Therefore, by choosing (c_1, c_2) in order to achieve the minimal rate $r(n_T)$ in Lemma B.5, we can guarantee that

$$\min_t \frac{1}{n_T} \|X(\theta_t - \theta^*)\|_2^2 \lesssim r(n_T).$$

Finally, due to target covariance concentration as guaranteed by the event in Lemma B.2, we thus have that the excess risk of the best predictor found by the algorithm on the target task is bounded as

$$\min_t \mathbb{E} [(x^\top \theta_t - x^\top \theta_t^*)^2] = \min_t \left\| \Sigma^{1/2}(\theta_t - \theta^*) \right\|_2^2 \lesssim \min(r_T^{(1)}(n), r_T^{(2)}(n), r_T^{(3)}(n)). \quad \square$$

C Linear Hard Case: Construction, Proofs, and Experiments

C.1 Hard Case Construction

In the following section, we will provide the construction of a task distribution family that is used in proving Theorem 3.2 in the main text.

C.1.1 Formal Construction

We now proceed with a formal construction that satisfies the conditions of Section 3.1. We provide the intermediate results used to prove Theorem 3.2, but leave their proofs for Section C.2 to simplify the presentation in this section.

Fix an $\varepsilon \in (0, 1)$, and let p be a Gaussian distribution on \mathbb{R}^d with block-diagonal covariance Σ

$$\Sigma := \begin{bmatrix} \varepsilon I_{d-k} & 0 \\ 0 & I_k \end{bmatrix}.$$

We define $E^*, E_k \subset \mathbb{R}^d$ to be the two eigenspaces of Σ determined by the two blocks, *i.e.*

$$E^* = \text{Col} \begin{bmatrix} \varepsilon I_{d-k} \\ 0 \end{bmatrix} \quad \text{and} \quad E_k = \text{Col} \begin{bmatrix} 0 \\ I_k \end{bmatrix}.$$

Then, for an orthogonal matrix $B \in \mathbb{R}^{d \times k}$ with $\text{Col } B^* \subseteq E^*$, define a corresponding task distribution given by

$$\theta = \frac{1}{\sqrt{2\varepsilon}} Bw + \delta, \tag{10}$$

where w and δ are uniformly sampled from the unit spheres in \mathbb{R}^k and E_k , respectively.

Recall that in the linear representation setting, FROZENREP optimizes the following objective to obtain a representation \hat{B} :

$$\hat{B} = \underset{B}{\text{argmin}} \min_{w_t} \frac{1}{2n_S T} \sum_{t \in [T]} \|y_t - X_t B w_t\|_2^2.$$

First, we characterize the span of \hat{B} in the limit of infinite source tasks and data. Intuitively, $\text{span}\{B^* w_t^*\}$ and $\text{span}\{\Delta_t^* w_t^*\}$ correspond to the green and red “vectors” in Figure 2, respectively, and thus we expect FROZENREP to learn E_k .

Lemma C.1 (FROZENREP learns incorrect space). *Fix an orthogonal matrix $B^* \in \mathbb{R}^{d \times k}$ with $\text{Col } B^* \subseteq E^*$, and assume that we sample task weights from the distribution in (10). Then, with infinitely many tasks and per-task samples ($n_S, T \rightarrow \infty$), $\text{Col } \hat{B} = E_k$.*

Although “incorrect”, it is unclear *a priori* that learning $\text{Col } E_k$ is undesirable performance-wise. The next result formalizes the resulting degradation in performance due to learning $\text{Col } E_k$.

Theorem C.1 (FROZENREP minimax bound). *For an orthogonal matrix $B^* \in \mathbb{R}^{d \times k}$ whose column space lies in E^* , let S_{B^*} be the set*

$$S_{B^*} := \left\{ \frac{1}{\sqrt{2\varepsilon}} B^* w + \delta \mid \|w\|_2 \leq 1, \|\delta\|_2 \leq 1, \delta \in E_k \right\}.$$

We consider the following procedure:

1. *We draw n_T samples for target time training, which are collected into a matrix X .*
2. *Player chooses a target-time estimator $\hat{\theta} = \bar{B}\hat{w} + \hat{\delta}$, where \bar{B} , \hat{w} , and $\hat{\delta}$ are measurable function of (X, y) , and $\text{Col } \bar{B} = \text{Col } \hat{B}$.*
3. *Player chooses target-time estimator $\hat{\theta} = \hat{B}\hat{w} + \hat{\delta}$, where \hat{w} and $\hat{\delta}$ are measurable functions of (\hat{B}, X, y) .*

4. Adversary chooses an orthogonal matrix $B^* \in \mathbb{R}^{d \times k}$ satisfying $\text{Col } B^* \subset E_k$, and a target time predictor $\theta^* \in S_{B^*}$.
5. FROZENREP returns a representation \hat{B} under the setting of Lemma C.1 with the task distribution determined by B^* .
6. Target time samples are generated using $y \sim \mathcal{N}(X\theta^*, \sigma^2 I_n)$, and the player estimator is evaluated.

Then, with probability at least $1 - \delta$ over the draw of X , we have that

$$\min_{\hat{w}, \hat{\delta}} \max_{\substack{B^* \\ \theta^* \in S_{B^*}}} \mathbb{E} \left[\frac{1}{n_T} \left\| X(\theta^* - \hat{B}\hat{w} - \hat{\delta}) \right\|_2^2 \right] \gtrsim \frac{\sigma^2 d}{n_T},$$

where the expectation is over the randomness in the labels $y \sim \mathcal{N}(X\theta^*, \sigma^2 I_{n_T})$. Note that the minimization over \hat{w} and $\hat{\delta}$ is performed over the set of measurable functions from (X, y) to \mathbb{R}^k and \mathbb{R}^d , respectively.

The result above comprises the minimax bound in Theorem 3.2. In contrast, by specializing the ADAPTREP performance bound in Theorem 3.1, we obtain the following corollary:

Corollary C.1. Set $k = \Theta(1)$, $d \gg k$, and $\varepsilon = k/d$. Furthermore, assume that $n_S T \gtrsim d^2$, and $n_S \geq n_T \asymp d$. Then, for a fixed target task in S_{B^*} as defined in Theorem C.1, with probability at least $1 - \delta$ over the draw of samples, the procedure in Section 3.2 achieves excess risk bounded as

$$\mathbb{E} \left[(x^\top \theta^* - x^\top \hat{\theta})^2 \right] \lesssim \min(\sigma/\sqrt{n_T}, \sigma^2 d/n_T).$$

Thus, we have constructed the desired task distribution family, proving Theorem 3.2.

C.2 Proofs for Section 3.4

Lemma C.1 (FROZENREP learns incorrect space). Fix an orthogonal matrix $B^* \in \mathbb{R}^{d \times k}$ with $\text{Col } B^* \subseteq E^*$, and assume that we sample task weights from the distribution in (10). Then, with infinitely many tasks and per-task samples ($n_S, T \rightarrow \infty$), $\text{Col } \hat{B} = E_k$.

Proof. Fix a B . For any $t \in [T]$, as $n_S \rightarrow \infty$, we have that

$$\begin{aligned} \min_{w_t} \frac{1}{n_S} \|y_t - X_t B w_t\|_2^2 &= \min_{w_t} \frac{1}{n_S} \|X_t(\theta_t^* - B w_t) + z_t\|_2^2 \\ &= \min_{w_t} \left\| \Sigma^{1/2}(\theta_t^* - B w_t) \right\|_2^2 + \sigma^2 \\ &= \left\| P_{\Sigma^{1/2} B}^\perp \Sigma^{1/2} \theta_t^* \right\|_2^2 + \sigma^2. \end{aligned}$$

Therefore, we can rewrite the objective defining \hat{B} as

$$\begin{aligned} \hat{B} &= \underset{B}{\operatorname{argmin}} \frac{1}{T} \sum_{t \in [T]} \left\| P_{\Sigma^{1/2} B}^\perp \Sigma^{1/2} \theta_t^* \right\|_2^2 \\ &= \underset{B}{\operatorname{argmin}} \operatorname{tr} P_{\Sigma^{1/2} B}^\perp \Sigma^{1/2} \left[\frac{1}{T} \sum_{t \in [T]} \theta_t^* (\theta_t^*)^\top \right] \Sigma^{1/2} \\ &= \underset{B}{\operatorname{argmin}} \operatorname{tr} P_{\Sigma^{1/2} B}^\perp \left[\frac{1}{2} P_{B^*} + P_{E_k} \right] \\ &= \underset{B}{\operatorname{argmin}} \frac{1}{2} \left\| P_{\Sigma^{1/2} B}^\perp P_{B^*} \right\|_F^2 + \left\| P_{\Sigma^{1/2} B}^\perp P_{E_k} \right\|_F^2, \end{aligned}$$

where the third equality uses the definition of the task distribution in (10). Now, let S be the space defined as $E_k + \text{Col } B^*$. Observe that

$$\begin{aligned}\|P_{\Sigma^{1/2}B}^\perp P_S\|_F^2 &= \text{tr } P_{\Sigma^{1/2}B}^\perp P_S \\ &= \text{tr } P_{\Sigma^{1/2}B}^\perp P_B^* + \text{tr } P_{\Sigma^{1/2}B}^\perp P_{E_k} \\ &= \|P_{\Sigma^{1/2}B}^\perp P_{B^*}\|_F^2 + \|P_{\Sigma^{1/2}B}^\perp P_{E_k}\|_F^2.\end{aligned}$$

By putting everything together, we thus see that

$$\hat{B} = \underset{B}{\text{argmin}} \|P_{\Sigma^{1/2}B}^\perp P_S\|_F^2 + \|P_{\Sigma^{1/2}B}^\perp P_{E_k}\|_F^2.$$

This objective is minimized if and only if the span of $\Sigma^{1/2}B$ is exactly E_k . By inverting, this is only possible if $\text{Col } \hat{B} = E_k$. \square

Theorem C.1 (FROZENREP minimax bound). *For an orthogonal matrix $B^* \in \mathbb{R}^{d \times k}$ whose column space lies in E^* , let S_{B^*} be the set*

$$S_{B^*} := \left\{ \frac{1}{\sqrt{2\varepsilon}} B^* w + \delta \mid \|w\|_2 \leq 1, \|\delta\|_2 \leq 1, \delta \in E_k \right\}.$$

We consider the following procedure:

1. We draw n_T samples for target time training, which are collected into a matrix X .
2. Player chooses a target-time estimator $\hat{\theta} = \bar{B}\hat{w} + \hat{\delta}$, where \bar{B} , \hat{w} , and $\hat{\delta}$ are measurable function of (X, y) , and $\text{Col } \bar{B} = \text{Col } \hat{B}$.
3. Player chooses target-time estimator $\hat{\theta} = \hat{B}\hat{w} + \hat{\delta}$, where \hat{w} and $\hat{\delta}$ are measurable functions of (\hat{B}, X, y) .
4. Adversary chooses an orthogonal matrix $B^* \in \mathbb{R}^{d \times k}$ satisfying $\text{Col } B^* \subset E_k$, and a target time predictor $\theta^* \in S_{B^*}$.
5. FROZENREP returns a representation \hat{B} under the setting of Lemma C.1 with the task distribution determined by B^* .
6. Target time samples are generated using $y \sim \mathcal{N}(X\theta^*, \sigma^2 I_n)$, and the player estimator is evaluated.

Then, with probability at least $1 - \delta$ over the draw of X , we have that

$$\min_{\hat{w}, \hat{\delta}} \max_{\substack{B^* \\ \theta^* \in S_{B^*}}} \mathbb{E} \left[\frac{1}{n_T} \|X(\theta^* - \hat{B}\hat{w} - \hat{\delta})\|_2^2 \right] \gtrsim \frac{\sigma^2 d}{n_T},$$

where the expectation is over the randomness in the labels $y \sim \mathcal{N}(X\theta^*, \sigma^2 I_{n_T})$. Note that the minimization over \hat{w} and $\hat{\delta}$ is performed over the set of measurable functions from (X, y) to \mathbb{R}^k and \mathbb{R}^d , respectively.

Proof. We instantiate the event guaranteed by Lemma B.2 over the draw of target samples, which guarantees that

$$0.9\Sigma \lesssim \frac{1}{n_T} X^\top X \lesssim 1.1\Sigma,$$

with probability at least $1 - \delta$.

Recall that by Lemma C.1, FROZENREP will always find an orthogonal matrix whose column space is E_k , no matter how B^* is chosen. Therefore, we can rewrite the minimax expression as

$$\begin{aligned}\min_{\hat{w}, \hat{\delta}} \max_{\substack{B^*, \theta^* \\ \theta^* \in S_{B^*}}} \mathbb{E} \left[\frac{1}{n_T} \|X(\theta^* - \hat{B}\hat{w} - \hat{\delta})\|_2^2 \right] &= \min_{\hat{\theta}} \max_{\substack{B^*, \theta^* \\ \theta^* \in S_{B^*}}} \mathbb{E} \left[\frac{1}{n_T} \|X(\theta^* - \hat{\theta})\|_2^2 \right] \\ &\gtrsim \min_{\hat{\theta}} \max_{\substack{B^*, \theta^* \\ \theta^* \in S_{B^*}}} \mathbb{E} \left[\|\theta^* - \hat{\theta}\|_\Sigma^2 \right],\end{aligned}$$

where $\hat{\theta}$ is simply a measurable function of (X, y) to \mathbb{R}^d . To further simplify the problem, observe that if we define the set

$$T = \left\{ \theta \in \mathbb{R}^d \mid \|P_{E^*}\theta\|_2^2 \leq \frac{1}{2\varepsilon}, \|P_{E_k}\theta\|_2^2 \leq 1 \right\} \subseteq S_{B^*}$$

then since the estimator only depends on B^* through θ^* ,

$$\min_{\hat{w}, \hat{\delta}} \max_{\substack{B^*, \theta^* \\ \theta^* \in S_{B^*}}} \mathbb{E} \left[\frac{1}{n_T} \|X(\theta^* - \hat{B}\hat{w} - \hat{\delta})\|_2^2 \right] \gtrsim \min_{\hat{\theta}} \max_{\theta^* \in T} \mathbb{E} \left[\|(\theta^* - \hat{\theta})\|_{\Sigma}^2 \right]. \quad (11)$$

To lower bound the right-hand side of (11), we use local coverings in the Σ -norm and apply the Fano bound for minimax risk. Let $B := \left\{ \theta \in \mathbb{R}^d \mid \varepsilon \|P_{E^*}\theta\|_2^2 + \|P_{E_k}\theta\|_2^2 \leq 1 \right\}$ be the unit ball in the Σ -norm, so that $\frac{1}{\sqrt{2}}B \subseteq T$. Using a known volumetric argument, there exists a $(1/2)$ -packing of B in the Σ -norm with at least 2^d elements. Equivalently, there exists a 2δ -packing of $4\delta B$ with at least 2^d elements – let this packing be S . Then, for any $\theta, \theta' \in S$,

$$\begin{aligned} \text{KL}(\mathcal{N}(X\theta, \sigma^2 I_{n_T}) \parallel \mathcal{N}(X\theta', \sigma^2 I_{n_T})) &= \frac{1}{2\sigma^2} \|X(\theta - \theta')\|_2^2 \lesssim \frac{n_T}{2\sigma^2} \|\theta - \theta'\|_{\Sigma}^2 \\ &\leq \frac{32n_T}{\sigma^2} \delta^2. \end{aligned}$$

Therefore, by Fano's inequality, for any $\delta^2 \leq 1/32$ (which ensures that $S \subseteq 4\delta B \subseteq \frac{1}{\sqrt{2}}B \subseteq T$),

$$\begin{aligned} \min_{\hat{\theta}} \max_{\theta^* \in T} \mathbb{E} \left[\frac{1}{n} \|(\theta^* - \hat{\theta})\|_{\Sigma}^2 \right] &\geq \delta^2 \left(1 - \frac{32n}{\sigma^2 d \log 2} \delta^2 - \frac{1}{d} \right) \\ &\geq \delta^2 \left(\frac{3}{4} - \frac{32n}{\sigma^2 d \log 2} \delta^2 \right). \quad (d \geq 4) \end{aligned}$$

As long as $n_T \geq \frac{\log 2}{2} \sigma^2 d$, we can set $\delta^2 = \left(\frac{\log 2}{64} \right) \frac{\sigma^2 d}{n_T}$, and thus putting everything together, we have that

$$\min_{\hat{w}, \hat{\delta}} \max_{\substack{B^*, \theta^* \\ \theta^* \in S_{B^*}}} \mathbb{E} \left[\frac{1}{n} \|X(\theta^* - \hat{B}\hat{w} - \hat{\delta})\|_2^2 \right] \gtrsim \frac{\sigma^2 d}{n_T}. \quad \square$$

C.3 Simulations

In this section, we experimentally verify the linear hard case presented in Section 3.4. Since the empirical success of MAML and its variants has already been demonstrated extensively in practice and in existing work, it is not the focus of this section.

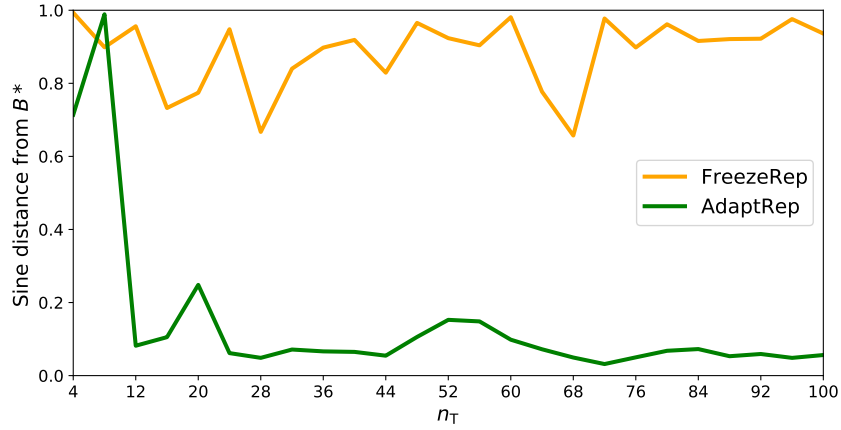


Figure 6: Sine distances of the representations learned by each method from the correct space B^* .

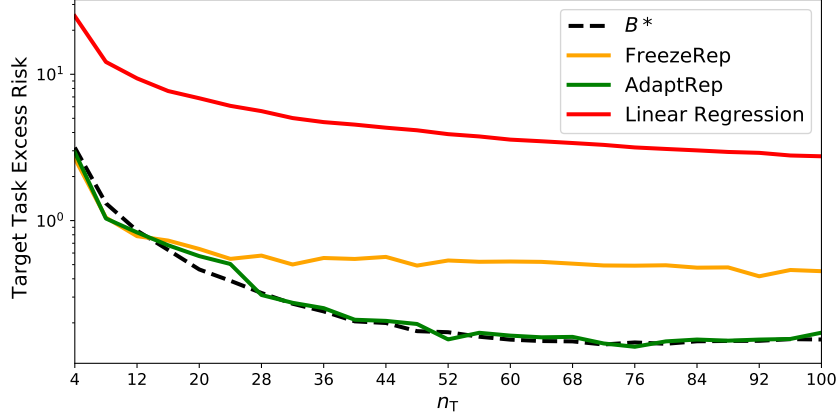


Figure 7: The worst-case average excess risk for several settings of n_T , plotted on a logarithmic y -axis.

We consider the high-dimensional setting where $d \asymp n_T$, and $k = \Theta(1)$. We fix $k = 2$ throughout the experiment. The matrix B^* spans the first k coordinates, while the residuals δ_t^* lie in the span of the last k coordinates. Finally, we set the Gaussian noise variance to $\sigma^2 = 1$.

During source training, both FROZENREP and ADAPTREP are provided with 1000 tasks and $10d$ samples per task from the task distribution in Section 3.4 (both are provided with the same samples). During target time, we evaluate the methods on the worst-case regression task from the same family.

Nonconvexity in Source Procedure. Rather than optimizing (4) or (6) during source training, we use an additional Frobenius-norm regularizer on $B^\top B - WW^\top$ to balance B and W . For FROZENREP, the regularized objective was shown to have a favorable optimization landscape in Tripuraneni et al. (2020a). We used L-BFGS to optimize these regularized objectives. To further mitigate optimization issues, we evaluated both methods with 10 random restarts, and report the best of the 10 restarts (as measured by the worst-case performance on the target task) for both methods.

(Subspace Alignment). We plotted the alignment of the learned representation (using the best of the 10 restarts) with B^* . We measure this via the sine of the largest principal angle between the two spaces, i.e.

$$\sin \Theta_1(\hat{B}, B^*) = \sqrt{1 - \sigma_1^2(\hat{B}^\top B^*)}.$$

We plot the results in Figure 6. As predicted by Lemma C.1, FROZENREP does not learn B^* , in contrast to ADAPTREP.

(Target Task Performance). We evaluated how the methods fare on their corresponding worst-case target tasks. We do so by training with the representation over 1000 i.i.d. draws of the target dataset, and averaging the excess risk over all obtained representations. We provide the results in Figure 7, and include a comparison with standard linear regression. As predicted, ADAPTREP performs much better than FROZENREP, with a gap that grows with n_T .

D Proof of Theorem 4.1

In this section, we will prove the guarantee provided in Theorem 4.1, along with all related intermediate results. Along these lines, we proceed as we did for Section B. More explicitly, the overall outline of the proof follows the four major steps below. We have placed in parenthesis the corresponding intermediate steps in the linear representation case (Section B) as an additional illustration of the procedure:

- (1) Provide a statistical rate for source training (Lemma B.3).
- (2) Bound the difference in performance between the solution found by the optimization algorithm and the ERM solution (Lemma B.6+Theorem H.1).
- (3) Prove uniform concentration over $\mathcal{A}_{C_T}(\theta_0)$ as a function of target sample size (Lemma B.5).
- (4) Connect the best-case performance in $\mathcal{A}_{C_T}(\theta_0)$ to the performance of the learner on the source tasks (Lemmas B.4, B.5 and B.8).

Note that step (4) is provided by the (ν, ε) -diversity condition. We will now proceed to demonstrate the remaining steps.

D.1 Bounding Average Source Task Performance (1)

The (ν, ε) -diversity condition implies that we can bound the best-case performance during target time training via control over the average source task performance. We proceed to provide such a bound using a standard uniform convergence argument. Recall that for a set of vector-valued functions \mathcal{H} mapping from $\mathbb{R}^m \rightarrow \mathbb{R}^n$, the Rademacher complexity of \mathcal{H} on n_S samples, denoted $\mathcal{R}_{n_S}(\mathcal{H})$, is given by

$$\mathcal{R}_{n_S}(\mathcal{H}) := \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{1}{n_S} \sum_{i \in [n_S]} \sum_{j \in [n]} \varepsilon_{ij} h(x_i)_j \right],$$

where the expectation is over the samples (x_i) and i.i.d. Rademacher random variables ε_{ij} .

Lemma D.1 (Source Training Bound). *Let $\theta \in \Theta_0$ be a minimizer of the training objective in (7). Then, with probability at least $1 - \delta$ over the random draw of inputs and labels,*

$$\frac{1}{T} \sum_{t \in [T]} \inf_{g_t \in \mathcal{A}_{C_S}(\theta)} \mathcal{L}_\infty^{\text{ex}}(g_t, g_t^*) \lesssim \frac{1}{\delta T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{C_S}(\theta)]^{\otimes T} \right] + \frac{B}{\delta \sqrt{n_S T}}.$$

Proof. For any $t \in [T]$, we define the empirical risk minimizer $\bar{g}_t := \operatorname{argmin}_{g \in \mathcal{A}_{C_S}(\theta)} \mathcal{L}(g, g_t^*)$. Then, since θ minimizes the training objective in (7),

$$\sum_{t \in [T]} \mathcal{L}(\bar{g}_t, g_t^*) - \mathcal{L}(g_t^*, g_t^*) \leq 0.$$

Following the canonical risk decomposition, we have that

$$\begin{aligned} \sum_{t \in [T]} \mathcal{L}_\infty^{\text{ex}}(\bar{g}_t, g_t^*) &\leq \sum_{t \in [T]} \mathcal{L}_\infty(\bar{g}_t, g_t^*) - \mathcal{L}(\bar{g}_t, g_t^*) + \sum_{t \in [T]} \mathcal{L}(\bar{g}_t, g_t^*) - \mathcal{L}(g_t^*, g_t^*) \\ &\quad + \sum_{t \in [T]} \mathcal{L}(g_t^*, g_t^*) - \mathcal{L}_\infty(g_t^*, g_t^*) \\ &\leq 2 \sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{C_S}(\theta)}} \left| \sum_{t \in [T]} \mathcal{L}_\infty(g_t, g_t^*) - \mathcal{L}(g_t, g_t^*) \right|. \end{aligned}$$

We give a high-probability bound on the right-hand side by first bounding its expectation, and then applying Markov's inequality to obtain the desired result. To bound the expectation, note that via

symmetrization,

$$\begin{aligned} & \mathbb{E} \left[\sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \sum_{t \in [T]} \mathcal{L}_\infty(g_t, g_t^*) - \mathcal{L}(g_t, g_t^*) \right| \right] \\ & \leq 2 \mathbb{E} \left[\sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \frac{1}{n_S} \sum_{t \in [T]} \sum_{i \in [n_S]} \varepsilon_{ij} \ell(g_t(x_{i,t}), y_{i,t}) \right| \right] \end{aligned}$$

Recentring about $\ell(0, y_{i,t})$ and using the constant-shift property of the Rademacher complexity,

$$\begin{aligned} & \mathbb{E} \left[\sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \frac{1}{n_S} \sum_{t \in [T]} \sum_{i \in [n_S]} \varepsilon_{ij} \ell(g_t(x_{i,t}), y_{i,t}) \right| \right] \\ & \leq \mathbb{E} \left[\sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \frac{1}{n_S} \sum_{t \in [T]} \sum_{i \in [n_S]} \varepsilon_{ij} [\ell(g_t(x_{i,t}), y_{i,t}) - \ell(0, y_{i,t})] \right| \right] + B \sqrt{\frac{T}{n_S}} \end{aligned}$$

Finally, since the loss is 1-Lipschitz, we can apply the Rademacher contraction principle, from which we find that

$$\begin{aligned} & \mathbb{E} \left[\sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \frac{1}{n_S} \sum_{t \in [T]} \sum_{i \in [n_S]} \varepsilon_{ij} [\ell(g_t(x_{i,t}), y_{i,t}) - \ell(0, y_{i,t})] \right| \right] \\ & \lesssim \mathbb{E} \left[\sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \frac{1}{n_S} \sum_{t \in [T]} \sum_{i \in [n_S]} \varepsilon_{ij} g_t(x_{i,t}) \right| \right] = \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right]. \end{aligned}$$

Putting everything together, we can thus bound the expectation of the maximum deviation as

$$\mathbb{E} \left[\sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \sum_{t \in [T]} \mathcal{L}_\infty(g_t, g_t^*) - \mathcal{L}(g_t, g_t^*) \right| \right] \lesssim \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] + B \sqrt{\frac{T}{n_S}}.$$

Therefore, by applying Markov's inequality, we have that with probability at least $1 - \delta$,

$$\begin{aligned} \frac{1}{T} \sum_{t \in [T]} \mathcal{L}_\infty^{\text{ex}}(\bar{g}_t, g_t^*) & \lesssim \frac{1}{T} \sup_{\substack{\theta \in \Theta_0 \\ g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)}} \left| \sum_{t \in [T]} \mathcal{L}_\infty(g_t, g_t^*) - \mathcal{L}(g_t, g_t^*) \right| \\ & \lesssim \frac{1}{\delta T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] + \frac{B}{\delta \sqrt{n_S T}}. \quad \square \end{aligned}$$

D.2 Bounding Optimization Performance (2)

Having provided a bound on source task performance, we now proceed to analyze the objective being optimized during target time training. We first note that the approximate linearity assumption in Assumption 4.4 and the norm-boundedness of \mathcal{C}_T via Assumption 4.5 ensure that the results from Section H apply, as long as we can show that the empirical loss satisfies an $(\alpha/\sqrt{n_T})$ -Lipschitz condition. We proceed to show that this is indeed a simple consequence of the 1-Lipschitz assumption on the loss given by Assumption 4.2.

Lemma D.2. *Define the function $\mathcal{L} : \mathbb{R}^{n_T} \rightarrow \mathbb{R}$ as*

$$\mathcal{L}(\hat{y}) := \frac{1}{n_T} \sum_{i \in [n_T]} \ell(\hat{y}_i, y_i).$$

for fixed $y_1, \dots, y_{n_T} \in \mathcal{Y}$. Then, for any \hat{y} , $\|\nabla_{\hat{y}} \mathcal{L}(\hat{y})\|_2^2 \leq 1/n_T$.

Proof. By direct computation,

$$\|\nabla_{\hat{y}} \mathcal{L}(\hat{y})\|_2^2 = \frac{1}{n_T^2} \sum_{i \in n_T} |\nabla_{\hat{y}_i} \ell(\hat{y}_i, y)|^2 \leq \frac{1}{n_T},$$

where we have used the fact that $\ell(\cdot, y)$ is 1-Lipschitz for any $y \in \mathcal{Y}$ by Assumption 4.2. \square

Having verified that the assumptions in Section H hold, it follows that we have the following performance bound on projected gradient descent during target time training:

Lemma D.3. *Assume that we run projected gradient descent during target training time for T_{PGD} iterations with step size η given by*

$$\eta = \frac{1}{\sqrt{T_{\text{PGD}}}} \left(\frac{R}{\sqrt{L^2 + \beta^2 R^2}} \right).$$

Let $g_0, \dots, g_{T_{\text{PGD}}}$ denote the sequence of predictors obtained, where $g_0 = g_{\theta_0}$. Then, for any $g \in \mathcal{A}_{\mathcal{C}_T}(\theta_0)$,

$$\min_t \mathcal{L}(g_t, g_t^*) - \mathcal{L}(g, g_t^*) \leq \beta R^2 + R \sqrt{\frac{L^2 + \beta^2 R^2}{T_{\text{PGD}}}}.$$

D.3 Bounding the Performance of ERM during Target Training (3)

We now proceed to bound the performance of the ERM solution during target time training. Via following the canonical risk decomposition as in Lemma D.1, we can prove such a bound simply by bounding the maximum deviation between empirical and population losses over $\mathcal{A}_{\mathcal{C}_T}(\theta_0)$.

Lemma D.4. *Let \mathcal{S} be the support of ρ . With probability at least $1 - \delta$ over the random draw of inputs and noise,*

$$\sup_{\substack{g^* \in \mathcal{S} \\ g \in \mathcal{A}_{\mathcal{C}_T}(\theta_0)}} |\mathcal{L}(g, g^*) - \mathcal{L}_\infty(g, g^*)| \leq \frac{1}{\delta} \sup_{\theta \in \Theta_0} \mathcal{R}_{n_T}[\mathcal{A}_{\mathcal{C}_T}(\theta)] + \frac{B}{\delta \sqrt{n_T}}.$$

Proof. The proof proceeds similarly to that of Lemma D.1. Note that the supremum over g^* does not affect the bound, since g^* only enters into the expression through the labels y_i , and no matter what choice of g^* is made, $|\ell(0, y_i)| \leq B$ for all i . The final result follows from taking a supremum over all possible initialization choices. \square

D.4 Concluding: Proving Theorem 4.1

Having completed all the steps for the outline, we now proceed to compile the main result. Intuitively, since the diversity condition allows us to bound the performance of the best predictor, and projected gradient descent can perform as well as any predictor in $\mathcal{A}_{\mathcal{C}_T}(\theta_0)$ (and thus, as well as the best predictor), we can obtain performance bounds on the iterates found while training on the target task.

Theorem 4.1 (General Performance Bound). *Assume that all assumptions in Section 4.2 hold. Let (θ_t) be the set of iterates generated by PGD following the procedure in Section 4.1. Then, with probability at least $1 - \delta$ over the random draw of samples,*

$$\begin{aligned} \mathbb{E}_{g^* \sim \rho} \left[\min_t \mathcal{L}_\infty^{\text{ex}}(g_{\theta_t}, g^*) \right] &\lesssim \underbrace{\beta R^2 + R \sqrt{\frac{L^2 + \beta^2 R^2}{T_{\text{PGD}}}}}_{\varepsilon_{\text{OPT}}} + \underbrace{\frac{1}{\delta} \sup_{\theta \in \Theta_0} \mathcal{R}_{n_T}[\mathcal{A}_{\mathcal{C}_T}(\theta)] + \frac{B}{\delta \sqrt{n_T}}}_{\varepsilon_{\text{EST}}} \\ &\quad + \underbrace{\frac{1}{\nu} \left\{ \frac{1}{\delta T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] + \frac{B}{\delta \sqrt{n_S T}} \right\}}_{\varepsilon_{\text{REPR}}} + \varepsilon. \end{aligned}$$

Note that the \mathcal{R}_{n_T} complexity term samples from p . Meanwhile, the \mathcal{R}_{n_S} complexity term samples from $p^{\otimes T}$, which concatenates T i.i.d. samples from p (one for each task) for every draw.

Proof. We instantiate the high-probability events in Lemmas D.1 and D.4 both with failure probability $\delta/2$, which simultaneously occur with probability at least $1 - \delta$ via a union-bound. Consequently, we have the bounds

$$\begin{aligned} \frac{1}{T} \sum_{t \in [T]} \inf_{g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)} \mathcal{L}_\infty^{\text{ex}}(g_t, g_t^*) &\lesssim \frac{1}{\delta T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] + \frac{B}{\delta \sqrt{n_S T}} \\ \sup_{\substack{g^* \in \mathcal{S} \\ g \in \mathcal{A}_{\mathcal{C}_T}(\theta_0)}} |\mathcal{L}(g, g^*) - \mathcal{L}_\infty(g, g^*)| &\lesssim \frac{1}{\delta} \sup_{\theta \in \Theta_0} \mathcal{R}_{n_T} [\mathcal{A}_{\mathcal{C}_T}(\theta)] + \frac{B}{\delta \sqrt{n_T}} \end{aligned}$$

Given g^* , let \bar{g} be the population risk minimizer in $\mathcal{A}_{\mathcal{C}_T}(\theta_0)$. Throughout this proof, we omit g^* when writing \mathcal{L} and \mathcal{L}_∞ whenever it is understood. Then, by Lemma D.3, projected gradient descent always generates iterates satisfying

$$\min_t \mathcal{L}(g_t) - \mathcal{L}(\bar{g}) \lesssim \beta R^2 + R \sqrt{\frac{L^2 + \beta^2 R^2}{T_{\text{PGD}}}}, \quad (12)$$

independent of g^* . Let \tilde{g} be the predictor achieving the minimum. We thus proceed to decompose the risk as

$$\begin{aligned} \min_t \mathcal{L}_\infty(g_t) - \mathcal{L}_\infty(g^*) &\lesssim \underbrace{\mathcal{L}_\infty(\tilde{g}) - \mathcal{L}(\tilde{g})}_{=: T_1} + \underbrace{\mathcal{L}(\tilde{g}) - \mathcal{L}(\bar{g})}_{=: T_2} \\ &\quad + \underbrace{\mathcal{L}(\bar{g}) - \mathcal{L}_\infty(\bar{g})}_{=: T_3} + \underbrace{\mathcal{L}_\infty(\bar{g}) - \mathcal{L}_\infty(g^*)}_{=: T_4}. \end{aligned}$$

We now bound the terms T_1, \dots, T_4 individually. As a result of uniform convergence as guaranteed by Lemma D.4,

$$T_1, T_3 \leq \sup_{\substack{g^* \in \mathcal{S} \\ g \in \mathcal{A}_{\mathcal{C}_T}(\theta_0)}} |\mathcal{L}(g) - \mathcal{L}_\infty(g)| \lesssim \frac{1}{\delta} \mathcal{R}_{n_T} [\mathcal{A}_{\mathcal{C}_T}(\theta_0)] + \frac{B}{\delta \sqrt{n_T}},$$

where \mathcal{S} is the support of the target task distribution ρ . Furthermore, we have demonstrated a bound on T_2 in (12). Finally, by applying the (ν, ε) -diversity condition as guaranteed by Assumption 4.3,

$$\begin{aligned} \mathbb{E}_{g^* \sim \rho} [T_4] &= \mathbb{E}_{g^* \sim \rho} \left[\inf_{g \in \mathcal{A}_{\mathcal{C}_T}(\theta_0)} \mathcal{L}_\infty^{\text{ex}}(g) \right] \\ &\leq \frac{1}{\nu} \left[\frac{1}{T} \sum_{t \in [T]} \inf_{g_t \in \mathcal{A}_{\mathcal{C}_S}(\theta)} \mathcal{L}_\infty^{\text{ex}}(g_t) \right] + \varepsilon \\ &\leq \frac{1}{\nu} \left\{ \frac{1}{\delta T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] + \frac{B}{\delta \sqrt{n_S T}} \right\} + \varepsilon. \end{aligned}$$

Taking the expectation with respect to $g^* \sim \rho$ and putting all of these bounds together, we thus have that

$$\begin{aligned} \mathbb{E}_{g^* \sim \rho} [\mathcal{L}_\infty^{\text{ex}}(g, g^*)] &\lesssim \beta R^2 + R \sqrt{\frac{L^2 + \beta^2 R^2}{T_{\text{PGD}}}} + \frac{1}{\delta} \sup_{\theta \in \Theta_0} \mathcal{R}_{n_T} [\mathcal{A}_{\mathcal{C}_T}(\theta)] + \frac{B}{\delta \sqrt{n_T}} \\ &\quad + \frac{1}{\nu} \left\{ \frac{1}{\delta T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] + \frac{B}{\delta \sqrt{n_S T}} \right\} + \varepsilon. \quad \square \end{aligned}$$

E Case Study: Two-Layer Neural Networks (Proofs)

E.1 Additional Assumptions on Initialization and Diversity

Before we proceed with the proofs for Section 5, we first formalize the additional assumptions that were referenced in the main text.

Assumption E.1 (Initialization Assumptions). *The initialization (B^*, w_0^*) satisfies the following assumptions:*

- (a) $B^* = [A^*, A^*]$ for $A^* \in \mathbb{R}^{d \times k}$. Furthermore, $w_0^* \in \{-1, 1\}^{2k}$, with $(w_0^*)_i = -(w_0^*)_{i+k}$ for $i \in [k]$.
- (b) The columns of B^* are norm-bounded by 1.
- (c) $c_1 I \preceq \mathbb{E} [\rho_{\theta_0^*}(x) \rho_{\theta_0^*}(x)^\top] \preceq c_2 I$ for $c_1 > 0$. We define $\kappa = c_2/c_1$.

We refer to initializations satisfying (a) as *antisymmetric initializations*. Note that under any antisymmetric initialization θ_0 , $f_{\theta_0}^\beta(x) \equiv 0$. Additionally, the assumption on the representation covariance in (c) ensures that that representation is well-conditioned, with condition number κ .

Finally, we need to impose diversity conditions on the source tasks. Recall that for every $t \in [T]$, there exists unit-norm $w_t^* \in \mathbb{R}^{2k}$ and $\delta_t^* \in \mathbb{R}^k$ such that the corresponding source predictor is parameterized by

$$\theta_t^* = \left(B^* + \frac{1}{\beta} \sum_{i \in [k]} (\delta_t^*)_i \Delta_i^*, w_0^* + \frac{1}{\beta} w_t^* \right).$$

We note that the fine-tuning step for $t \in [T]$ can be parametrized by $\omega_t^* = [w_t^*, \delta_t^*] \in \mathbb{R}^{3k}$ via a (natural) linear transformation. We assume that this parametrization has a matrix representation Γ/β , so that $\theta_t^* = \theta_0^* + \Gamma \omega_t^*/\beta$. Then, we impose the following condition on $(\omega_t^*)_t$:

Assumption E.2. Let $\Omega \in \mathbb{R}^{2k \times T}$ be the matrix $[\omega_1^*, \dots, \omega_T^*]$. Then, $\sigma_{2k}(\Omega) \gtrsim T/k$.

The assumption above is analogous to the diversity conditions assumed in the previous sections.

E.2 Formalizing Approximate Linearity

In this section, we demonstrate that under Assumption E.1, $f_{\theta_t^*}^\beta$ behaves approximately like its linearization for large enough β .

Recall that the linearization of f_θ^β involves two feature vectors, $\phi_{B_0}(x)$ and $\psi_{B_0, w_0}(x)$. We now provide the following formal definitions for these vectors:

Definition E.1 (Feature vectors ϕ, ψ, ρ). Let $\theta_0 = (B_0, w_0)$ be an antisymmetric parameter, i.e. satisfying Assumption E.1(c). Then, for every x , there exists feature vectors $\phi_{B_0}(x)$ and $\psi_{B_0, w_0}(x)$ such that

$$f_{\theta_0 + (\Delta, w)}^\beta(x) = \beta w^\top \phi_{B_0}(x) + \beta \Delta^\top \psi_{B_0, w_0}(x) + \beta \zeta_{B_0, w_0}^{\Delta, w}(x).$$

We interpret the features $\phi_{B_0}(x)$ and $\psi_{B_0, w_0}(x)$ to be the gradients of $(B, w) \mapsto f_{(B, w)}^\beta(x)$ evaluated at θ_0 , which have closed forms

$$\phi_{B_0}(x) = \sigma(B_0^\top x) \quad \text{and} \quad \psi_{B_0, w_0}(x) = w_0^{\text{diag}} \sigma'(B_0^\top x) x^\top.$$

Additionally, $\zeta_{B_0, w_0}^{\Delta, w}(x)$ is the Taylor error. By Taylor's theorem, there exists $\bar{B} = B_0 + \alpha_1 \Delta$ and $\bar{w} = w_0 + \alpha_2 w$ for some $\alpha_1, \alpha_2 \in [0, 1]$ such that $\zeta_{(B_0, w_0)}^{\Delta, w}(x)$ can be written as

$$\zeta_{(B_0, w_0)}^{\Delta, w}(x) = \sum_{i=1}^{2k} \bar{w}_i \sigma''(\bar{B}_i^\top x) (\Delta_i^\top x)^2 + 2w_i \sigma'(\bar{B}_i^\top x) (\Delta_i^\top x).$$

Finally, we define $\rho_{B_0, w_0}(x)$ to be the concatenation of $\phi_{B_0}(x)$ and $\psi_{B_0, w_0}(x)$. ◇

For $\theta_0 = (B_0, w_0)$ and $\delta = (\Delta, w)$, we will frequently abuse notation and let $\rho_{\theta_0}(x)\delta$ denote the linearization of $f_{\theta_0+\delta}$, i.e.

$$\rho_{\theta_0}(x)\delta = w^\top \phi_{B_0}(x) + \langle \psi_{\theta_0}(x), \Delta \rangle.$$

We will show that if $\|\Delta\|_F$ and $\|w\|_F$ are both $O(1/\beta)$, then the remainder term $\zeta_{B_0, w_0}^{\Delta, w}(x)$ is $O(1/\beta)$, and thus the function class is approximately linear in (Δ, w) with feature functions $\phi_{B_0}(\cdot)$ and $\psi_{B_0, w_0}(\cdot)$. We first bound the Hessian term, in order to control the remainder term $\zeta_{(B_0, w_0)}^{\Delta, w}(x)$, which depends on the Hessian.

Lemma E.1 (Hessian Bound). *Fix a matrix $\bar{B} \in \mathbb{R}^{2k \times d}$ and a vector $\bar{w} \in \mathbb{R}^{2k}$. Assume that all rows of \bar{B} are 2-norm-bounded, and that $\|\bar{w}\|_2 \leq 2$. Then, for $x \sim p$ almost surely, $\left\| \nabla_{\theta}^2 f_{(\bar{B}, \bar{w})}^\beta(x) \right\|_2 \lesssim \beta(\mu + L)$.*

Proof. Throughout the proof, we will use the fact that $\|x\|_2 \leq 1$ for $x \sim p$ almost surely. By definition,

$$\left\| \nabla_{\theta}^2 f_{(\bar{B}, \bar{w})}^\beta(x) \right\|_2 = \sup_{\|\Delta\|_F^2 + \|w\|_2^2 \leq 1} \beta \left| \sum_{i \in [2k]} \bar{w}_i \sigma''(\bar{B}_i^\top x) \langle \Delta_i, x \rangle^2 + 2 \sum_{i \in [2k]} w_i \sigma'(\bar{B}_i^\top x) \langle \Delta_i, x \rangle \right|$$

To bound the first term, observe that by applying Hölder's inequality,

$$\begin{aligned} \sum_{i \in [2k]} \bar{w}_i \sigma''(\bar{B}_i^\top x) \langle \Delta_i, x \rangle^2 &\leq \left(\sup_{i \in [2k]} |w_i \sigma''(\bar{B}_i^\top x)| \right) \sum_{i \in [2k]} \langle \Delta_i, x \rangle^2 \\ &= \left(\sup_{i \in [2k]} |\bar{w}_i \sigma''(\bar{B}_i^\top x)| \right) \|\Delta x\|_2^2 \leq 2\mu. \end{aligned}$$

To bound the second term, we apply Cauchy-Schwarz, from which we see that

$$\begin{aligned} 2 \sum_{i \in [2k]} w_i \sigma'(\bar{B}_i^\top x) \langle \Delta_i, x \rangle &\leq 2 \left(\sum_{i \in [2k]} [w_i \sigma'(\bar{B}_i^\top x)]^2 \right)^{1/2} \|\Delta x\|_2 \\ &\leq 2L \|w\|_2 \|\Delta x\|_2 \leq 2L. \end{aligned}$$

Altogether, we thus have that almost surely for $x \sim p$,

$$\left\| \nabla_{\theta}^2 f_{(\bar{B}, \bar{w})}^\beta(x) \right\|_2 \lesssim \beta(\mu + L). \quad \square$$

As an immediate corollary, since $\zeta_{\theta_0}^\delta(x)$ evaluates the Hessian at a point satisfying the preconditions of Lemma E.1 during both source and target training time, we have the following result:

Corollary E.1. *For any $\theta_0 \in \Theta_0$ and any (Δ, w) with $\|\Delta\|_F \leq c_1/\beta$, $\|w\|_2 \leq c_2/\beta$, we have that $\left| \beta \zeta_{\theta_0}^{\Delta, w}(x) \right| \leq c_1 c_2 (\mu + L)/\beta$ for $x \sim p$ almost surely.*

Therefore, we see that for large enough β , the remainder term is close to 0, and thus the family is indeed approximately linear. Finally, we provide a norm bound on the combined representation vector $\rho_{\theta_0}(x)$.

Lemma E.2 (Representation Norm Bound). *For $x \sim p$, $\|\rho_{\theta_0}(x)\|_2 \leq 2L\sqrt{k}$ almost surely for any $\theta_0 = (B_0, w_0) \in \Theta_0$.*

Proof. Recall from definitions that $\|\rho_{\theta_0}(x)\|_2^2 = \|\phi_{B_0}(x)\|_2^2 + \|\psi_{\theta_0}(x)\|_F^2$. To bound the activation features, since $\sigma(0) = 0$ and $|\sigma'(x)| \leq L$ for $x \in [-1, 1]$,

$$\|\phi_{B_0}(x)\|_2^2 = 2k \max_{i \in [2k]} |\sigma(B_{0,i}^\top x)|^2 \leq 2k \max_{i \in [2k]} \left| \int_0^{B_{0,i}^\top x} \sigma'(z) dz \right|^2 \leq 2kL^2.$$

To bound the gradient features, we once again use the boundedness of σ' to obtain

$$\|\psi_{\theta_0}(x)\|_F^2 = \left\| w_0^{\text{diag}} \sigma'(B_0^\top x) x^\top \right\|_F^2 = \left\| w_0^{\text{diag}} \sigma'(B_0^\top x) \right\|_2^2 \|x\|_2^2 = \|\sigma'(B_0^\top x)\|_2^2 \|x\|_2^2 \leq 2kL^2.$$

Putting the bounds together, we obtain the desired overall bound. \square

E.3 Verifying Assumptions of Section 4.2

Having shown that the function class is approximately linear as well as the norm bounds above, we now can proceed to verify the required assumptions. First, we verify the assumptions on the loss function; clearly, the squared error loss is convex, so we simply need to verify that the loss is Lipschitz over the prediction, and that $\ell(0, y)$ is bounded. Note that we have to prove separate bounds for source and target training, due to the change in the function class.

Lemma E.3 (Validating Loss Assumptions). *Assume that $\gamma^2 \geq 2\kappa$. Define the quantities*

$$\alpha_1 := L\sqrt{k} + \frac{\mu + L}{\beta} \quad \text{and} \quad \alpha_2 := \kappa^{1/2}L\sqrt{k} + \frac{\kappa}{\gamma}(\mu + L).$$

Then, during both source and target training time, $|\ell(0, y)| \leq \alpha_1^2$ for any y , and $\ell(\cdot, y)$ is α_1 -Lipschitz during source training time and α_2 -Lipschitz during target training time for any y , all up to universal constants.

Proof. For the squared error loss, $\ell(0, y) = y^2$ and $\nabla_{\hat{y}}\ell(\hat{y}, y) = 2(\hat{y} - y)$ for any y . Therefore, since the additive noise is $O(1)$ bounded, the claims hold as long as we can prove a bound on the predictions of any feasible predictor.

We first consider bounding the Lipschitz constant during source-time training. That is, we need to provide a uniform bound on $x \mapsto \beta\rho_{\theta_0}(x)^\top\delta + \beta\zeta_{\theta_0}^\delta(x)$ for $\theta_0 \in \Theta_0$, $\delta \in \mathcal{C}_S$, and x in the support of p . To this end, we apply the bounds in Corollary E.1 and Lemma E.2, from which we obtain

$$|\beta\rho_{\theta_0}(x)^\top\delta + \beta\zeta_{\theta_0}^\delta(x)| \lesssim L\sqrt{k} + \frac{\mu + L}{\beta} =: \alpha_1.$$

By applying a similar argument for target-time training, we find that

$$|\gamma\rho_{\theta_0}(x)^\top\delta + \gamma\zeta_{\theta_0}^\delta(x)| \lesssim \kappa^{1/2}L\sqrt{k} + \frac{\kappa}{\gamma}(\mu + L) =: \alpha_2.$$

We thus can conclude that during both source and target training time, $|\ell(0, y)| \leq \alpha_1^2$ for any y , and that $\ell(\cdot, y)$ is α_1 -Lipschitz during source training time and α_2 -Lipschitz during target training time for any feasible y . \square

Lemma E.4 (Neural network diversity). *Assume that $\gamma > \max(\kappa\beta, \sqrt{2\kappa})$. Then, the source tasks satisfy a $(1, (L + \mu)^2/\beta^2)$ -diversity condition with respect to the target task distribution ρ .*

Proof. Observe that we can write the excess risk of a predictor $f_{(B_0, w_0)}^\gamma$ as

$$\mathbb{E} \left[[(\gamma\rho_{B_0, w_0}(x)\delta - \beta\rho_{B^*, w_0^*}(x)\delta^*) + (\gamma\zeta_{B_0, w_0}^\delta(x) - \beta\zeta_{B^*, w_0^*}^{\delta^*}(x))]^2 \right].$$

We proceed to upper bound the averaged best-case target performance, and lower bound the best-case performance averaged over source tasks.

Upper bounding the expected best-case target performance.

Fix a δ^* . Then, since $\gamma^2 > 2\kappa$ guarantees that Lemma E.1 holds,

$$\begin{aligned} \mathbb{E} \left[(f_{\theta_0+\delta}^\gamma(x) - f_{\theta_0^*+\delta^*}^\beta(x))^2 \right] &\lesssim \mathbb{E} \left[(\gamma\rho_{\theta_0}(x)\delta - \rho_{\theta_0^*}(x)\delta^*)^2 \right] + \mathbb{E} \left[(\gamma\zeta_{\theta_0}^\delta(x) - \beta\zeta_{\theta_0^*}^{\delta^*}(x))^2 \right] \\ &\lesssim \mathbb{E} \left[(\gamma\rho_{\theta_0}(x)\delta - \beta\rho_{\theta_0^*}(x)\delta^*)^2 \right] + \left[\frac{\kappa}{\gamma}(L + \mu) \right]^2 + \frac{(L + \mu)^2}{\beta^2} \\ &\lesssim \mathbb{E} \left[(\gamma\rho_{\theta_0}(x)\delta - \beta\rho_{\theta_0^*}(x)\delta^*)^2 \right] + \frac{(L + \mu)^2}{\beta^2}. \end{aligned}$$

Therefore, by taking the infimum over $\delta \in \mathcal{C}_T$ and noting that the unconstrained infimum of the right-hand side is feasible by Proposition I.7 and the conditioning constraint imposed by Θ_0 ,

$$\inf_{\delta \in \mathcal{C}_T} \mathbb{E} \left[(f_{\theta_0+\delta}^\gamma(x) - f_{\theta_0^*+\delta^*}^\beta(x))^2 \right] \leq \beta^2 \left\| \Lambda(\rho_{\theta_0}, \rho_{\theta_0^*})^{1/2} \delta^* \right\|_2^2 + \left(\frac{L + \mu}{\beta} \right)^2$$

Now, by the choice of target task distribution, we can write δ^* as $\Gamma\omega^*/\beta$, where we note that $\mathbb{E}[\omega^*(\omega^*)^\top] = (1/k)I_{2k}$. Thus, by taking expectations, we have that

$$\mathbb{E}_{\delta^* \sim \rho} \left[\inf_{\delta \in \mathcal{C}_T^\gamma} \mathbb{E} \left[(f_{\theta_0+\delta}^\gamma(x) - f_{\theta_0^*+\delta^*}^\beta(x))^2 \right] \right] \leq \frac{1}{k} \left\| \Lambda(\rho_{\theta_0}, \rho_{\theta_0^*})^{1/2} \Gamma \right\|_F^2 + \left(\frac{L+\mu}{\beta} \right)^2.$$

Lower bounding the source task-averaged best-case target performance.

Fix a source task $t \in [T]$. By Corollary I.1, we have the bound

$$\begin{aligned} \beta^2 \mathbb{E} \left[(\rho_{\theta_0}(x)\delta - \rho_{\theta_0^*}(x)\delta_t^*)^2 \right] &\lesssim \mathbb{E} \left[(f_{\theta_0+\delta}^\beta(x) - f_{\theta_t^*}^\beta(x))^2 \right] + \mathbb{E} \left[(\zeta_{\theta_0}^\delta(x) - \zeta_{\theta_t^*}^{\delta^*}(x))^2 \right] \\ &\lesssim \mathbb{E} \left[(f_{\theta_0+\delta}^\beta(x) - f_{\theta^*}^\beta(x))^2 \right] + \left(\frac{L+\mu}{\beta} \right)^2. \end{aligned}$$

Therefore, by taking the infimum on both sides and applying Proposition I.7,

$$\left\| \Lambda(\rho_{\theta_0}, \rho_{\theta^*})^{1/2} \Gamma \omega_t^* \right\|_2^2 \leq \inf_{\delta \in \mathcal{C}_S} \mathbb{E} \left[(f_{\theta_0+\delta}^\beta(x) - f_{\theta^*}^\beta(x))^2 \right] + \left(\frac{L+\mu}{\beta} \right)^2.$$

By averaging over the set of source tasks,

$$\frac{1}{T} \left\| \Lambda(\rho_{\theta_0}, \rho_{\theta^*})^{1/2} \Gamma \Omega \right\|_F^2 \leq \frac{1}{T} \sum_{t \in [T]} \inf_{\delta \in \mathcal{C}_S} \mathbb{E} \left[(f_{\theta_0+\delta}^\beta(x) - f_{\theta^*}^\beta(x))^2 \right] + \left(\frac{L+\mu}{\beta} \right)^2.$$

Finally, by applying Assumption E.2 to lower bound the left-hand side, we find that

$$\frac{1}{k} \left\| \Lambda(\rho_{\theta_0}, \rho_{\theta^*})^{1/2} \Gamma \right\|_F^2 \leq \frac{1}{T} \sum_{t \in [T]} \inf_{\delta \in \mathcal{C}_S} \mathbb{E} \left[(f_{\theta_0+\delta}^\beta(x) - f_{\theta^*}^\beta(x))^2 \right] + \left(\frac{L+\mu}{\beta} \right)^2.$$

Concluding.

Consequently, by putting the two parts together, we have that

$$\begin{aligned} \mathbb{E}_{\delta^* \sim \rho} \left[\inf_{\delta \in \mathcal{C}_T^\gamma} \mathbb{E} \left[(f_{\theta_0+\delta}^\gamma(x) - f_{\theta_0^*+\delta^*}^\beta(x))^2 \right] \right] \\ \leq \frac{1}{T} \sum_{t \in [T]} \inf_{\delta \in \mathcal{C}_S} \mathbb{E} \left[(f_{\theta_0+\delta}^\beta(x) - f_{\theta^*}^\beta(x))^2 \right] + \left(\frac{L+\mu}{\beta} \right)^2. \end{aligned} \quad \square$$

Lemma E.5 (Neural network approximate linearity). *Assume that $\gamma^2 \geq 2\kappa$. For any $\theta_0 \in \Theta$,*

$$\sup_{\delta \in \mathcal{C}_T} \frac{1}{n_T} \sum_{i \in [n_T]} \left\| \nabla_{\theta}^2 f_{\theta_0+\delta}^\gamma(x_i) \right\|_2^2 \lesssim \gamma^2(L+\mu)^2 \quad \text{and} \quad \frac{1}{n_T} \sum_{i \in [n_T]} \left\| \nabla_{\theta} f_{\theta_0}^\gamma(x_i) \right\|_2^2 \lesssim \gamma^2 k L^2.$$

Proof. The first bound is a trivial consequence of Lemma E.1. Meanwhile, we note that the second inequality, the quantity inside the norm is exactly $\beta\rho_{\theta_0}(x)$, and thus the bound follows by Lemma E.2. \square

E.4 Computations

Lemma E.6 (Source Rademacher Bound). *We have the following bound on the Rademacher complexity:*

$$\frac{1}{T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] \lesssim \frac{kL}{(n_S T)^2} + \left[(L+\mu) \frac{k}{\sqrt{n_S T}} + L \sqrt{\frac{k}{n_S}} \right] \log(n_S T) + \frac{\mu+L}{\beta}$$

Proof. We first consider bounding the empirical Rademacher complexity for a fixed set of inputs $(x_{i,t})$, after which we obtain the desired result by taking expectations over the sampled inputs. Expanding the definition, we have that

$$\begin{aligned}
& \frac{1}{T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] \\
&= \frac{1}{nT} \mathbb{E} \left[\sup_{\substack{(B_0, w_0) \in \Theta_0 \\ (\Delta_t, w_t) \in \mathcal{C}_S}} \left| \sum_{\substack{i \in [n_S] \\ t \in [T]}} \varepsilon_{i,t} \left[w_t^\top \phi(B^\top x_{i,t}) + \langle \Delta_t, \psi_B(x_{i,t}) \rangle + \zeta_{B, w_0}^{\Delta_t, w_t}(x_{i,t}) \right] \right| \right] \\
&\leq \frac{1}{T} \mathcal{R}_{n_S}(\mathcal{F}_\Phi^{\otimes T} \circ \Phi) + \frac{1}{T} \mathcal{R}_{n_S}(\mathcal{F}_\Psi^{\otimes T} \circ \Psi) + \frac{\mu + L}{\beta},
\end{aligned}$$

where

$$\begin{aligned}
\mathcal{F}_\Phi &:= \{z \mapsto \beta \langle w, z \rangle \mid \|w\|_2 \leq 1/\beta\} \\
\mathcal{F}_\Psi &:= \{Z \mapsto \beta \langle \Delta, Z \rangle \mid \|\Delta\|_F \leq 1/\beta\} \\
\Phi &:= \{\phi_B \mid (\exists w) (B, w) \in \Theta_0\} \\
\Psi &:= \{\psi_{B,w} \mid (B, w) \in \Theta_0\}.
\end{aligned}$$

We proceed to bound the two complexity terms by noting that for any function class \mathcal{F} , $\mathcal{R}_{n_S}(\mathcal{F}) \lesssim \mathcal{G}_{n_S}(\mathcal{F})$, and making use of the Gaussian complexity chain rule from Proposition I.6. First, note that for a fixed set of latent vectors z_1, \dots, z_{n_S} ,

$$\mathcal{G}_Z(\mathcal{F}_\Phi) = \frac{1}{n_S} \mathbb{E} \left[\sup_{\|w\|_2 \leq 1} |\langle z, Zw \rangle| \right] \leq \frac{1}{n_S} \sqrt{\text{tr } ZZ^\top} = \frac{1}{n_S} \sqrt{\sum_{i \in [n_S]} \|z_i\|_2^2}.$$

Then, applying the representation norm bound Lemma E.2 to the latents,

$$\mathbb{E} \left[\max_{Z \in \mathcal{Z}} \mathcal{G}_Z(\mathcal{F}_\Phi) \right] \leq L \sqrt{\frac{k}{n_S}}.$$

Similarly, for a fixed set of latent matrices Z_1, \dots, Z_{n_S}

$$\mathcal{G}_Z(\mathcal{F}_\Psi) \leq \frac{1}{n_S} \sqrt{\sum_{i \in [n_S]} \|Z_i\|_F^2} \implies \mathbb{E} \left[\max_{Z \in \mathcal{Z}} \mathcal{G}_Z(\mathcal{F}_\Psi) \right] \leq L \sqrt{\frac{k}{n_S}},$$

where we once again used the representation bound from Lemma E.2.

Finally, we bound the complexity of the classes Φ and Ψ . Firstly, we have

$$\begin{aligned}
\frac{1}{T} \mathcal{G}_{n_S}(\Phi) &= \frac{1}{n_S T} \mathbb{E} \left[\sup_{(B_0, w_0) \in \Theta_0} \left| \sum_{r \in [2k]} \sum_{i \in [n_S]} \sum_{t \in [T]} z_{r,i,t} \sigma(B_r^\top x_{i,t}) \right| \right] \\
&\leq \frac{2k}{n_S T} \mathbb{E} \left[\sup_{\|b\|_2 \leq 1} \left| \sum_{i \in [n_S]} \sum_{t \in [T]} z_{i,t} \sigma(b^\top x_{i,t}) \right| \right] \\
&\lesssim \frac{kL}{n_S T} \mathbb{E} \left[\left\| \sum_{i \in [n_S T]} z_i x_i \right\|_2 \right] \leq \frac{kL}{n_S T} \sqrt{\sum_{i \in [n_S T]} \|x_i\|_2^2} \leq \frac{kL}{\sqrt{n_S T}}.
\end{aligned}$$

where the second inequality makes use of the Gaussian contraction result in Ledoux & Talagrand (1991, Corollary 3.17). Note that we reindex at the third line, which is made possible by the fact

that the source tasks have the same input distributions. Secondly,

$$\begin{aligned}
\frac{1}{T} \mathcal{G}_{n_S}(\Psi) &= \frac{1}{n_S T} \mathbb{E} \left[\sup_{(B_0, w_0) \in \Theta_0} \left| \sum_{r \in [2k]} \sum_{i \in [n_S]} \sum_{t \in [T]} \sigma'(B_r^\top x_{i,t}) \langle x_{i,t}, z_{r,i,t} \rangle \right| \right] \\
&= \frac{1}{n_S T} \mathbb{E} \left[\sup_{(B_0, w_0) \in \Theta_0} \left| \sum_{r \in [2k]} \sum_{i \in [n_S]} \sum_{t \in [T]} [\sigma'(B_r^\top x_{i,t}) - \sigma'(0)] \langle x_{i,t}, z_{r,i,t} \rangle \right| \right] + \frac{L}{\sqrt{n_S T}} \\
&= \frac{2k}{n_S T} \mathbb{E} \left[\sup_{\|b\|_2 \leq 1} \left| \sum_{i \in [n_S]} \sum_{t \in [T]} [\sigma'(b^\top x_{i,t}) - \sigma'(0)] \langle x_{i,t}, z_{r,i,t} \rangle \right| \right] + \frac{L}{\sqrt{n_S T}} \\
&\lesssim \frac{k\mu}{n_S T} \mathbb{E} \left[\left\| \sum_{i \in [n_S T]} \langle x_{i,t}, z_{i,t} \rangle x_{i,t} \right\|_2 \right] + \frac{L}{\sqrt{n_S T}} \leq \frac{k\mu}{n_S T} \sqrt{\sum_{i \in [n_S T]} \|x_i\|_2^4} + \frac{L}{\sqrt{n_S T}} \\
&\leq \frac{k\mu}{\sqrt{n_S T}} + \frac{L}{\sqrt{n_S T}},
\end{aligned}$$

following the same reasoning as before. Therefore, by applying Proposition I.6, we obtain the bound

$$\frac{1}{T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{C_S}(\theta)]^{\otimes T} \right] \lesssim \frac{kL}{(n_S T)^2} + \left[(L + \mu) \frac{k}{\sqrt{n_S T}} + L \sqrt{\frac{k}{n_S}} \right] \log(n_S T) + \frac{\mu + L}{\beta}. \quad \square$$

Lemma E.7 (Target Rademacher bound). *We have the following Rademacher complexity bound:*

$$\sup_{\theta \in \Theta_0} \mathcal{R}_{n_T}[\mathcal{A}_{C_T}^\gamma(\theta)] \lesssim L\sqrt{\kappa} \sqrt{\frac{k}{n_T}} + \frac{\kappa}{\gamma} (L + \mu).$$

Proof. Following definitions, we have that

$$\begin{aligned}
&\sup_{\theta \in \Theta_0} \mathcal{R}_{n_T}[\mathcal{A}_{C_T}^\gamma(\theta)] \\
&= \frac{\gamma}{n_T} \mathbb{E} \left[\sup_{\delta \in \mathcal{C}_T^\gamma} |\langle \varepsilon, \rho(X)\delta + \zeta_{\theta_0}^\delta(X) \rangle| \right] \leq \frac{\gamma}{n_T} \mathbb{E} \left[\sup_{\delta \in \mathcal{C}_T^\gamma} |\langle \varepsilon, \rho(X)\delta \rangle| \right] + \frac{\kappa}{\gamma} (L + \mu) \\
&= \frac{\sqrt{\kappa}}{n_T} \mathbb{E} [\|\rho(X)^\top \varepsilon\|_2] + \frac{\kappa}{\gamma} (L + \mu) \leq \frac{\sqrt{\kappa}}{n_T} \sqrt{\mathbb{E} [\|\rho(X)^\top \varepsilon\|_2^2]} + \frac{\kappa}{\gamma} (L + \mu) \\
&\leq L\sqrt{\kappa} \sqrt{\frac{k}{n_T}} + \frac{\kappa}{\gamma} (L + \mu),
\end{aligned}$$

where the final inequality uses the bound on the representation given in Lemma E.2. \square

E.5 Compiling the Bound

Theorem 5.1 (Neural net performance bound). *Assume that Assumptions 5.1, E.1 and E.2 hold. Then, if $n_S \geq n_T$, there exists a setting of the training parameters (see Section E) such that with probability at least $1 - \delta$, the iterates (θ_t) satisfy*

$$\mathbb{E}_{f^* \sim \rho} \left[\min_t \mathcal{L}_\infty^{\text{ex}}(f_{\theta_t}^\gamma, f^*) \right] \lesssim \kappa L^2 \frac{k}{\sqrt{n_T}} + L(L + \mu) \frac{k^{3/2}}{\sqrt{n_S T}} + \left(\frac{\mu + L}{\beta} \right) L\sqrt{k}.$$

Proof. By incorporating previous bounds and invoking Theorem 4.1, we have that for large enough γ ,

$$\begin{aligned}
&\mathbb{E}_{f^* \sim \rho} \left[\min_t \mathcal{L}_\infty^{\text{ex}}(f_{\theta_t}^\gamma, f^*) \right] \\
&\lesssim \sqrt{\frac{k\kappa(L + \mu)^2}{T_{\text{PGD}}}} + \kappa L^2 \frac{k}{\sqrt{n_T}} + L(L + \mu) \frac{k^{3/2}}{\sqrt{n_S T}} + \left(\frac{\mu + L}{\beta} \right) L\sqrt{k}.
\end{aligned}$$

Therefore, by running enough projected gradient descent iterations so that the first term matches $kL^2 k / \sqrt{n_T}$, we obtain the desired bound. \square

F Case Study: Logistic Regression

To further illustrate our framework, we analyze the performance of ADAPTRP on logistic regression, as done by Tripuraneni et al. (2020b). In this setting, we let $\theta = (B, w)$ for $B \in \mathbb{R}^{d \times k}$ and $w \in \mathbb{R}^k$, and define the predictor corresponding to θ to be $g_\theta(x) = x^\top Bw$.

F.1 Statistical Assumptions

As in the linear setting, we consider an input distribution p with covariance Σ . We restrict the set of labels \mathcal{Y} to $\{0, 1\}$, and consider the conditional distribution $q(y | g_\theta(x)) = \text{Ber}(\sigma(g_\theta(x)))$, where $\sigma(y) = 1/(1 + e^{-y})$ is the sigmoid function.

We define the optimal parameters for tasks $t \in [T]$ to be $(B^* + \Delta_t^*, w_t^*)$, where B^* is orthogonal and $\|\Delta_t^*\|_F \leq \delta_0$. As before, we define $\delta_t^* := \Delta_t^* w_t^*$ for any $t \in [T]$ and $W^* = [w_1, \dots, w_T] \in \mathbb{R}^{k \times T}$. Having defined the prior quantities, we make use of the statistical assumptions presented in Section 3.1, reproduced below for convenience:

Assumption F.1 (Sub-Gaussian input). *There exists $\rho > 0$ such that if $x \sim p_t$, then $\Sigma^{-1/2}x$ is ρ^2 -sub-Gaussian.*

Assumption F.2 (Source task diversity). *For any $t \in [T]$, $\|w_t^*\|_2 \leq r$, and $\sigma_k^2(W^*) = \Omega(r^2 T/k)$.*

Finally, we define the target task distribution ρ by sampling w^*, δ^* uniformly from the r - and δ_0 -balls of \mathbb{R}^k , respectively, and letting $\theta^* = B^* w^* + \delta^*$.

F.2 Training Procedure

We use the standard logistic loss $\ell(\hat{y}, y) = -y \log[\sigma(\hat{y})] - (1 - y) \log[1 - \sigma(\hat{y})]$. During source training, we optimize over initializations $\Theta_0 := \{(B, 0) \mid B \text{ is orthogonal}\}$. Let $B_0 \in \mathbb{R}^{d \times k}$ be the obtained representation. To adapt to the target task, we initialize the learner at $\theta_0 = ([B_0, B_0], [w_0, -w_0])$ for some unit-norm vector w_0 , and scale the predictor by a fixed parameter β to be chosen later. Finally, we set the feasible sets for optimization to be

$$\mathcal{C}_S := \{(\Delta, w) \mid \|\Delta\|_F \leq \delta_0, \|w\|_2 \leq r\} \quad \text{and} \quad \mathcal{C}_T^\beta := \left\{(\Delta, w) \mid \|\Delta\|_F \leq \frac{\delta_0}{\beta}, \|w\|_2 \leq \frac{r\kappa^{1/2}}{\beta}\right\},$$

where $\kappa := \lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma)$. Note the similarity of this procedure to that of the linear setting.

F.3 Performance Guarantee

Having described the statistical assumptions and the training procedure, we now specialize the guarantee of Theorem 4.1 to this setting.

Theorem F.1 (Performance Guarantee for Logistic Regression). *Assume that Assumption F.1 and Assumption F.2 both hold. Set the parameters for target time training to be*

$$\beta = (\kappa r^2 + \delta_0^2) \max\left(1, \frac{\sqrt{\text{tr } \Sigma}}{1/\delta \sqrt{n_T}}\right) \quad \text{and} \quad T_{\text{PGD}} = \frac{(\kappa r^2 + \delta_0^2) \text{tr } \Sigma}{1/\delta^2 n_T}.$$

Then, for $n_S, n_T \gtrsim \rho^4 d$, we have that with probability at least $1 - \delta$ over the random draw of samples, the iterates (θ_t) satisfy

$$\begin{aligned} & \mathbb{E}_{g^* \sim \rho} \left[\min_t \mathcal{L}_\infty^{\text{ex}}(g_{\theta_t}, g^*) \right] \\ & \lesssim \frac{1}{\delta} \left\{ r\kappa^{1/2} \|\Sigma\|_2^{1/2} \sqrt{\frac{k}{n_T}} + \frac{\delta_0}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma} \right. \\ & \quad \left. + \exp\left[\rho^3(r + \delta_0) \|\Sigma\|_2^{1/2}\right] \left[\frac{r \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log n_S} + \frac{\delta_0}{\sqrt{n_S}} \sqrt{\text{tr } \Sigma} \right] \right\}. \end{aligned}$$

F.4 Proofs

In this section, we prove the performance bound provided for logistic regression. First, we verify that the setting satisfies the assumptions of our general framework. Subsequently, we compute the quantities required to instantiate our bounds.

F.4.1 Verifying Assumptions of Section 4.2

In this section, we verify that the logistic regression setting, as described in Section F, satisfies the assumptions required by the general framework in Section 4.2.

It is easily verified that the logistic loss is 1-Lipschitz, convex, and that $|\ell(0, y)| \leq 1$ for $y \in \{0, 1\} = \mathcal{Y}$. Furthermore, as we have already characterized the approximate linearity of the function class in Lemma B.6, the approximate linearity assumption in Assumption 4.4 holds with high probability. Finally, \mathcal{C}_T is norm-bounded by $(\kappa r^2 + \delta_0^2)/\beta^2$.

Therefore, all that remains is verifying that a (ν, ε) -diversity condition holds in this setting. In what follows, we will do so by connecting the logistic loss to squared error loss via leveraging smoothness and local strong convexity, as was done by Tripuraneni et al. (2020b). Consequently, we can utilize the same argument as in the linear setting to obtain the desired diversity condition.

Lemma F.1 (Diversity condition, logistic regression). *Under the assumptions above, the source tasks satisfy a $(\Omega(\exp[-\rho^3(r + \delta_0) \|\Sigma\|_2^{1/2}]), 0)$ -diversity condition.*

Proof. We remark that under the choice of q and the logistic loss, we have that

$$\mathbb{E}_{x,y} [\ell(g_\theta(x), y) - \ell(g_{\theta'}(x), y)] = \mathbb{E}_x [\text{KL}(\text{Ber}(\sigma(g_\theta(x))) \| \text{Ber}(\sigma(g_{\theta'}(x))))]$$

Using the results in Tripuraneni et al. (2020b, Lemmas 2 and 3),

$$\frac{1}{8} \mathbb{E}_x [\exp(-\max(|g_\theta(x)|, |g_{\theta'}(x)|)) (g_\theta(x) - g_{\theta'}(x))^2] \leq \mathbb{E}_{x,y} [\ell(g_\theta(x), y) - \ell(g_{\theta'}(x), y)] \quad (13)$$

$$\mathbb{E}_{x,y} [\ell(g_\theta(x), y) - \ell(g_{\theta'}(x), y)] \leq \frac{1}{8} \mathbb{E}_x [(g_\theta(x) - g_{\theta'}(x))^2]. \quad (14)$$

In what follows, we will make use of (13) to lower bound the task-averaged best-case performance, and (14) to upper bound the expected best-case target performance. Note that this results in bounds in terms of best-case mean squared error – we can then proceed to use arguments similar to that of the linear case to connect the two inequalities.

Upper bounding the expected best-case target performance.

Assume that the optimal predictor is given by $\theta^* = B^* w^* + \delta^*$. Then, defining

$$\hat{\Delta} = \frac{1}{\beta} \delta^* w_0^\top \quad \text{and} \quad \bar{w} = \frac{1}{\beta} (B_0^\top \Sigma B_0)^\dagger B_0^\top \Sigma B^* w^*,$$

we note that $\hat{\delta} := ([\hat{\Delta}, 0], [0, \bar{w}^\top]^\top) \in \mathcal{C}_T^\beta$, following the argument in Lemma B.8, and achieves a squared-error excess risk of

$$\mathbb{E} [(\beta g_{\theta_0 + \hat{\delta}}(x) - g_{\theta^*}(x))^2] = \left\| \Sigma^{1/2} [(B_0 \bar{w} + \delta^*) - (B^* w^* - \delta^*)] \right\|_2^2 = \left\| P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* w^* \right\|_2^2.$$

Therefore, by taking the infimum of the excess risk over feasible δ and applying the excess risk upper bound in (14),

$$\begin{aligned} \inf_{\delta \in \mathcal{C}_T^\beta} \mathbb{E}_{x,y} [\ell(\beta g_{\theta_0 + \delta}(x), y) - \ell(g_{\theta^*}(x), y)] &\leq \mathbb{E}_{x,y} [\ell(\beta g_{\theta_0 + \hat{\delta}}(x), y) - \ell(g_{\theta^*}(x), y)] \\ &\leq \frac{1}{8} \mathbb{E} [(\beta g_{\theta_0 + \hat{\delta}}(x) - g_{\theta^*}(x))^2] \\ &\leq \frac{1}{8} \left\| P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* w^* \right\|_2^2. \end{aligned}$$

Finally, recall that $\mathbb{E}_{\theta^* \sim \rho} [(w^*)(w^*)^\top] = (r^2/k)I$. Therefore, by taking expectations with respect to ρ ,

$$\begin{aligned} & \mathbb{E}_{\theta^* \sim \rho} \left[\inf_{\delta \in \mathcal{C}_T^\beta} \mathbb{E}_{x,y} [\ell(\beta g_{\theta_0 + \delta}(x), y) - \ell(g_{\theta^*}(x), y)] \right] \\ & \lesssim \mathbb{E}_{\theta^* \sim \rho} \left[\text{tr} P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* (w^*)(w^*)^\top \left(P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* \right)^\top \right] \\ & = \frac{r^2}{k} \left\| P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* \right\|_F^2. \end{aligned}$$

Lower bounding the task-averaged best-case performance.

We proceed to lower bound the task-averaged best-case performance. To proof proceeds similarly to that of Tripuraneni et al. (2020b), lower bounding the excess risk by a constant multiple of the squared-error excess risk. Then, the result follows by an application of the transfer lemma in Lemma B.4.

Fix a task $t \in [T]$. Define $Z_1 := x^\top (B^* w_t^* + \delta_t^*)$ and $Z_2 := x^\top (B_0 w_t + \delta_t)$. We see that Z_1 is sub-Gaussian with parameter $\rho^2 \|\Sigma^{1/2} (B^* w_t^* + \delta_t^*)\|_2^2 \leq \rho^2 (r^2 + \delta_0^2) \|\Sigma\|_2 =: \sigma^2$. Note that Z_2 is also sub-Gaussian with proxy σ^2 . Then, by applying the lower bound in (15),

$$\mathbb{E}_{x,y} [\ell(g_{\hat{\theta}_t}(x), y) - \ell(g_{\theta_t^*}(x), y)] \geq \frac{1}{8} \mathbb{E}_x [e^{-\max(|Z_1|, |Z_2|)} (g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2]. \quad (15)$$

We will now show that the right-hand side is lower bounded by a constant multiple of the squared error excess risk. Intuitively, since Z_1 and Z_2 concentrate around 0, $e^{-\max(|Z_1|, |Z_2|)}$ concentrates around 1.

Formally, consider the event $E_\alpha := \{\max(|Z_1|, |Z_2|) \leq \alpha\sigma\}$ for any α . Since the quantity inside the expectation of (15) is non-negative,

$$\begin{aligned} & \frac{1}{8} \mathbb{E}_x [e^{-\max(|Z_1|, |Z_2|)} (g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2] \\ & \geq \frac{1}{8} \mathbb{E} [\mathbb{1}[E_\alpha] e^{-\alpha\sigma} (g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2] \\ & \geq \frac{e^{-\alpha\sigma}}{8} \left\{ \mathbb{E} [(g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2] - \mathbb{E} [\mathbb{1}[E_\alpha^C] (g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2] \right\}. \end{aligned}$$

Now, we upper bound the last term. By Cauchy-Schwarz,

$$\mathbb{E} [\mathbb{1}[E_\alpha^C] (g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2] \leq \sqrt{P(E_\alpha^C)} \sqrt{\mathbb{E} [(g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^4]},$$

where we have used the fact that $\text{Var}[X] \leq \sigma^2$ for a σ^2 -sub-Gaussian random variable. We use properties of sub-Gaussian random variables to bound both factors. First, by Chebyshev's inequality,

$$P(E_\alpha^C) \leq 2P(|Z_1| \geq \alpha\sigma) \leq \frac{2\text{Var}[Z_1]}{\alpha^2\sigma^2} \leq \frac{2}{\alpha^2}.$$

For the second factor, $g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x)$ is sub-Gaussian with proxy $\rho^2 \mathbb{E} [(g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2]$, and so via an equivalent definition of sub-Gaussian random variables,

$$\sqrt{\mathbb{E} [(g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^4]} \lesssim \rho^2 \mathbb{E} [(g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2].$$

Putting these bounds together, we have that

$$\mathbb{E} [\mathbb{1}[E_\alpha^C] (g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2] \lesssim \frac{\rho^2}{\alpha} \mathbb{E} [(g_{\hat{\theta}_t}(x) - g_{\theta_t^*}(x))^2].$$

Therefore, by setting $\alpha \asymp \rho^2$, we can ensure that

$$\begin{aligned} \mathbb{E}_{x,y} \left[\ell(g_{\hat{\theta}_t}(x), y) - \ell(g_{\theta_t^*}(x), y) \right] &\gtrsim e^{-\rho^3(r+\delta_0)\|\Sigma\|_2^{1/2}} \mathbb{E} \left[(x^\top \hat{\theta}_t - x^\top \theta_t^*)^2 \right] \\ \implies \frac{1}{T} \sum_{t \in [T]} \mathbb{E}_{x,y} \left[\ell(g_{\hat{\theta}_t}(x), y) - \ell(g_{\theta_t^*}(x), y) \right] \\ &\gtrsim \frac{e^{-\rho^3(r+\delta_0)\|\Sigma\|_2^{1/2}}}{T} \sum_{t \in [T]} \left\| \Sigma^{1/2} (B_0 w_t + \delta_t - B^* w_t^* + \delta_t^*) \right\|_2^2. \end{aligned}$$

By applying Lemma B.4, we finally obtain the final lower bound

$$\frac{e^{\rho^3(r+\delta_0)\|\Sigma\|_2^{1/2}}}{T} \sum_{t \in [T]} \inf_{\delta_t \in \mathcal{C}_S} \mathbb{E}_{x,y} \left[\ell(g_{\theta_0+\delta_t}(x), y) - \ell(g_{\theta_t^*}(x), y) \right] \geq \frac{r^2}{k} \left\| P_{\Sigma^{1/2} B_0}^\perp \Sigma^{1/2} B^* \right\|_F^2.$$

Putting everything together, we thus have the desired diversity condition

$$\begin{aligned} &\inf_{\|\delta\| \leq \delta_0} \mathbb{E}_{x,y} \left[\ell(\beta g_{\theta_0+\delta}(x), y) - \ell(g_{\theta^*}(x), y) \right] \\ &\lesssim \frac{e^{\rho^3(r+\delta_0)\|\Sigma\|_2^{1/2}}}{T} \sum_{t \in [T]} \inf_{\|\delta_t\| \leq \delta_0} \mathbb{E}_{x,y} \left[\ell(g_{\theta_0+\delta_t}(x), y) - \ell(g_{\theta_t^*}(x), y) \right]. \quad \square \end{aligned}$$

F.4.2 Computations

Having demonstrated that the assumptions in Section 4.2 hold, we proceed to calculate the relevant quantities required for establishing a performance bound in this setting. First, we compute the Rademacher complexity for source task training.

Lemma F.2 (Source Rademacher Bound). *Assume that $n_S \gtrsim \rho^4 d$. Then, we can bound the source Rademacher complexity as*

$$\frac{1}{T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] \lesssim \frac{r \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log n_S} + \frac{\delta_0}{\sqrt{n_S}} \sqrt{\text{tr } \Sigma}.$$

Proof. Observe that if $(\Delta_t, w_t) \in \mathcal{C}_S$, then $\delta_t := \Delta_t w_t$ satisfies $\|\delta_t\|_2 \leq \delta_0$. Therefore, we have that

$$\begin{aligned} &\frac{1}{T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{\mathcal{C}_S}(\theta)]^{\otimes T} \right] \\ &= \mathbb{E} \left[\sup_{\substack{B \\ (\Delta_t, w_t) \in \mathcal{C}_S}} \frac{1}{n_S T} \sum_{t \in [T]} \langle \varepsilon_t, X_t (B + \Delta_t) w_t \rangle \right] \\ &= \underbrace{\mathbb{E} \left[\sup_{\substack{B \\ \|w_t\|_2 \leq r}} \frac{1}{n_S T} \sum_{t \in [T]} \langle \varepsilon_t, X_t B w_t \rangle \right]}_{=:(I)} + \underbrace{\mathbb{E} \left[\sup_{\|\delta_t\|_2 \leq \delta_0} \frac{1}{n_S T} \sum_{t \in [T]} \langle \varepsilon_t, X \delta_t \rangle \right]}_{=:(II)} \end{aligned}$$

We proceed to bound these two quantities separately.

Bounding (I) via discretization.

In this section, we bound the complexity by discretizing the set $\mathcal{O}^{k \times d}$ of $(k \times d)$ orthogonal matrices. We remark that the argument is similar in form to that of Lemma B.3.

Let \mathcal{S} be an ε -covering of $\mathcal{O}^{k \times d}$ in the Frobenius norm with at most $(6\sqrt{k}/\varepsilon)^{dk}$ elements guaranteed by Proposition I.1. Then,

$$(I) \leq \underbrace{\mathbb{E} \left[\sup_{\substack{B \in \mathcal{S} \\ \|w_t\|_2 \leq r}} \frac{1}{n_S T} \sum_{t \in [T]} \langle \varepsilon_t, X_t B w_t \rangle \right]}_{=:(A)} + \underbrace{\mathbb{E} \left[\sup_{\substack{B, \bar{B} \in \mathcal{S} \\ \|B - \bar{B}\|_F \leq \varepsilon \\ \|w_t\|_2 \leq r}} \frac{1}{n_S T} \sum_{t \in [T]} \langle \varepsilon_t, X_t (B - \bar{B}) w_t \rangle \right]}_{=:(B)}$$

To bound (A), we bound the corresponding Gaussian complexity and use the fact that $\mathcal{R}(\cdot) \lesssim \mathcal{G}(\cdot)$. Consequently, via multiple applications of Cauchy-Schwarz,

$$\begin{aligned} (A) &\leq \sqrt{\frac{\pi}{2}} \mathbb{E} \left[\sup_{\substack{B \in \mathcal{S} \\ \|w_t\|_2 \leq r}} \frac{1}{n_S T} \sum_{t \in [T]} \langle P_{X_t B} z_t, X_t B w_t \rangle \right] \\ &\lesssim \sqrt{\mathbb{E} \left[\sup_{B \in \mathcal{S}} \frac{1}{n_S T} \sum_{t \in [T]} \|P_{X_t B} z_t\|_2^2 \right]} \sqrt{\mathbb{E} \left[\sup_{\substack{B \in \mathcal{S} \\ \|w_t\|_2 \leq r}} \frac{1}{n_S T} \sum_{t \in [T]} \|X_t B w_t\|_2^2 \right]}. \end{aligned}$$

Conditioned on X_t , $\sum_{t \in [T]} \|P_{X_t B} z_t\|_2^2$ is distributed as a chi-squared random variable with kT degrees of freedom, with mean kT . Therefore, using known bounds on expectations of finite maxima of subexponential random variables,

$$\begin{aligned} \mathbb{E} \left[\sup_{B \in \mathcal{S}} \sum_{t \in [T]} \|P_{X_t B} z_t\|_2^2 - kT \right] &\lesssim \sqrt{kT \log |\mathcal{S}|} + \log |\mathcal{S}| \\ \implies \mathbb{E} \left[\sup_{B \in \mathcal{S}} \sum_{t \in [T]} \|P_{X_t B} z_t\|_2^2 \right] &\lesssim kT + \log |\mathcal{S}|. \end{aligned} \tag{16}$$

Furthermore, by applying the expectation bound on the empirical spectral norm in Proposition I.5,

$$\sqrt{\mathbb{E} \left[\sup_{\substack{B \in \mathcal{S} \\ \|w_t\|_2 \leq r}} \frac{1}{n_S T} \sum_{t \in [T]} \|X_t B w_t\|_2^2 \right]} = r \sqrt{\mathbb{E} \left[\sup_{B \in \mathcal{S}} \frac{1}{T} \sum_{t \in [T]} \lambda_{\max} \left(\frac{B^\top X_t^\top X_t B}{n_S} \right) \right]} \tag{17}$$

$$\begin{aligned} &\leq r \sqrt{\mathbb{E} \left[\lambda_{\max} \left(\frac{X_t^\top X_t}{n_S} \right) \right]} \\ &\lesssim r \|\Sigma\|_2^{1/2}. \end{aligned} \tag{18}$$

Therefore, by combining the inequalities from (16) and (18),

$$(A) \lesssim \frac{\|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + \log |\mathcal{S}|}.$$

We proceed to bound (B), the error arising from discretization. We have that

$$(B) \leq \mathbb{E} \left[\sup_{\substack{B, \bar{B} \in \mathcal{S} \\ \|B - \bar{B}\|_F \leq \varepsilon \\ \|w_t\|_2 \leq r}} \frac{1}{n_S T} \sum_{t \in [T]} \langle \varepsilon_t, X_t (B - \bar{B}) w_t \rangle \right] \lesssim r \varepsilon \left\| \Sigma^{1/2} \right\|_2.$$

Finally, by setting $\varepsilon = \sqrt{k/n_S}$, we have an overall bound of

$$(I) \leq \frac{r \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log n_S}.$$

Bounding (II).

By bounding by the Gaussian complexity,

$$(II) \leq \frac{1}{n_S T} \mathbb{E} \left[\sup_{\|\delta_t\|_2 \leq \delta_0} \sum_{t \in [T]} \langle z_t, X_t \delta_t \rangle \right] \lesssim \frac{\delta_0}{\sqrt{n_S}} \mathbb{E} \left[\left\| \frac{1}{\sqrt{n}} X_1^\top z \right\|_2 \right] \leq \frac{\delta_0}{\sqrt{n_S}} \sqrt{\text{tr } \Sigma}.$$

Putting the bounds on (I) and (II) together, we thus have that

$$\frac{1}{T} \mathcal{R}_{n_S} \left[\bigcup_{\theta \in \Theta_0} [\mathcal{A}_{C_S}(\theta)]^{\otimes T} \right] \lesssim \frac{r \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log n_S} + \frac{\delta_0}{\sqrt{n_S}} \sqrt{\text{tr } \Sigma}. \quad \square$$

Now, we compute the Rademacher complexity term associated with target task training.

Lemma F.3 (Target Rademacher Bound). *For $\beta \gtrsim r\kappa^{1/2}$, the Rademacher complexity of the feasible set during target time training is bounded by*

$$\sup_{\theta \in \Theta_0} \mathcal{R}_{n_T} [\mathcal{A}_{C_T}(\theta_0)] \lesssim r\kappa^{1/2} \|\Sigma\|_2^{1/2} \sqrt{\frac{k}{n_T}} + \frac{\delta_0}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma}.$$

Proof. We can write the Rademacher complexity for a fixed $\theta_0 \in \Theta_0$ as

$$\mathcal{R}_{n_T} [\mathcal{A}_{C_T}(\theta_0)] = \beta \mathbb{E} \left[\sup_{(\Delta, w) \in \mathcal{C}_T^\beta} \frac{1}{n_T} \langle \varepsilon, X[B_0, B_0]w + X\Delta w_0 + X\Delta w \rangle \right].$$

Firstly, by converting to Gaussian complexity and following standard arguments, we can bound the first term as

$$\begin{aligned} \beta \mathbb{E} \left[\sup_{\|w\|_2 \leq r\kappa^{1/2}/\beta} \frac{1}{n_T} \langle z, X[B_0, B_0]w \rangle \right] &\leq \beta \mathbb{E} \left[\sup_{\|w\|_2 \leq \kappa/\beta} \frac{1}{n_T} \langle [B_0, B_0]^\top X^\top z, w \rangle \right] \\ &\leq \frac{r\kappa^{1/2}}{\sqrt{n_T}} \mathbb{E} \left[\left\| B_0^\top \frac{X^\top}{\sqrt{n_T}} z \right\|_2 \right] \\ &\leq r\kappa^{1/2} \|\Sigma\|_2^{1/2} \sqrt{\frac{k}{n_T}}. \end{aligned}$$

For the second term, note that we only need to consider the set of rank-1 matrices of the form $\Delta = v[w_0, -w_0]^\top$, where $\|v\|_2 \leq \delta_0/\beta\sqrt{2}$. As such,

$$\beta \mathbb{E} \left[\sup_{\|v\|_2 \leq \delta_0/\beta\sqrt{2}} \frac{1}{n_T} \langle z, Xv \rangle \right] \leq \frac{\delta_0}{\sqrt{n_T}} \sqrt{\mathbb{E} \left[\left\| \frac{X}{\sqrt{n_T}} z \right\|_2^2 \right]} \lesssim \frac{\delta_0}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma}.$$

Finally, to bound the Hessian term, note that $\|\Delta w\|_2 \leq r\kappa^{1/2}\delta_0/\beta^2$. Therefore, by a similar argument to the previous term,

$$\beta \mathbb{E} \left[\sup_{(\Delta, w) \in \mathcal{C}_T} \frac{1}{n_T} \langle \varepsilon, X\Delta w \rangle \right] \leq \frac{\kappa\delta_0}{\beta\sqrt{n_T}} \sqrt{\text{tr } \Sigma}.$$

Thus, as long as $\beta \geq r\kappa^{1/2}$, we can take suprema and obtain the overall bound of

$$\sup_{\theta \in \Theta_0} \mathcal{R}_{n_T} [\mathcal{A}_{C_T}(\theta)] \lesssim r\kappa^{1/2} \|\Sigma\|_2^{1/2} \sqrt{\frac{k}{n_T}} + \frac{\delta_0}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma}. \quad \square$$

F.4.3 Compiling the Bound

Having performed all required computations, we now prove the provided performance guarantee.

Theorem F.1 (Performance Guarantee for Logistic Regression). *Assume that Assumption F.1 and Assumption F.2 both hold. Set the parameters for target time training to be*

$$\beta = (\kappa r^2 + \delta_0^2) \max \left(1, \frac{\sqrt{\text{tr } \Sigma}}{1/\delta \sqrt{n_T}} \right) \quad \text{and} \quad T_{\text{PGD}} = \frac{(\kappa r^2 + \delta_0^2) \text{tr } \Sigma}{1/\delta^2 n_T}.$$

Then, for $n_S, n_T \gtrsim \rho^4 d$, we have that with probability at least $1 - \delta$ over the random draw of samples, the iterates (θ_t) satisfy

$$\begin{aligned} & \mathbb{E}_{g^* \sim \rho} \left[\min_t \mathcal{L}_\infty^{\text{ex}}(g_{\theta_t}, g^*) \right] \\ & \lesssim \frac{1}{\delta} \left\{ r \kappa^{1/2} \|\Sigma\|_2^{1/2} \sqrt{\frac{k}{n_T}} + \frac{\delta_0}{\sqrt{n_T}} \sqrt{\text{tr } \Sigma} \right. \\ & \quad \left. + \exp \left[\rho^3 (r + \delta_0) \|\Sigma\|_2^{1/2} \right] \left[\frac{r \|\Sigma\|_2^{1/2}}{\sqrt{n_S T}} \sqrt{kT + kd \log n_S} + \frac{\delta_0}{\sqrt{n_S}} \sqrt{\text{tr } \Sigma} \right] \right\}. \end{aligned}$$

Proof. With probability at least $1 - \delta/2$, the approximate linearity property of the function class holds via Lemma B.2 and Lemma B.6. Therefore, by combining all prior calculations and instantiating Theorem 4.1 with failure probability $\delta/2$, projected gradient descent finds a predictor g satisfying the desired bound in the theorem statement. \square

G The Nonlinear Hard Case

G.1 Construction

In what follows, we establish the existence of a nonlinear setting where there exists a sample complexity separation between ADAPTREP and FROZENREP, similar to Section 3.4. As before, fix $k, d \in \mathbb{N}$ with $2k < d$. The construction relies on the observation that linear predictors lying in a rank- k space are representable as linear functions of $2k$ appropriately chosen ReLU neurons.

Following the discussion in Section E, note that when we take $\beta \rightarrow \infty$, the resulting function class can be expressed as

$$f_{(B+\Delta, w)} = w^\top \sigma(B^\top x) + \langle x \sigma'(x^\top B), \Delta \rangle,$$

where $B, \Delta \in \mathbb{R}^{d \times 2k}$, and $w \in \mathbb{R}^{2k}$. We further constrain B so that the first k columns are equal to the negation of the last k columns. Finally, we choose $\sigma(x) = \max(x, 0)$, which we will also write as x_+ for convenience⁸. For convenience, we follow the convention in Section E of writing ϕ_B, ψ_B and ρ_B for the activation, gradient, and concatenated features corresponding to B , respectively, as defined in Definition E.1.

We briefly review the construction in Section 3.4. Let the input distribution p be a Gaussian distribution on \mathbb{R}^d with covariance

$$\Sigma = \begin{bmatrix} \varepsilon I_{d-k} & 0 \\ 0 & I_k \end{bmatrix}$$

for a fixed $\varepsilon \in (0, 1)$. Furthermore, we define $E^*, E_k \subset \mathbb{R}^d$ to be the two eigenspaces of Σ determined by the two blocks, *i.e.*

$$E^* = \text{Col} \begin{bmatrix} \varepsilon I_{d-k} \\ 0 \end{bmatrix} \quad \text{and} \quad E_k = \text{Col} \begin{bmatrix} 0 \\ I_k \end{bmatrix}.$$

Then, for any orthogonal matrix $A \in \mathbb{R}^{d \times k}$ with $\text{Col } A \subseteq E^*$, define a distribution over θ given by

$$\theta = \frac{1}{\sqrt{2\varepsilon}} Av + \delta, \tag{19}$$

where v and δ are sampled uniformly at random from the unit spheres in \mathbb{R}^k and E_k , respectively.

Now, we lift this linear task distribution setting into the ReLU setting. In particular, we sample the source tasks by fixing orthogonal $A \in \mathbb{R}^{d \times k}$, sampling v and δ as before, and letting the optimal predictor be $f_{(B+\Delta, w)}$, where

$$B := [A, -A], \quad w := \frac{1}{\sqrt{2\varepsilon}}[v, -v], \quad \Delta := \frac{1}{k} \mathbb{1} \delta^\top. \tag{20}$$

One can easily verify algebraically that

$$f_{(B+\Delta, w)}(x) = x^\top \left(\frac{1}{\sqrt{2\varepsilon}} Av + \delta \right),$$

as desired. As before, we consider the family of task distributions induced by any A^* with $\text{Col } A^* \subseteq E^*$.

With the above task distribution, we can then prove the following hardness result on FROZENREP:

Theorem G.1 (FROZENREP Minimax Bound, ReLU). *For an orthogonal matrix $A^* \in \mathbb{R}^{d \times k}$ such that $\text{Col } A^* \subset E^*$, let $B^* = [A^*, -A^*]$, and define S_{A^*} to be the set*

$$S_{A^*} = \left\{ \frac{1}{\sqrt{2\varepsilon}} A^* v + \delta \mid \|v\|_2 \leq 1, \|\delta\|_2 \leq 1, \delta \in E_k \right\}.$$

Furthermore, let \hat{B} be the output of FROZENREP with access to infinite per-task samples and tasks, with the task distribution determined by B^ . Then, with high probability over the draw of $n_T \gg d$ samples during target training, we have that*

$$\min_{\hat{B}, \hat{w}, \hat{\Delta}} \max_{\theta^* \in S_{A^*}} \mathbb{E} \left[\frac{1}{n_T} \left\| X \theta^* - f_{(\hat{B} + \hat{\Delta}, \hat{w})}(X) \right\|_2^2 \right] \gtrsim \frac{\sigma^2 d}{n_T},$$

⁸Note that $\sigma'(x) = \mathbb{1}[x > 0]$.

where \bar{B} , \hat{w} , and $\hat{\Delta}$ are all measurable functions of (X, y) to $\mathbb{R}^{d \times 2k}$, \mathbb{R}^{2k} , and $\mathbb{R}^{d \times 2k}$, respectively, and $\text{Col } \bar{B} = \text{Col } \hat{B}$. Furthermore, the expectation is over the randomness in the labels $y \sim \mathcal{N}(X\theta^*, \sigma^2 I_{n_T})$.

In contrast, we have the following upper bound on the performance of ADAPTREP:

Lemma G.1 (Adaptation target performance, ReLU). *Let $B_0 = [A_0, -A_0]$, and fix a $\theta^* \in S_{A^*}$. Consider a learner which solves*

$$\min_B \min_{\substack{(w_t, \Delta_t) \\ \|\Delta_t\|_F \leq 1}} \frac{1}{n_S} \sum_{t \in [T]} \|\rho_B(X)[w, \Delta] - y\|_2^2.$$

during source training, and

$$\min_{\substack{w, \Delta \\ \|\Delta\|_F \leq 1}} \frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - y\|_2^2$$

during target training, where B_0 is the representation obtained from source training. Finally, we set $k = \Theta(1)$ and $\varepsilon = k/d$. Then, with access to infinite per-task samples and tasks during source training, the learner achieves target loss bounded as

$$\frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - X\theta^*\|_2^2 \lesssim \frac{\sigma^2}{n_T} \left(1 + \log \frac{1}{\delta}\right) + \frac{\sigma}{\sqrt{n_T}} \sqrt{1 + \log \frac{1}{\delta}}$$

with probability at least $1 - \delta$ over the draw of target samples.

Proofs of the above results are provided in the following section. As before, we compare the two methods when $n_T = \Theta(d)$. Then, from the results above, the lower bound on the loss of FROZEN-REP is $\Omega(1)$, while the upper bound on the loss of ADAPTREP is $O(1/\sqrt{n_T})$. Therefore, we also see a strict separation between the two methods within this setting as well, which grows with $n_T \rightarrow \infty$.

G.2 Proofs

Throughout this section, we write $\theta_t^* = A^*v_t^* + \delta_t^*$ for the linear predictor parameter for source task $t \in [T]$. Furthermore, we assume that this linear predictor corresponds to ReLU predictor $f_{(B^* + \Delta^*, w^*)}$. First, we prove the following intermediate technical result which will be used throughout this section.

Lemma G.2 (Optimal ReLU Predictor is Linear). *Let $\theta \in \mathbb{R}^d$, and assume that $B = [A, -A]$. Then, if we set*

$$w = \operatorname{argmin}_{w \in \mathbb{R}^{2k}} \mathbb{E} [(w^\top (B^\top x)_+ - x^\top \theta_t^*)^2],$$

then the predictor $x \mapsto w^\top (B^\top x)$ is in fact a linear function of x , and takes the form of $x^\top A v$ for some $v \in \mathbb{R}^k$.

Proof. Write $w = [w_+, w_-]$, and thus the objective defining w can be written as

$$\min_{w_+, w_- \in \mathbb{R}^k} \mathbb{E} [(w_+^\top (A^\top x)_+ + w_-^\top (-A^\top x)_+ - x^\top \theta_t^*)^2].$$

Now, define the matrices

$$\begin{aligned} \Omega &:= \mathbb{E} [(A^\top x)_+ (x^\top A)_+] = \mathbb{E} [(-A^\top x)_+ (-x^\top A)_+] \\ \Gamma &:= \mathbb{E} [(-A^\top x)_+ (x^\top A)_+] = \mathbb{E} [(A^\top x)_+ (-x^\top A)_+], \end{aligned}$$

where the equalities follow from the fact that Ax and $-Ax$ are equal in distribution. Additionally,,

$$\mathbb{E} [(A^\top x)_+ x^\top] = \frac{1}{2} \mathbb{E} [(A^\top x)_+ x^\top] - \frac{1}{2} \mathbb{E} [(-A^\top x)_+ x^\top] = \frac{1}{2} A^\top \mathbb{E} [xx^\top] = \frac{1}{2} A^\top \Sigma,$$

and thus once again since Ax and $-Ax$ are equal in distribution,

$$\mathbb{E} [(-A^\top x)_+ x^\top] = -\mathbb{E} [(A^\top x)_+ x^\top] = -\frac{1}{2} A^\top \Sigma.$$

Thus, by explicitly solving for the optimum of the convex objective,

$$\begin{aligned} \begin{bmatrix} w_+ \\ w_- \end{bmatrix} &= \frac{1}{2} \begin{bmatrix} \Omega & \Gamma \\ \Gamma & \Omega \end{bmatrix}^{-1} \begin{bmatrix} A^\top \Sigma \theta_t^* \\ -A^\top \Sigma \theta_t^* \end{bmatrix} \\ &= \begin{bmatrix} (\Omega - \Gamma \Omega \Gamma)^{-1} & 0 \\ 0 & (\Omega - \Gamma \Omega \Gamma)^{-1} \end{bmatrix} \begin{bmatrix} I & -\Gamma \Omega^{-1} \\ -\Gamma \Omega^{-1} & I \end{bmatrix} \begin{bmatrix} A^\top \Sigma \theta_t^* \\ -A^\top \Sigma \theta_t^* \end{bmatrix}, \end{aligned}$$

where we have applied blockwise matrix inversion to obtain the second inequality. From this, we see that $w_+ = w_-$, and thus the corresponding optimal predictor is linear, and can be shown to be $x \mapsto x^\top A w_+$. \square

G.2.1 Hardness Result for FROZENREP

Lemma G.3 (FROZENREP learns incorrect neurons). *Fix an orthogonal matrix $A^* \in \mathbb{R}^{d \times k}$ with $\text{Col } A^* \subseteq E^*$, and assume that we sample tasks from the distribution in (20). Assume that $\hat{B} = [\hat{A}, -\hat{A}]$ is the representation found by FROZENREP. Then, with infinitely many tasks and per-task samples (i.e. $n_S, T \rightarrow \infty$), $\hat{A} = \text{Col } E_k$.*

Proof. Intuitively, the result follows from the fact that the optimal predictor is equivalent to a linear predictor, and thus the result follows by Lemma C.1. More formally, by Lemma G.2,

$$\begin{aligned} \mathcal{L}(B) &= \frac{1}{T} \sum_{t \in [T]} \min_w \mathbb{E} [(w^\top (Bx)_+ - x^\top \theta_t^*)^2] = \frac{1}{T} \sum_{t \in [T]} \min_v \mathbb{E} [(x^\top A v - x^\top \theta_t^*)^2] \\ &= \frac{1}{T} \sum_{t \in [T]} \min_v \left\| \Sigma^{1/2} (A v - \theta_t^*) \right\|_2^2 = \frac{1}{T} \sum_{t \in [T]} \left\| P_{\Sigma^{1/2} A}^\perp \Sigma^{1/2} \theta_t^* \right\|_2^2. \end{aligned}$$

At this point, we recognize that the objective is equivalent to the one analyzed in Lemma C.1, and thus the same characterization of global optima holds. That is, the solution \hat{B} found by FROZENREP can be expressed as $\hat{B} = [\hat{A}, -\hat{A}]$, where $\text{Col } \hat{A} = E_k$. \square

Theorem G.1 (FROZENREP Minimax Bound, ReLU). *For an orthogonal matrix $A^* \in \mathbb{R}^{d \times k}$ such that $\text{Col } A^* \subset E^*$, let S_{A^*} be the set*

$$S_{A^*} = \left\{ \frac{1}{\sqrt{2\varepsilon}} A^* v + \delta \mid \|v\|_2 \leq 1, \|\delta\|_2 \leq 1, \delta \in E_k \right\}.$$

We consider the following procedure:

1. We draw n_T samples for target time training, which are collected into a matrix X .
2. Player chooses target-time estimator $x \mapsto \hat{w}^\top (\bar{B}^\top x) + \mathbb{1}[x^\top \bar{B} \geq 0] \hat{\Delta}^\top x$, where \bar{B} , \hat{w} and $\hat{\Delta}$ are measurable functions of (X, y) , and \bar{B} is an orthogonal matrix with $\text{Col } \bar{B} = \text{Col } \hat{B}$.
3. Adversary chooses an orthogonal matrix $A^* \in \mathbb{R}^{d \times k}$ satisfying $\text{Col } A^* \subset E_k$, and a target time predictor $\theta^* \in S_{A^*}$.
4. Compute the representation \hat{B} returned by FROZENREP under the setting of Lemma G.3 with the task distribution determined by A^* .
5. Target time samples are generated using $y \sim \mathcal{N}(X\theta^*, \sigma^2 I_n)$, and the player estimator is evaluated.

Then, with probability at least $1 - \delta$ over the draw of X , we have that

$$\min_{\bar{B}, \hat{w}, \hat{\Delta}} \max_{\theta^* \in S_{A^*}} \mathbb{E} \left[\frac{1}{n_T} \sum_{i=1}^{n_T} \left(x_i^\top \theta^* - \hat{w}^\top (\bar{B}^\top x_i) - \mathbb{1}[x_i^\top \bar{B} > 0] \hat{\Delta}^\top x_i \right)^2 \right] \gtrsim \frac{\sigma^2 d}{n_T},$$

where the expectation is over the randomness in the labels $y \sim \mathcal{N}(X\theta^, \sigma^2 I_{n_T})$.*

Proof. Throughout this proof, we assume the high-probability event in Lemma B.2, which guarantees that with probability at least $1 - \delta$,

$$0.9\Sigma \preceq \frac{1}{n_T} X^\top X \preceq 1.1\Sigma.$$

Let S denote the intersection of the d -dimensional unit sphere with E_k . For any infinite-dimensional vector f indexed by S and a measure μ on S , we define $\mu^\top f$ to denote integration with respect to μ , i.e.

$$\mu^\top f = \int_S f \, d\mu.$$

We then define the infinite-dimensional vectors η and $\zeta^{(i)}$ for $i = 1, \dots, d$, both indexed by S , as

$$\eta(x) := (v^\top x)_+ \quad \text{and} \quad \zeta^{(i)}(x) := \mathbb{1}[v^\top x > 0] x_i.$$

Recall that by Lemma G.3, $\text{Col } \hat{B} = E_k$. With the preceding discussion in mind, we can equivalently think of the player as choosing $d + 1$ signed measures $\alpha, \beta_1, \dots, \beta_d$ over S , all with a common support of $2k$ elements in S , as a function of (X, y) . The player then plays the predictor

$$x \mapsto \alpha^\top \eta(x) + \sum_{i=1}^d \beta_i^\top \zeta^{(i)}(x).$$

Then, if we let $T = \left\{ \theta \mid \|P_{E^*} \theta\|_2^2 \leq 1/2\varepsilon, \|P_{E_k} \theta\|_2^2 \leq 1 \right\} \subseteq S_{A^*}$, then we have the inequality

$$\begin{aligned} & \min_{\hat{w}, \hat{\Delta}} \max_{\theta^* \in S_{A^*}} \mathbb{E} \left[\frac{1}{n_T} \sum_{i=1}^{n_T} \left(x_i^\top \theta^* - \hat{w}^\top (\hat{B}^\top x_i)_+ - \mathbb{1}[x_i^\top \hat{B} > 0] \hat{\Delta}^\top x \right)^2 \right] \\ &= \min_{\alpha, \beta_1, \dots, \beta_d} \max_{\theta^* \in S_{A^*}} \mathbb{E} \left[\frac{1}{n_T} \sum_{i=1}^{n_T} \left(x_i^\top \theta^* - \alpha^\top \eta(x) - \sum_{j=1}^d \beta_j^\top \zeta^{(j)}(x) \right)^2 \right] \\ &= \min_{\alpha, \beta_1, \dots, \beta_d} \max_{\theta^* \in T} \mathbb{E} \left[\frac{1}{n_T} \sum_{i=1}^{n_T} \left(x_i^\top \theta^* - \alpha^\top \eta(x) - \sum_{j=1}^d \beta_j^\top \zeta^{(j)}(x) \right)^2 \right], \end{aligned}$$

where the second equality follows from the fact that the expression has no direct dependence on A^* .

The rest of the argument follows makes use of Fano's minimax bound, applied similarly as in the linear setting. First, note that for any $\theta \in T$, there exists $d + 1$ signed measures on S with common support size $2k$, which we denote with α^θ and $\beta_1^\theta, \dots, \beta_d^\theta$, such that $x^\top \theta = (\alpha^\theta)^\top \eta(x) + \sum_j (\beta_j^\theta)^\top \zeta^{(j)}(x)$. Using this construction, we can lift the local packing set from the linear setting into this setting, where the associated seminorm (defined over vector measures) is given by

$$\|(\alpha, \beta_1, \dots, \beta_d)\|^2 = \frac{1}{n_T} \sum_{i=1}^{n_T} \left(\alpha^\top \eta(x) + \sum_{j=1}^d \beta_j^\top \zeta^{(j)}(x) \right)^2.$$

More formally, let B be the unit ball under the Σ -norm, which we observe satisfies $(1/\sqrt{2})B \subseteq T$. Recall that there exists a $(1/2)$ -packing of the unit Σ -ball in the Σ -norm with at least 2^d elements, via a standard volumetric argument - let this set be P . Equivalently, there exists a $(2\delta/\sqrt{0.9})$ -packing of $(4\delta/\sqrt{0.9})B$ with at least 2^d elements, which we denote as P . Note that for any $\theta, \theta' \in P$,

$$\left\| (\alpha^\theta, \beta_1^\theta, \dots, \beta_d^\theta) - (\alpha^{\theta'}, \beta_1^{\theta'}, \dots, \beta_d^{\theta'}) \right\|^2 = \frac{1}{n_T} \|X(\theta - \theta')\|_2^2 \geq 0.9 \|\theta - \theta'\|_\Sigma^2 > 4\delta^2,$$

which implies that the vector measures corresponding to the elements of P are 2δ -separated in the associated seminorm. Furthermore, for any $\theta, \theta' \in P$,

$$\text{KL}(\mathcal{N}(X\theta, \sigma^2 I_{n_T}) \parallel \mathcal{N}(X\theta', \sigma^2 I_{n_T})) = \frac{1}{2\sigma^2} \|X(\theta - \theta')\|_2^2 \lesssim \frac{n_T}{2\sigma^2} \|\theta - \theta'\|_\Sigma^2 \leq \frac{32n_T}{0.9\sigma^2} \delta^2.$$

Therefore, for any $\delta^2 \leq 1/32$, which ensures that $P \subseteq 4\delta B \subseteq (1/\sqrt{2})B \subseteq T$, Fano's inequality implies that

$$\begin{aligned} \min_{\alpha, \beta_1, \dots, \beta_d} \max_{\theta^* \in T} \mathbb{E} \left[\left\| (\alpha^{\theta^*}, \beta_1^{\theta^*}, \dots, \beta_d^{\theta^*}) - (\alpha, \beta_1, \dots, \beta_d) \right\|^2 \right] &\geq \delta^2 \left(1 - \frac{32n_T}{0.9\sigma^2 \log 2} \delta^2 - \frac{1}{d} \right) \\ &\geq \delta^2 \left(\frac{3}{4} - \frac{32n_T}{0.9\sigma^2 d \log 2} \delta^2 \right). \end{aligned}$$

Therefore, as long as $n_T \geq \left(\frac{0.9 \log 2}{2} \right) \sigma^2 d$, we can set $\delta^2 = \left(\frac{0.9 \log 2}{64} \right) \frac{\sigma^2 d}{n_T}$, which implies that

$$\min_{\hat{w}, \hat{\Delta}} \max_{\theta^* \in \mathcal{S}_{A^*}} \mathbb{E} \left[\frac{1}{n_T} \sum_{i=1}^{n_T} \left(x_i^\top \theta^* - \hat{w}^\top (\bar{B}^\top x_i) - \mathbb{1} [x_i^\top \bar{B} > 0] \hat{\Delta}^\top x \right)^2 \right] \gtrsim \frac{\sigma^2 d}{n_T}. \quad \square$$

G.2.2 Adaptation Upper Bound

Having proven the minimax result for FROZENREP, we now proceed to prove a corresponding upper bound on the performance of ADAPTRREP. To do so, we need to prove a result analogous to Lemma B.4 for the ReLU setting.

Before we proceed to the proof, we need to find proper generalizations for relevant objects in the proof of Lemma B.4. In particular, we recall the prominent use of the projector $P_{\Sigma B_0}$, which, loosely speaking, can be thought of as representing the ‘‘average’’ component of the signal that can be represented by a parameter in $\text{Col } B_0$.

Recall that the inputs (which are element of \mathbb{R}^d) are sampled from a distribution p . This input distribution induces the $L_2(p)$ -norm⁹ on vector-valued functions of \mathbb{R}^d and its associated inner product via

$$\|\zeta\|_{L_2} = \mathbb{E} \left[\|\zeta(x)\|_2^2 \right] \quad \text{and} \quad \langle \zeta, \xi \rangle_{L_2} = \mathbb{E} [\zeta(x)^\top \xi(x)].$$

Now, let $\zeta : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and $\xi : \mathbb{R}^d \rightarrow \mathbb{R}^q$ be two representation functions on \mathbb{R}^d . Then, we can define the *linear projector onto ζ* as the linear operator P_ζ taking representations $\mathbb{R}^d \rightarrow \mathbb{R}^q$ onto itself via

$$[P_\zeta \xi](x) := \mathbb{E} [\xi(x) \zeta(x)^\top] \mathbb{E} [\zeta(x) \zeta(x)^\top]^\dagger \zeta(x).$$

We denote the corresponding orthogonal projection as $P_\zeta^\perp \xi := \xi - P_\zeta \xi$. We note that these operators satisfy the orthogonality property

$$\mathbb{E} [[P_\zeta \xi](x)^\top [P_\zeta^\perp \xi](x)] = 0,$$

which can be easily verified algebraically.

To understand the operator P_ζ , consider the problem of approximating a linear function of ξ via a linear function of ζ , as measured via the input distribution p . More formally, for $v \in \mathbb{R}^q$, we want to find

$$w^* = \underset{w \in \mathbb{R}^p}{\text{argmin}} \left\| \xi(\cdot)^\top v - \zeta(\cdot)^\top w \right\|_{L_2}^2 \quad \text{and} \quad f^* = \zeta(\cdot)^\top w^*.$$

As the problem is differentiable and convex in w , we can simply use standard optimality conditions to find that

$$w^* = \mathbb{E} [\zeta(x) \zeta(x)^\top]^\dagger \mathbb{E} [\zeta(x) \xi(x)^\top] v \quad \text{and} \quad f^* = \zeta(\cdot)^\top w^*.$$

That is, $P_\zeta \xi$ is performing exactly the transformation required on ξ such that $[P_\zeta \xi(\cdot)]^\top v$ is the best approximation to $\xi(\cdot)^\top v$ via linear functions of ζ in L_2 -norm. To further connect this construction to the linear setting, observe that if $\zeta(x) = B_0^\top x$ and $\xi(x) = x$, then for any $\theta, v \in \mathbb{R}^d$,

$$\langle \xi(\cdot)^\top v, [P_\zeta^\perp \xi](\cdot)^\top \theta \rangle_{L_2} = v^\top \Sigma^{1/2} P_{\Sigma B_0}^\perp \Sigma^{1/2} \theta,$$

and thus $P_\zeta \xi$ is indeed the desired generalization of the projection operators used in the proof of Lemma B.4. Having introduced the required mathematical tools, we now prove the corresponding transfer lemma for this setting.

⁹we write L_2 throughout as a shorthand for $L_2(p)$.

Lemma G.4 (Transfer Lemma, ReLU). *Assume that $B_0 = [A_0, -A_0]$. Then, we have that*

$$\frac{1}{T} \left\| P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* W^* \right\|_F^2 \lesssim \frac{1}{T} \sum_{i=1}^T \min_{w, \Delta} \mathbb{E} \left[(x^\top \theta_t^* - \rho_{B_0}(x)[w, \Delta])^2 \right].$$

Proof. This proof follows the outline of Lemma B.4, with all inner products computed with respect to $\langle \cdot, \cdot \rangle_{L_2}$. Throughout the proof, we will write P and P^\perp as shorthand for the operators $P_{\phi_{B_0}}$ and $P_{\phi_{B_0}}^\perp$, respectively. First, we decompose the task-averaged population risk as for any choices of $(w_t), (\Delta_t)$ as

$$\begin{aligned} & \frac{1}{T} \sum_{i=1}^T \left\| (\cdot)^\top \theta_t^* - \rho_{B_0}(\cdot)[w_t, \Delta_t] \right\|_{L_2}^2 \\ & \gtrsim \frac{1}{T} \sum_{t \in [T]} \left\| \phi_{B^*}(\cdot)^\top w_t^* - [P \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\|_{L_2}^2 \\ & \quad + \frac{1}{T} \sum_{t \in [T]} \left\| \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\|_{L_2}^2 \\ & \quad - \frac{1}{T} \sum_{t \in [T]} \left| \left\langle \phi_{B^*}(\cdot)^\top w_t^* - [P \rho_{B_0}](\cdot)[w_t, \Delta_t], \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\rangle_{L_2} \right| \\ & \gtrsim \frac{1}{T} \sum_{t \in [T]} \left\| P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* v_t^* \right\|_2^2 + \frac{1}{T} \sum_{t \in [T]} \left\| \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\|_{L_2}^2 \\ & \quad - \frac{1}{T} \sum_{t \in [T]} \left| \left\langle \phi_{B^*}(\cdot)^\top w_t^* - [P \rho_{B_0}](\cdot)[w_t, \Delta_t], \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\rangle_{L_2} \right|. \end{aligned}$$

The second inequality follows from noting that $\phi_{B^*}(x)^\top w_t^*$ is a linear function of x , and thus by Lemma G.2, the best linear predictor on ϕ is a linear function of x . Since the L_2 -norm on linear functions of x is equivalent to the Σ -norm on the parameters, we obtain the inequality above¹⁰.

Now, we proceed to prove a bound on the inner product above. For any fixed $t \in [T]$,

$$\begin{aligned} & \left| \left\langle \phi_{B^*}(\cdot)^\top w_t^* - [P \rho_{B_0}](\cdot)[w, \Delta], \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w, \Delta] \right\rangle_{L_2} \right| \\ & \leq \left| \left\langle \phi_{B^*}(\cdot)^\top w_t^* - [P \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*], \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*] \right\rangle_{L_2} \right| \\ & \quad + \left| \left\langle \phi_{B^*}(\cdot)^\top w_t^* - [P \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*], [P^\perp \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*] - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\rangle_{L_2} \right| \\ & \quad + \left| \left\langle [P \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*] - [P \rho_{B_0}](\cdot)[w_t, \Delta_t], \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\rangle_{L_2} \right| \\ & \leq \left| \left\langle \phi_{B^*}(\cdot)^\top w_t^* - [P \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*], \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*] \right\rangle_{L_2} \right| \\ & \quad + \left| \left\langle [P^\perp \phi_{B^*}](\cdot)^\top w_t^*, [P^\perp \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*] - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\rangle_{L_2} \right| \\ & \quad + \left| \left\langle [P \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*] - [P \rho_{B_0}](\cdot)[w_t, \Delta_t], \psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t] \right\rangle_{L_2} \right|, \end{aligned}$$

where adding and subtracting $[P \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*]$ and $[P^\perp \rho_{B^*}](\cdot)[w_t^*, \Delta_t^*]$ in the first and second arguments of the inner product, respectively, results in the first inequality. The second inequality then uses the orthogonality properties of P and P^\perp in the second term. Furthermore, observe that

$$\begin{aligned} & \left| \left\langle \phi_{B^*}(x)^\top w_t^* - [P \rho_{B^*}](x)[w_t^*, \Delta_t^*], \psi_{B^*}(x)^\top \Delta_t^* - [P^\perp \rho_{B^*}](x)[w_t^*, \Delta_t^*] \right\rangle_{L_2} \right| \\ & = \left| \left\langle [P^\perp \phi_{B^*}](x)^\top w_t^* - [P \psi_{B^*}](x)^\top \Delta_t^*, [P^\perp \psi_{B^*}](x)^\top \Delta_t^* - [P \phi_{B^*}](x)[w_t^*, \Delta_t^*] \right\rangle_{L_2} \right| \\ & = \left| \left\langle \phi_{B^*}(x)^\top w_t^*, \psi_{B^*}(x)^\top \Delta_t^* \right\rangle_{L_2} \right| = |w_t^* A^* \Sigma \delta_t^*| \\ & = 0, \end{aligned}$$

¹⁰This is exactly the argument used in Lemma B.4.

where we have used the algebraically-verifiable fact that $[P\rho_{B^*}](\cdot)[w, \Delta] = [P\phi_{B^*}](x)^\top + [P\psi_{B^*}](x)[\Delta]$, and that $[P^\perp \rho_{B^*}][w, \Delta]$ decomposes in a similar fashion. Finally, since $\phi_B^*(x)^\top w_t^*$ is linear, we can apply Lemma G.2, and thus

$$\|[P^\perp \phi_{B^*}](x)^\top w_t^*\|_{L_2} = \|P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* v_t^*\|_2.$$

From here, we note that we have an analogous quadratic inequality to that of Lemma B.4 in the terms

$$\frac{1}{T} \sum_{t \in [T]} \|P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* v_t^*\|_2^2 \quad \text{and} \quad \frac{1}{T} \sum_{t \in [T]} \|\psi_{B^*}(\cdot)^\top \Delta_t^* - [P^\perp \rho_{B_0}](\cdot)[w_t, \Delta_t]\|_{L_2}^2$$

upon applying Cauchy-Schwarz as before. Furthermore, note that via orthogonality, we have the Pythagorean identity

$$\|P\rho\|_{L_2}^2 + \|P^\perp \rho\|_{L_2}^2 \leq \|\rho\|_{L_2}^2,$$

and thus by following the exact same algebraic argument as in Lemma B.4 of using Proposition I.2,

$$\begin{aligned} \frac{1}{T} \left\| P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* V^* \right\|_F^2 &= \frac{1}{T} \sum_{t \in [T]} \left\| P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* v_t^* \right\|_2^2 \\ &\leq \frac{1}{T} \sum_{t \in [T]} \mathbb{E} [(x^\top \theta_t^* - \rho_{B_0}(x)[w_t, \Delta_t])^2]. \end{aligned}$$

Since the result holds for any (w_t) and (Δ_t) , it holds for the minimizers. \square

Lemma G.1 (Adaptation target performance, ReLU). *Let $B_0 = [A_0, -A_0]$, and fix a $\theta^* \in S_{A^*}$. Consider a learner which solves*

$$\min_B \min_{\substack{(w_t, \Delta_t) \\ \|\Delta_t\|_F \leq 1}} \frac{1}{n_S} \sum_{t \in [T]} \|\rho_B(X)[w, \Delta] - y\|_2^2.$$

during source training, and

$$\min_{\substack{w, \Delta \\ \|\Delta\|_F \leq 1}} \frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - y\|_2^2$$

during target training, where B_0 is the representation obtained from source training. Then, with probability at least $1 - \delta$,

$$\begin{aligned} \frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - X\theta^*\|_2^2 &\lesssim \left\| P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* \right\|_2^2 + \frac{\sigma^2 k}{n_T} \left(1 + \log \frac{1}{\delta} \right) \\ &\quad + \frac{\sigma}{\sqrt{n_T}} \sqrt{k \operatorname{tr} \Sigma \left(1 + \log \frac{1}{\delta} \right)}. \end{aligned}$$

In particular, with $k = \Theta(1)$ and $\varepsilon = k/d$, and assuming access to infinite per-task samples and tasks during source training, the learner achieves target loss bounded as

$$\frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - X\theta^*\|_2^2 \lesssim \frac{\sigma^2}{n_T} \left(1 + \log \frac{1}{\delta} \right) + \frac{\sigma}{\sqrt{n_T}} \sqrt{1 + \log \frac{1}{\delta}}$$

with probability at least $1 - \delta$ over the draw of target samples.

Proof. Throughout the proof, we instantiate the high-probability event in Lemma B.2, which guarantees that with probability at least $1 - \delta/9$,

$$0.9\Sigma \preceq \frac{1}{n_T} X^\top X \preceq 1.1\Sigma.$$

By forming the least-squares basic inequality, we have that

$$\begin{aligned} & \frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - X(A^*v^* + \delta^*)\|_2^2 \\ & \lesssim \underbrace{\frac{1}{n_T} \|P_{XA_0}^\perp X A^* v^*\|_2^2}_{=:\zeta} + \frac{1}{n_T} |\langle z, \rho_{B_0}(X)[w, \Delta] - P_{XA_0} X A^* v^* - X \delta^* \rangle|. \end{aligned}$$

Then, since $P_{XA_0} X A^* v^*$ is in the span of $\phi_{B_0}(X)$, we can bound the right-hand side of the basic inequality by

$$\begin{aligned} & \zeta + \frac{1}{n_T} \left| \left\langle z, \underbrace{P_{\phi_{B_0}(X)} (\rho_{B_0}(X)[w, \Delta] - P_{XA_0} X A^* v^* - X \delta^*)}_{=:T_{\text{low-rank}}} \right\rangle \right| \\ & + \frac{1}{n_T} \left| \left\langle z, P_{\phi_{B_0}(X)}^\perp (\psi_{B_0}(X)[\Delta] - X \delta^*) \right\rangle \right|. \end{aligned}$$

Furthermore, we can lower bound the left-hand side by

$$\frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - X(A^*v^* + \delta^*)\|_2^2 \geq \frac{1}{n_T} \|T_{\text{low-rank}}\|_2^2 + \frac{2}{n_T} \langle P_{XA_0}^\perp X A^* v^*, T_{\text{low-rank}} \rangle.$$

Note that $\phi_{B_0}(X)$ is a matrix with $\text{rank} \leq 2k$, and therefore, $P_{\phi_{B_0}(X)} z$ is a chi-squared random variable, with at most $2k$ degrees of freedom. Subsequently, by applying known bounds on chi-squared random variables, we have that with probability at least $1 - (4/9)\delta$,

$$\begin{aligned} \frac{1}{n_T} \|T_{\text{low-rank}}\|_2^2 & \leq \zeta + \left(\sqrt{\zeta} + \frac{2\sigma k}{\sqrt{n_T}} \sqrt{1 + \log \frac{1}{\delta}} \right) \frac{1}{\sqrt{n_T}} \|T_{\text{low-rank}}\|_2 \\ & + \frac{1}{n_T} \left| \left\langle z, P_{\phi_{B_0}(X)}^\perp (\psi_{B_0}(X)[\Delta] - X \delta^*) \right\rangle \right|, \end{aligned}$$

which thus implies via Proposition I.2 that

$$\begin{aligned} & \frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - X(A^*v^* + \delta^*)\|_2^2 \\ & \lesssim \zeta + \frac{\sigma^2 k}{n_T} \left(1 + \log \frac{1}{\delta} \right) + \frac{1}{n_T} \left| \left\langle z, P_{\phi_{B_0}(X)}^\perp (\psi_{B_0}(X)[\Delta] - X \delta^*) \right\rangle \right|. \end{aligned}$$

To bound the final term, note that

$$\begin{aligned} & \frac{1}{n_T} \left| \left\langle z, P_{\phi_{B_0}(X)}^\perp (\psi_{B_0}(X)[\Delta] - X \delta^*) \right\rangle \right| \\ & \leq \frac{1}{\sqrt{n_T}} \left[\left\| \frac{1}{\sqrt{n_T}} \psi_{B_0}(X)^\top P_{\phi_{B_0}(X)}^\perp z \right\|_2 + \left\| \frac{1}{\sqrt{n_T}} X^\top P_{\phi_{B_0}(X)}^\perp z \right\|_2 \right], \end{aligned}$$

by applying the norm bounds on Δ and δ^* , and thus by applying the Hanson-Wright inequality, with probability at least $1 - (4/9)\delta$,

$$\begin{aligned} & \frac{1}{n_T} \left| \left\langle z, P_{\phi_{B_0}(X)}^\perp (\psi_{B_0}(X)[\Delta] - X \delta^*) \right\rangle \right| \\ & \leq \frac{\sigma}{\sqrt{n_T}} \left[\sqrt{\frac{1}{n_T} \text{tr} \psi_{B_0}(X)^\top \psi_{B_0}(X)} + \sqrt{\frac{1}{n_T} \text{tr} X^\top X} \right] \sqrt{1 + \log \frac{1}{\delta}} \\ & \lesssim \frac{\sigma}{\sqrt{n_T}} \sqrt{k \text{tr} \Sigma \left(1 + \log \frac{1}{\delta} \right)}. \end{aligned}$$

Therefore, since $\|v^*\|_2 \leq 1$, we have an overall bound of

$$\begin{aligned} \frac{1}{n_T} \|\rho_{B_0}(X)[w, \Delta] - X \theta^*\|_2^2 & \lesssim \left\| P_{\Sigma^{1/2} A_0}^\perp \Sigma^{1/2} A^* \right\|_2^2 + \frac{\sigma^2 k}{n_T} \left(1 + \log \frac{1}{\delta} \right) \\ & + \frac{\sigma}{\sqrt{n_T}} \sqrt{k \text{tr} \Sigma \left(1 + \log \frac{1}{\delta} \right)}. \end{aligned}$$

To prove the second part of the statement, note that with infinite source tasks and samples, global optimality with respect to the source loss together with Lemma G.4 implies that

$$\frac{1}{T} \left\| P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* V^* \right\|_F^2 = 0 \implies \left\| P_{\Sigma^{1/2}A_0}^\perp \Sigma^{1/2} A^* \right\|_2 = 0,$$

since $(k/T) [V^{*\top} V^*] = I_k$. Furthermore, with the choice of d and ε , $\text{tr } \Sigma = \Theta(1)$. Therefore, we obtain the second claim. \square

H A Performance Bound for Projected Gradient Descent

In this section, we provide a performance bound for projected gradient descent on the objective

$$\mathcal{L}(\theta) := \frac{1}{n} \sum_{i \in [n]} \ell(g_\theta(x_i), y_i)$$

for $\theta \in \Theta$, where $\Theta \subseteq \mathbb{R}^q$. We assume that Θ is norm-bounded by D , and that Θ is convex and contains 0. Furthermore, we assume that g_θ is twice-differentiable as a function of θ .

Key to the performance bound that we will demonstrate is that g_θ is “approximately linear” in the parameter θ , which we formally define below. Under this assumption, we demonstrate \mathcal{L} is approximately convex over Θ if \mathcal{L} is Lipschitz as a function of the vector of predictions and ℓ is convex in the first argument. Therefore, with slight modifications to the online analysis of projected gradient descent, we obtain the desired performance bound.

Following the discussion above, we make the following assumptions:

Assumption H.1 (Approximate linearity). *There exists β and L such that*

$$\sup_{\theta \in \Theta} \frac{1}{n} \sum_{i \in [n]} \|\nabla_{\theta}^2 g_\theta(x_i)\|_2^2 \leq \beta^2 \quad \text{and} \quad \frac{1}{n} \sum_{i \in [n]} \|\nabla_{\theta} g_0(x_i)\|_2^2 \leq L^2.$$

Assumption H.2 (Assumptions on ℓ). *We assume that ℓ is convex in the first argument. Furthermore, if \mathcal{L} is viewed as a function of the vector of predictions $g_\theta(X)$, then we have that*

$$\|\nabla_g \mathcal{L}(g_\theta(X))\|_2^2 \leq \frac{\alpha^2}{n}$$

for any $\theta \in \Theta$, i.e. \mathcal{L} is (α/\sqrt{n}) -Lipschitz as a function of the vector of predictions.

Note that we abuse notation in Assumption H.2, using \mathcal{L} to reference both the function of the parameter, and of the vector of predictions. Given these two assumptions, we now proceed to demonstrate that the loss landscape of \mathcal{L} has several desirable properties.

Lemma H.1 (Approximate convexity in Θ). *Let $\theta_1, \theta_2 \in \Theta$. Then,*

$$\langle \nabla_{\theta} \mathcal{L}(\theta_1), \theta_2 - \theta_1 \rangle \leq \mathcal{L}(\theta_2) - \mathcal{L}(\theta_1) + 4\alpha\beta D^2.$$

Proof. Note that \mathcal{L} is a convex function of the vector of predictions $g_\theta(X)$, as it is a sum of convex functions by Assumption H.2. Therefore,

$$\mathcal{L}(\theta_2) - \mathcal{L}(\theta_1) \geq \langle \nabla_g \mathcal{L}(g_{\theta_1}(X)), g_{\theta_2}(X) - g_{\theta_1}(X) \rangle.$$

Furthermore, by the chain rule,

$$\langle \nabla_{\theta} \mathcal{L}(\theta_1), \theta_2 - \theta_1 \rangle = \langle \nabla_g \mathcal{L}(g_{\theta_1}(X)), [\nabla_{\theta} g_{\theta_1}(X)](\theta_2 - \theta_1) \rangle.$$

Putting the two statements together and applying Cauchy-Schwarz,

$$\begin{aligned} \mathcal{L}(\theta_2) - \mathcal{L}(\theta_1) - \langle \nabla_{\theta} \mathcal{L}(\theta_1), \theta_2 - \theta_1 \rangle &\geq \langle \nabla_g \mathcal{L}(g_{\theta_1}(X)), g_{\theta_2}(X) - g_{\theta_1}(X) - [\nabla_{\theta} g_{\theta_1}(X)](\theta_2 - \theta_1) \rangle \\ &\geq -\|\nabla_g \mathcal{L}(g_{\theta_1}(X))\|_2 \|g_{\theta_2}(X) - g_{\theta_1}(X) - [\nabla_{\theta} g_{\theta_1}(X)](\theta_2 - \theta_1)\|_2 \\ &\geq -\frac{\alpha}{\sqrt{n}} \|g_{\theta_2}(X) - g_{\theta_1}(X) - [\nabla_{\theta} g_{\theta_1}(X)](\theta_2 - \theta_1)\|_2 \end{aligned}$$

Now, by Taylor’s theorem, there exists $\bar{\theta} = \lambda\theta_1 + (1-\lambda)\theta_2$ such that

$$g_{\theta_2}(X) - g_{\theta_1}(X) - \nabla_{\theta} g_{\theta_1}(X)(\theta_2 - \theta_1) = [(\theta_2 - \theta_1)^\top \nabla_{\theta}^2 g_{\bar{\theta}}(x_i)(\theta_2 - \theta_1)]_{i \in [n]}.$$

Consequently, by rearranging,

$$\begin{aligned} \langle \nabla_{\theta} \mathcal{L}(\theta_1), \theta_2 - \theta_1 \rangle &\leq \mathcal{L}(\theta_2) - \mathcal{L}(\theta_1) + \frac{\alpha}{\sqrt{n}} \sqrt{\sum_{i \in [n]} [(\theta_2 - \theta_1)^\top \nabla_{\theta}^2 g_{\bar{\theta}}(x_i)(\theta_2 - \theta_1)]^2} \\ &\leq \mathcal{L}(\theta_2) - \mathcal{L}(\theta_1) + 4\alpha\beta D^2. \end{aligned}$$

□

Lemma H.2 (Gradient bound in Θ). *For any $\theta \in \Theta$, we have that*

$$\|\nabla_{\theta} \mathcal{L}(\theta)\|_2^2 \lesssim \alpha^2 (L^2 + \beta^2 D^2)$$

Proof. Using the Lipschitz assumption in Assumption H.2,

$$\begin{aligned} \|\nabla_{\theta} \mathcal{L}(\theta)\|_2^2 &= \|[\nabla_{\theta} g_{\theta}(X)]^{\top} \nabla_g \mathcal{L}(g_{\theta}(X))\|_2^2 \leq \|\nabla_{\theta} g_{\theta}(X)\|_F^2 \|\nabla_g \mathcal{L}(g_{\theta}(X))\|_2^2 \\ &\leq \frac{\alpha^2}{n} \sum_{i \in [n]} \|\nabla_{\theta} g_{\theta}(x_i)\|_2^2 \end{aligned}$$

Now, by integrating, we have that for any fixed $i \in [n]$,

$$\nabla_{\theta} g_{\theta}(x_i) = \nabla_{\theta} g_0(x_i) + \left[\int_0^1 \nabla_{\theta}^2 g_{\alpha\theta}(x_i) d\alpha \right] \theta.$$

Therefore,

$$\begin{aligned} \|\nabla_{\theta} \mathcal{L}(\theta)\|_2^2 &\lesssim \frac{\alpha^2}{n} \sum_{i \in [n]} \|\nabla_{\theta} g_0(x_i)\|_2^2 + \frac{\alpha^2}{n} \left\| \int_0^1 \nabla_{\theta}^2 g_{\alpha\theta}(x_i) d\alpha \right\|_2^2 \|\theta\|_2^2 \\ &\leq \frac{\alpha^2}{n} \sum_{i \in [n]} \|\nabla_{\theta} g_{\theta_0}(x_i)\|_2^2 + \frac{\alpha^2 D^2}{n} \int_0^1 \frac{1}{n} \sum_{i \in [n]} \|\nabla_{\theta}^2 g_{\alpha\theta}(x_i)\|_2^2 d\alpha \\ &\leq \alpha^2 (L^2 + \beta^2 D^2). \end{aligned} \quad \square$$

Having proven the results above, we now proceed to the main claim of this section.

Theorem H.1 (Bound on PGD performance). *Assume we run projected gradient descent on \mathcal{L} with constraint set Θ for T_{PGD} iterations with step size η given by*

$$\eta = \frac{1}{\sqrt{T_{\text{PGD}}}} \left(\frac{D}{\alpha \sqrt{L^2 + \beta^2 D^2}} \right).$$

Let $(\theta_t)_{t \in [T_{\text{PGD}}]}$ denote the sequence of PGD iterates obtained, where $\theta_0 = 0$. Then, for any $\theta \in \Theta$,

$$\min_t \mathcal{L}(\theta_t) - \mathcal{L}(\theta) \lesssim \alpha \beta D^2 + \alpha D \sqrt{\frac{L^2 + \beta^2 D^2}{T_{\text{PGD}}}}.$$

Proof. For any $t \in [T_{\text{PGD}}]$,

$$\begin{aligned} r_{t+1}^2 &= \|\theta_{t+1} - \theta\|_2^2 \\ &\leq \|\theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t) - \theta\|_2^2 \\ &= r_t^2 + 2\eta \langle \nabla_{\theta} \mathcal{L}(\theta_t), \theta - \theta_t \rangle + \eta^2 \|\nabla_{\theta} \mathcal{L}(\theta_t)\|_2^2 \\ &\leq r_t^2 + 2\eta [\mathcal{L}(\theta_0 + \delta) - \mathcal{L}(\theta_0 + \delta_t) + 4\alpha\beta D^2] + \eta^2 \alpha^2 [L^2 + \beta^2 D^2], \end{aligned}$$

where the first inequality follows from the nonexpansive property of projections onto convex sets. Furthermore, the last inequality makes use of the approximate convexity property from Lemma H.1 and the gradient bound over Θ from Lemma H.2. Therefore, via telescoping,

$$r_{T_{\text{PGD}}}^2 \leq r_0^2 + 2\eta \left[\sum_{t=0}^{T_{\text{PGD}}-1} \mathcal{L}(\theta) - \mathcal{L}(\theta_t) \right] + 8\eta T_{\text{PGD}} \alpha \beta D^2 + \eta^2 T_{\text{PGD}} \alpha^2 [L^2 + \beta^2 D^2],$$

or by rearranging,

$$\begin{aligned} \frac{r_0^2 - r_{T_{\text{PGD}}}^2}{2\eta T_{\text{PGD}}} + 4\alpha\beta D^2 + \frac{\eta\alpha^2}{2} [L^2 + \beta^2 D^2] &\geq \frac{1}{T_{\text{PGD}}} \sum_{t=0}^{T_{\text{PGD}}-1} \mathcal{L}(\theta_0 + \delta_t) - \mathcal{L}(\theta_0 + \delta) \\ &\geq \min_{t=0, \dots, T_{\text{PGD}}-1} \mathcal{L}(\theta_0 + \delta_t) - \mathcal{L}(\theta_0 + \delta). \end{aligned}$$

Now, using the choice of step size, observe that

$$\begin{aligned} \frac{r_0^2 - r_{T_{\text{PGD}}}^2}{2\eta T_{\text{PGD}}} + \frac{\eta}{2} [\alpha^2 L^2 + \beta^2 D^2] &\leq \frac{D^2}{2\eta T_{\text{PGD}}} + \frac{\eta\alpha^2}{2} [L^2 + \beta^2 D^2] \\ &= \alpha D \sqrt{\frac{L^2 + \beta^2 D^2}{T_{\text{PGD}}}}, \end{aligned}$$

from which the desired claim easily follows. □

I Technical Lemmas

Proposition I.1 (Du et al. (2020), Lemma A.5). *Let $\mathcal{O}^{d_1 \times d_2}$ be the set of matrices in $\mathbb{R}^{d_1 \times d_2}$ with orthonormal columns, $d_1 \geq d_2$. Then, there exists an ε -covering of $\mathcal{O}^{d_1 \times d_2}$ with at most $(6\sqrt{d_2}/\varepsilon)^{d_1 d_2}$ elements.*

Proposition I.2 (Solving quadratic inequalities). *Assume that $ax^2 \leq bx + c$ for $a, b, c > 0$. Then, $bx + c \lesssim (b^2/a) + c$.*

Proof. Since $a > 0$, the solution set to the inequality is given by the interval $[r_1, r_2]$, where r_1 and r_2 are the roots of $ax^2 - bx - c$. By the quadratic formula, the larger root r_2 is given by

$$r_2 = \frac{b + \sqrt{b^2 + 4ac}}{2a} \leq \frac{b}{a} + \sqrt{\frac{c}{a}}.$$

Therefore,

$$x \leq r_2 \leq \frac{b}{a} + \sqrt{\frac{c}{a}} \implies bx + c \leq \frac{b^2}{a} + \left(\frac{b}{\sqrt{a}}\right) \sqrt{c} + c \lesssim \frac{b^2}{a} + c,$$

where the last inequality makes use of the Cauchy-Schwarz inequality. \square

Corollary I.1. *Let X and Y be random variables.*

$$\mathbb{E}[X^2] \lesssim \mathbb{E}[(X + Y)^2] + \mathbb{E}[Y^2]$$

Proof. We have that

$$\mathbb{E}[(X + Y)^2] = \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2] \geq \mathbb{E}[X^2] - 2\mathbb{E}[|XY|].$$

Therefore, by applying Cauchy-Schwarz,

$$\mathbb{E}[X^2] \leq \mathbb{E}[(X + Y)^2] + 2\mathbb{E}[|XY|] \leq \mathbb{E}[(X + Y)^2] + 2\sqrt{\mathbb{E}[X^2]\mathbb{E}[Y^2]}.$$

Finally, by applying Proposition I.2,

$$\mathbb{E}[X^2] \lesssim \mathbb{E}[(X + Y)^2] + \mathbb{E}[Y^2]. \quad \square$$

Proposition I.3. *Let A, B be matrices with compatible dimensions, and assume that $\text{rank } A = r > 0$. Then,*

$$\|P_A B\|_F \leq \frac{1}{\sigma_r(A)} \|A^\top B\|_F.$$

Proof. Let (U, Σ, V^\top) be the compact singular value decomposition of A , i.e. we only retain positive singular values in Σ . By rotational invariance,

$$\|A^\top B\|_F^2 = \|V \Sigma U^\top B\|_F^2 = \|\Sigma U^\top B\|_F^2.$$

Furthermore, by definition,

$$\|\Sigma U^\top B\|_F^2 = \sum_i \|\Sigma U^\top B e_i\|_2^2 \geq \sigma_r^2(\Sigma) \sum_i \|U^\top B e_i\|_2^2 = \sigma_r^2(A) \|U^\top B\|_F^2.$$

Finally, by applying rotational invariance once more,

$$\|A^\top B\|_F^2 \geq \sigma_r^2(A) \|U U^\top B\|_F^2 = \sigma_r^2(A) \|P_A B\|_F^2,$$

from which the desired claim follows. \square

Proposition I.4. *Let $\lambda, \gamma > 0$, and fix a vector y . Then,*

$$\min_{\substack{A, x \\ Ax=y}} \frac{\lambda}{2} \|A\|_F^2 + \frac{\gamma}{2} \|x\|_2^2 = \sqrt{\lambda\gamma} \|y\|_2.$$

Proof. We proceed by cases. If $y = 0$, then the result is trivial.

Otherwise, if $y \neq 0$, note that $x^* \neq 0$. Now, for any fixed $x \neq 0$, the minimizing choice for A is zx^\top for some z . To see this, observe that if A is not rank-1, then we can achieve a lower Frobenius norm by reducing its rank. Consequently, for a given x , the minimizing choice for A is $yx^\top / \|x\|_2^2$ necessarily. Therefore,

$$\min_{\substack{A, x \\ Ax=y}} \frac{\lambda}{2} \|A\|_F^2 + \frac{\gamma}{2} \|x\|_2^2 = \min_x \frac{\lambda}{2} \left(\frac{\|y\|_2^2}{\|x\|_2^2} \right) + \frac{\gamma}{2} \|x\|_2^2 = \min_{z>0} \frac{\lambda}{2} \left(\frac{\|y\|_2^2}{z} \right) + \frac{\gamma z}{2}.$$

This final optimization problem is convex in z – using first-order optimality conditions, we can thus easily see that $z^* = \sqrt{\lambda/\gamma} \|y\|_2$, and therefore

$$\min_{\substack{A, x \\ Ax=y}} \frac{\lambda}{2} \|A\|_F^2 + \frac{\gamma}{2} \|x\|_2^2 = \sqrt{\lambda\gamma} \|y\|_2. \quad \square$$

Proposition I.5 (Expectation bound on empirical spectral norm). *Let $X \in \mathbb{R}^{n \times d}$ be a matrix with rows drawn i.i.d. from a zero-mean distribution with covariance Σ . Furthermore, assume that the whitened distribution is ρ^2 -sub-Gaussian. Then, whenever $n \gtrsim \rho^4 d$,*

$$\mathbb{E} \left[\lambda_{\max} \left(\frac{X^\top X}{n} \right) \right] \lesssim \|\Sigma\|_2.$$

Proof. By Weyl's inequality,

$$\mathbb{E} \left[\lambda_{\max} \left(\frac{X^\top X}{n} \right) \right] \leq \|\Sigma\|_2 + \mathbb{E} \left[\left| \lambda_{\max} \left(\frac{X^\top X}{n} \right) - \Sigma \right| \right] \leq \|\Sigma\|_2 + \mathbb{E} \left[\left\| \frac{X^\top X}{n} - \Sigma \right\|_2 \right].$$

Thus, by applying the result in Vershynin (2017, Theorem 4.4.1), we have that as long as $n_S \gtrsim \rho^4 d$,

$$\mathbb{E} \left[\lambda_{\max} \left(\frac{X^\top X}{n} \right) \right] \lesssim \|\Sigma\|_2. \quad \square$$

Proposition I.6 (Gaussian complexity chain rule, Tripuraneni et al. (2020b), Theorem 7). *Assume that \mathcal{F} is a class of functions $\mathbb{R}^k \rightarrow \mathbb{R}$ such that every $f \in \mathcal{F}$ is L -Lipschitz in the L_2 -norm. Furthermore, assume that Φ is a class of functions $\mathbb{R}^d \rightarrow \mathbb{R}^k$ such that for any $\phi \in \Phi$, $\phi(x)$ is norm-bounded by D for any x in the support of the input distribution. Then, we have the bound*

$$\frac{1}{T} \mathcal{G}_n(\mathcal{F}^{\otimes T} \circ \Phi) \leq \frac{8D}{(nT)^2} + 128 \left(\frac{L}{T} \mathcal{G}_n(\Phi) + \mathbb{E} \left[\sup_{Z \in \mathcal{Z}} \mathcal{G}_Z(\mathcal{F}) \right] \right) \log(nT),$$

where \mathcal{Z} is the random set $\{(\phi(x_{i_1}), \dots, \phi(x_{i_n})) \mid i_1, \dots, i_n \in [nT]\}$ and $\mathcal{G}_Z(\mathcal{F})$ is the empirical Gaussian complexity on samples Z . Note that the inner expectation is over the nT input samples, and that we have assumed that all input samples come from a single distribution.

Proposition I.7 (Tripuraneni et al. (2020b), Lemma 6). *Let $h, h^* : \mathbb{R}^d \rightarrow \mathbb{R}^k$ be representation functions, and define*

$$\Lambda(h, h^*) := \mathbb{E} [h^*(x)h^*(x)^\top] - \mathbb{E} [h^*(x)h(x)^\top] (\mathbb{E} [h(x)h(x)^\top])^\dagger \mathbb{E} [h(x)h^*(x)^\top].$$

Then, $\inf_v \mathbb{E} [(h(x)^\top v - h^(x)^\top v^*)^2] = (v^*)^\top \Lambda(h, h^*) v^*$. Furthermore, if*

$$\sigma_{\min}(\mathbb{E} [h(x)h(x)^\top]) \geq c_1 > 0 \quad \text{and} \quad \sigma_{\max}(\mathbb{E} [h^*(x)h^*(x)^\top]) \leq c_2,$$

then this infimum is achieved within the ball of radius $\|v^\|_2 \sqrt{c_2/c_1}$.*

Proof. The calculation of the infimum is provided in Tripuraneni et al. (2020b), and is thus omitted. However, we prove the sharper radius bound below.

Define $F_{h, h^*} := \mathbb{E} [h(x)h'(x)^\top]$, so that $\Lambda(h, h^*) = F_{h^*, h^*} - F_{h^*, h} F_{h, h}^\dagger F_{h, h^*}$. Then, since $\Lambda(h, h^*) \succeq 0$, and recalling that the infimum is achieved at $v = F_{h, h}^\dagger F_{h, h^*} v^*$,

$$\left\| F_{h, h}^\dagger F_{h, h^*} v^* \right\|_2^2 \leq \frac{1}{c_1} \left\| F_{h, h}^{1/2} F_{h, h}^\dagger F_{h, h^*} v^* \right\|_2^2 \leq \frac{1}{c_1} \left\| F_{h^*, h^*}^{1/2} v^* \right\|_2^2 \leq \frac{c_1}{c_2} \|v^*\|_2^2,$$

as desired. \square