

## A. Appendix Overview

The Appendix is structured as follows

- Section B proves the result for linear representation learning with the tr-tr objective, Theorem 5.1, that shows that “bad” full rank solutions exist arbitrarily close to the optimal value of the tr-tr objective  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}$ . The proof works by arguing that smaller values of  $\lambda$  are desirable for the tr-tr objective, and this can be simulated by effectively making the norm of the representation layer very high. Furthermore a higher rank representation is preferable over lower rank ones to fit the noise in the labels in the training data better.
- Section C proves the main result for linear representation learning with tr-val objective, Theorem 5.4, that proves that the optimal solutions to the tr-val objective  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2))$  for most  $n_1$  and  $\sigma$  will be *low-rank* representations that are also *expressive enough*. The result is also extended to solutions that are  $\tau$ -optimal in the  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}$  objective for a small enough  $\tau$ . The crux of the proof for this result is in Theorem C.1 that provides a closed form expression for the tr-val objective that disentangles the expressivity and the low-rankness of the representation.
- Section D presents additional experimental details and results, including results for the MiniImageNet dataset.

## B. More on Train-Train split

### B.1. Proof of main result

**Theorem 5.1.** *For every  $\lambda, n > 0$ , for every  $\tau > 0$ , there exists a “bad” representation layer  $\mathbf{A}_{\text{bad}} \in \mathbb{R}^{d \times d}$  that satisfies  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}_{\text{bad}}; n) \leq \inf_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}; n) + \tau$ , but has the following lower bound on meta-testing loss*

$$\inf_{\bar{\lambda} > 0} \mathcal{L}_{\bar{\lambda}, \text{rep}}^{\text{test}}(\mathbf{A}_{\text{bad}}; \bar{n}_1) \geq \sigma^2 + \min \left\{ 1 - \frac{\bar{n}_1}{d(1 + \sigma^2)}, \frac{d\sigma^2}{(\bar{n}_1 + d\sigma^2)} \right\}$$

*Proof.* For most of the proof, we will leave out the  $n$  in the expression for  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}$ , i.e. we will denote the tr-tr loss as  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ . We prove this result by using Lemma 5.2 first, and later prove this lemma.

**Lemma 5.2.** *For every  $\lambda > 0$  and  $\mathbf{A} \in \mathbb{R}^{d \times d}$  with rank  $r$ ,*

$$\begin{aligned} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}; n) &\geq \lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}; n) \geq \sigma^2 \frac{(n-r)_+}{n} \\ &\& \lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{I}_d; n) = \sigma^2 \frac{(n-d)_+}{n} \end{aligned}$$

This tells us that for any matrix  $\mathbf{A}$ ,  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) \geq \sigma^2 \frac{(n - \text{rank}(\mathbf{A}))_+}{n} \geq \sigma^2 \frac{(n-d)_+}{n}$ . Also the lower bound of  $\sigma^2 \frac{(n-d)_+}{n}$  can be achieved by  $\kappa \mathbf{I}_d$  in the limit of  $\kappa \rightarrow \infty$ , thus  $\sigma^2 \frac{(n-d)_+}{n} = \inf_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ . Also since  $\lim_{\kappa \rightarrow \infty} \kappa \mathbf{I}_d = \sigma^2 \frac{(n-d)_+}{n} = \inf_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ , for a large enough  $\kappa = \kappa(\tau)$ ,  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa(\tau) \mathbf{I}_d)$  can be made lesser than  $\inf_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) + \tau$ . Thus we pick  $\mathbf{A}_{\text{bad}} = \kappa(\tau) \mathbf{I}_d$ . With this choice, the new task is essentially linear regression in  $d$ -dimension with isotropic Gaussians.

To show that  $\mathbf{A}_{\text{bad}}$  is indeed bad, we will use the lower bound for ridge regression on isotropic Gaussian linear regression from Theorem 4.2(a) in Saunshi et al. (2020). They show that the excess risk for  $\mathbf{I}_d$  (and thus  $\kappa(\tau) \mathbf{I}_d$ ), regardless of the choice of regularizer  $\bar{\lambda}$ , for a new task  $\rho_v$  will be lower bounded by

$$\inf_{\bar{\lambda} > 0} \mathbb{E}_{S \sim \rho_v^{\bar{n}_1}} [\|\mathbf{A} \mathbf{w}_{\bar{\lambda}}(\mathbf{A}; S) - \mathbf{v}\|^2] \geq \begin{cases} \frac{d \|\mathbf{v}\|^2 \sigma^2}{\bar{n}_1 \|\mathbf{v}\| + \sigma^2 d} & \text{if } \bar{n}_1 \geq d \\ \frac{\bar{n}_1}{d} \frac{\|\mathbf{v}\|^2 \sigma^2}{\|\mathbf{v}\| + \sigma^2} + \frac{d - \bar{n}_1}{d} \|\mathbf{v}\|^2 & \text{if } \bar{n}_1 < d \end{cases}$$

Their proof can be easily modified to replace  $\|\mathbf{v}\|^2$  with  $\mathbb{E}_{\mathbf{v} \sim \bar{\mu}} \|\mathbf{v}\|^2$ . The lower bound can be simplified for the the  $\bar{n}_1 < d$  case to  $\frac{\bar{n}_1}{d} \frac{\|\mathbf{v}\|^2 \sigma^2}{\|\mathbf{v}\| + \sigma^2} + \frac{d - \bar{n}_1}{d} \|\mathbf{v}\|^2 = \|\mathbf{v}\|^2 - \frac{\bar{n}_1}{d} \frac{\|\mathbf{v}\|^2}{\|\mathbf{v}\|^2 + \sigma^2}$ . Plugging in  $\|\mathbf{v}\| = 1$  completes the proof.  $\square$

We now prove the lemma

**Lemma 5.2.** For every  $\lambda > 0$  and  $\mathbf{A} \in \mathbb{R}^{d \times d}$  with rank  $r$ ,

$$\begin{aligned} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}; n) &\geq \lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}; n) \geq \sigma^2 \frac{(n-r)_+}{n} \\ &\& \lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{I}_d; n) = \sigma^2 \frac{(n-d)_+}{n} \end{aligned}$$

*Proof.* We first prove that having  $\lambda = 0$  will lead to the smallest loss  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$  for every  $\mathbf{A}$ . We then observe that  $\lambda = 0$  can be simulated by increasing the norm of  $\mathbf{A}$ . These claims mathematically mean that, (a)  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) \geq \mathcal{L}_{\lambda', \text{rep}}^{\text{tr-tr}}(\mathbf{A})$  whenever  $\lambda \geq \lambda'$  and (b)  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}) = \mathcal{L}_{\frac{\lambda}{\kappa}, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ . This will give us that  $\lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}) = \lim_{\lambda \rightarrow 0} \mathcal{L}_{0, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}) \geq \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ .

Intuitively,  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}$  is trying to learn a linear classifier on top of data that is linear transformed by  $\mathbf{A}$  with the goal of fitting the same data well.

**Lemma B.1.** For any representation layer  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and  $\lambda > 0$ , we have the following

$$\lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}) = \lim_{\lambda \rightarrow 0} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) \leq \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$$

Fitting the data is better when there is less restriction on the norm of the classifier, which in this case means when  $\lambda$  is smaller. Furthermore, increasing the norm of the representation layer  $\mathbf{A}$  effectively reduces the impact the regularizer will have. We first prove this lemma later, first we use it to prove Lemma 5.2 that shows that the loss for low rank matrices will be high.

Lemma B.1 already shows that  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}; n) \geq \lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}; n)$ . Also using Lemma B.1, we can replace  $\lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A})$  with  $\lim_{\lambda \rightarrow 0} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ . Using Equation (9) and from central limit theorem, we have

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) = \lim_{T \rightarrow \infty} \widehat{\mathcal{L}}^{\text{tr-tr}}(\mathbf{A}) = \mathbb{E}_{\rho_{\mathbf{v}} \sim \mu} \left[ \mathbb{E}_{S \sim \rho_{\mathbf{v}}^n} \left[ \frac{1}{n} \|\mathbf{X} \mathbf{A} \mathbf{w}_{\lambda}(\mathbf{A}; S) - \mathbf{Y}\|^2 \right] \right] \quad (14)$$

$$= \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{A}^* \mathbf{A}^* \top)} \left[ \mathbb{E}_{S \sim \rho_{\mathbf{v}}^n} \left[ \frac{1}{n} \|\mathbf{X} \mathbf{A} \mathbf{w}_{\lambda}(\mathbf{A}; S) - \mathbf{Y}\|^2 \right] \right] \quad (15)$$

This is because  $\widehat{\mathcal{L}}^{\text{tr-tr}}$  is an average loss for  $T$  train tasks, and the limit when  $T \rightarrow \infty$  it converges to the expectation over the task distribution  $\mu$ . We first observe that  $S \sim \rho_{\mathbf{v}}$  is equivalent to sampling  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ ,  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)$  which gives us  $\mathbf{Y} = \mathbf{X} \mathbf{v} + \boldsymbol{\eta}$ , where  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $\boldsymbol{\eta} \in \mathbb{R}^n$ ,  $\mathbf{Y} \in \mathbb{R}^n$ . Using the definition of  $\mathbf{w}_{\lambda}(\mathbf{A}; S)$  from Equation (8) the standard KKT condition for linear regression, we can write a closed form solution for

$$\mathbf{w}_{\lambda}(\mathbf{A}; S) = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \|\mathbf{X} \mathbf{A} \mathbf{w} - \mathbf{Y}\|^2 + \lambda \|\mathbf{w}\|^2 = \left( \frac{\mathbf{A}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{A}}{n} + \lambda \mathbf{I}_d \right)^{-1} \frac{\mathbf{A}^{\top} \mathbf{X}^{\top}}{n} \mathbf{Y} \quad (16)$$

$$\lim_{\lambda \rightarrow 0} \mathbf{w}_{\lambda}(\mathbf{A}; S) = \left( \frac{\mathbf{A}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{A}}{n} \right)^{\dagger} \frac{\mathbf{A}^{\top} \mathbf{X}^{\top}}{n} \mathbf{Y} \quad (17)$$

where the last step is folklore that the limit of ridge regression as regularization coefficient goes to 0 is the minimum  $\ell_2$ -norm linear regression solution. Plugging this into Equation (20) and taking the limit

$$\begin{aligned} \lim_{\lambda \rightarrow 0} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) &= \mathbb{E}_{\mathbf{v}, S} \left[ \frac{1}{n} \left\| \mathbf{X} \mathbf{A} \left( \lim_{\lambda \rightarrow 0} \mathbf{w}_{\lambda}(\mathbf{A}; S) \right) - \mathbf{Y} \right\|^2 \right] \\ &= \mathbb{E}_{\mathbf{v}, S} \left[ \frac{1}{n} \left\| \left( \mathbf{X} \mathbf{A} \left( \frac{\mathbf{A}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{A}}{n} \right)^{\dagger} \frac{\mathbf{A}^{\top} \mathbf{X}^{\top}}{n} - \mathbf{I}_n \right) \mathbf{Y} \right\|^2 \right] \\ &= \frac{1}{n} \mathbb{E}_{\mathbf{v}, S} \left[ \left\| P_{\mathbf{X} \mathbf{A}}^{\perp} \mathbf{Y} \right\|^2 \right] \end{aligned}$$

where for any matrix  $\mathbf{B} \in \mathbb{R}^{n \times d}$ , we denote  $P_{\mathbf{B}} \in \mathbb{R}^{n \times n} = \mathbf{B}\mathbf{B}^\dagger$  to denote the projection matrix onto the span of columns of  $\mathbf{B}$ , and  $P_{\mathbf{B}}^\perp = I_n - P_{\mathbf{B}}$  is the projection matrix onto the orthogonal subspace. Note that if  $\text{rank}(\mathbf{B}) = n$ , then  $P_{\mathbf{B}} = I_n$  and  $P_{\mathbf{B}}^\perp = 0$ . Thus the error incurred is the amount of the label  $\mathbf{Y}$  that the representation  $\mathbf{X}\mathbf{A}$  cannot predict linearly. We further decompose this

$$\begin{aligned} \lim_{\lambda \rightarrow 0} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) &= \frac{1}{n} \mathbb{E}_{\mathbf{v}, S} \left[ \left\| P_{\mathbf{X}\mathbf{A}}^\perp \mathbf{Y} \right\|^2 \right] = \frac{1}{n} \mathbb{E}_{\mathbf{v}, S} \left[ \left\| P_{\mathbf{X}\mathbf{A}}^\perp \mathbf{X}\mathbf{v} + P_{\mathbf{X}\mathbf{A}}^\perp \boldsymbol{\eta} \right\|^2 \right] \\ &= \underbrace{\frac{1}{n} \mathbb{E}_{\mathbf{v}, S} \left[ \left\| P_{\mathbf{X}\mathbf{A}}^\perp \mathbf{X}\mathbf{v} \right\|^2 \right]}_{\text{fitting signal}} + \underbrace{\frac{1}{n} \mathbb{E}_{\mathbf{v}, S} \left[ \left\| P_{\mathbf{X}\mathbf{A}}^\perp \boldsymbol{\eta} \right\|^2 \right]}_{\text{fitting noise}} + \underbrace{\frac{2}{n} \mathbb{E}_{\mathbf{v}, S} \left[ \boldsymbol{\eta}^\top P_{\mathbf{X}\mathbf{A}}^\perp \mathbf{X}\mathbf{v} \right]}_{\text{cross term}} \end{aligned} \quad (18)$$

$$\begin{aligned} &= \frac{1}{n} \mathbb{E}_S \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{A}^* \mathbf{A}^{*\top})} \left[ \left\| P_{\mathbf{X}\mathbf{A}}^\perp \mathbf{X}\mathbf{v} \right\|^2 \right] + \mathbb{E}_S \frac{1}{n} \text{tr} \left( P_{\mathbf{X}\mathbf{A}}^\perp \mathbb{E}_{\boldsymbol{\eta}} \left[ \boldsymbol{\eta} \boldsymbol{\eta}^\top \right] P_{\mathbf{X}\mathbf{A}}^\perp \right) \\ &= \underbrace{\frac{1}{n} \mathbb{E}_S \left[ \left\| P_{\mathbf{X}\mathbf{A}}^\perp \mathbf{X}\mathbf{A}^* \right\|^2 \right]}_{\alpha(\mathbf{A})} + \underbrace{\mathbb{E}_S \frac{\sigma^2}{n} \text{tr} \left( P_{\mathbf{X}\mathbf{A}}^\perp \right)}_{\beta(\mathbf{A})} \end{aligned} \quad (19)$$

We first note, due to independence of  $\boldsymbol{\eta}$  and  $\mathbf{X}$  that the cross term in Equation (18) is 0 in expectation. Using the distributions for  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{A}^* \mathbf{A}^{*\top})$  and  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 I_n)$ , we get the final expressions. Firstly, we note that  $\alpha(\mathbf{A}) \geq 0$  for every  $\mathbf{A}$ , and a sufficient condition for  $\alpha(\mathbf{A}) = 0$  is that  $\mathbf{A}^*$  lies in the span of  $\mathbf{A}$ , i.e.  $P_{\mathbf{A}}^\perp \mathbf{A}^* = 0$ . More importantly, it is clear that  $\alpha(\mathbf{I}_d) = 0$ . Next we look at the  $\beta(\mathbf{A})$  term which is proportional to the trace of  $P_{\mathbf{X}\mathbf{A}}^\perp$  which, for a projection matrix, is also equal to the rank of the matrix. Note that since the rank of  $\mathbf{X}\mathbf{A} \in \mathbb{R}^{d \times n}$  is at most  $\min\{n, d\}$ . Thus we get

$$\beta(\mathbf{A}) = \sigma^2 \frac{\text{rank}(P_{\mathbf{X}\mathbf{A}}^\perp)}{n} = \sigma^2 \frac{(n - \text{rank}(P_{\mathbf{X}\mathbf{A}}))}{n} \geq \sigma^2 \frac{(n - \text{rank}(\mathbf{A}))_+}{n}$$

where  $(x)_+ = \mathbb{1}_{x \geq 0} x$ . This along with  $\alpha(\mathbf{A}) \geq 0$  proves the first part of the result. For the second part  $\alpha(\mathbf{I}_d) = 0$  along with noticing that  $\text{tr}(P_{\mathbf{X}\mathbf{I}_d}^\perp) = n - \text{rank}(\mathbf{X}) = n - \min\{n, d\}$  since a Gaussian matrix is full rank with measure 1, thus giving us  $\beta(\mathbf{I}_d) = \sigma^2 \frac{(n-d)_+}{n}$  and completing the proof.  $\square$

Thus Lemma 5.2 shows that  $\lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{I}_d) = \inf_{\mathbf{A}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$  and so picking a large enough  $\kappa(\tau)$  can always give us  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa(\tau) \mathbf{I}_d) \leq \inf_{\mathbf{A}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) + \tau$  which completes the proof of Theorem 5.1. The only thing left to prove is Lemma B.1 which we do below

*Proof of Lemma B.1.* Using Equation (16) we get the following expression for  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}$

$$\begin{aligned} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) &= \mathbb{E}_{\mathbf{v}, S} \left[ \frac{1}{n} \left\| \mathbf{X}\mathbf{A}\mathbf{w}_\lambda(\mathbf{A}; S) - \mathbf{Y} \right\|^2 \right] \\ &= \mathbb{E}_{\mathbf{v}, S} \left[ \frac{1}{n} \left\| \left( \mathbf{X}\mathbf{A} \left( \frac{\mathbf{A}^\top \mathbf{X}^\top \mathbf{X}\mathbf{A}}{n} + \lambda \mathbf{I}_d \right)^{-1} \frac{\mathbf{A}^\top \mathbf{X}^\top}{n} - \mathbf{I}_n \right) \mathbf{Y} \right\|^2 \right] \end{aligned} \quad (20)$$

From this we get

$$\begin{aligned} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}) &= \mathbb{E}_{\mathbf{v}, S} \left[ \frac{1}{n} \left\| \left( \kappa \mathbf{X}\mathbf{A} \left( \frac{\kappa^2 \mathbf{A}^\top \mathbf{X}^\top \mathbf{X}\mathbf{A}}{n} + \lambda \mathbf{I}_d \right)^{-1} \frac{\kappa \mathbf{A}^\top \mathbf{X}^\top}{n} - \mathbf{I}_n \right) \mathbf{Y} \right\|^2 \right] \\ &= \mathbb{E}_{\mathbf{v}, S} \left[ \frac{1}{n} \left\| \left( \mathbf{X}\mathbf{A} \left( \frac{\mathbf{A}^\top \mathbf{X}^\top \mathbf{X}\mathbf{A}}{n} + \frac{\lambda}{\kappa^2} \mathbf{I}_d \right)^{-1} \frac{\mathbf{A}^\top \mathbf{X}^\top}{n} - \mathbf{I}_n \right) \mathbf{Y} \right\|^2 \right] \\ &= \mathcal{L}_{\frac{\lambda}{\kappa^2}, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) \end{aligned}$$

Thus  $\lim_{\kappa \rightarrow \infty} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\kappa \mathbf{A}) = \lim_{\kappa \rightarrow \infty} \mathcal{L}_{\frac{\lambda}{\kappa^2}, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) = \lim_{\lambda \rightarrow 0} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ . Note that we have used the fact that  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}$  is continuous in  $\lambda$  for  $\lambda > 0$ . To prove the remaining result, i.e.  $\lim_{\lambda \rightarrow 0} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A}) \leq \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$ , we just need to prove that  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-tr}}(\mathbf{A})$  is a

increasing function of  $\lambda$ . Suppose  $\frac{\mathbf{X}\mathbf{A}}{\sqrt{n}} = \mathbf{U}_\mathbf{X}\mathbf{S}_\mathbf{X}\mathbf{V}_\mathbf{X}^\top$  is the singular value decomposition. For any  $\lambda' < \lambda$ , we can rewrite  $\mathcal{L}_{\lambda,\text{rep}}^{\text{tr-tr}}(\mathbf{A})$  from Equation (20) as follows

$$\begin{aligned}
 \mathcal{L}_{\lambda,\text{rep}}^{\text{tr-tr}}(\mathbf{A}) &= \mathbb{E}_{\mathbf{v},\mathbf{S}} \left[ \frac{1}{n} \left\| \left( \mathbf{X}\mathbf{A} \left( \frac{\mathbf{A}^\top \mathbf{X}^\top \mathbf{X}\mathbf{A}}{n} + \lambda \mathbf{I}_d \right)^{-1} \frac{\mathbf{A}^\top \mathbf{X}^\top}{n} - \mathbf{I}_n \right) \mathbf{Y} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{v},\mathbf{S}} \left[ \frac{1}{n} \left\| \left( \mathbf{V}_\mathbf{X} \mathbf{S}_\mathbf{X} \mathbf{U}_\mathbf{X}^\top \left( \mathbf{U}_\mathbf{X} \mathbf{S}_\mathbf{X} \mathbf{V}_\mathbf{X}^\top \mathbf{V}_\mathbf{X} \mathbf{S}_\mathbf{X} \mathbf{U}_\mathbf{X}^\top + \lambda \mathbf{I}_d \right)^{-1} \mathbf{U}_\mathbf{X} \mathbf{S}_\mathbf{X} \mathbf{V}_\mathbf{X}^\top - \mathbf{I}_n \right) \mathbf{Y} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{v},\mathbf{S}} \left[ \frac{1}{n} \left\| \mathbf{V}_\mathbf{X} \left( \mathbf{S}_\mathbf{X} \left( \mathbf{S}_\mathbf{X}^2 + \lambda \mathbf{I}_d \right)^{-1} \mathbf{S}_\mathbf{X} - \mathbf{I}_n \right) \mathbf{V}_\mathbf{X}^\top \mathbf{Y} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{v},\mathbf{S}} \left[ \frac{1}{n} \sum_{i=1}^d \left[ \left( \frac{\mathbf{S}_\mathbf{X}[i]^2}{\lambda + \mathbf{S}_\mathbf{X}[i]^2} - 1 \right)^2 \mathbf{V}_\mathbf{X}^\top \mathbf{Y}[i]^2 \right] \right] \\
 &= \mathbb{E}_{\mathbf{v},\mathbf{S}} \left[ \frac{1}{n} \sum_{i=1}^d \left[ \left( \frac{\lambda}{\lambda + \mathbf{S}_\mathbf{X}[i]^2} \right)^2 \mathbf{V}_\mathbf{X}^\top \mathbf{Y}[i]^2 \right] \right] \\
 &\geq \mathbb{E}_{\mathbf{v},\mathbf{S}} \left[ \frac{1}{n} \sum_{i=1}^d \left[ \left( \frac{\lambda'}{\lambda' + \mathbf{S}_\mathbf{X}[i]^2} \right)^2 \mathbf{V}_\mathbf{X}^\top \mathbf{Y}[i]^2 \right] \right] \\
 &= \mathcal{L}_{\lambda',\text{rep}}^{\text{tr-tr}}(\mathbf{A})
 \end{aligned}$$

where the only inequality in the above sequence follows from the observation that  $\frac{\lambda}{\lambda+a}$  is an increasing function for  $a, \lambda > 0$ . This completes the proof.  $\square$

## C. More on Train-Validation split

### C.1. Proof of main results

We first prove the result that for  $n_1 = \bar{n}_1$  and  $\lambda = \bar{\lambda}$ ,  $\mathcal{L}_{\lambda,\text{rep}}^{\text{tr-val}} \equiv \mathcal{L}_{\bar{\lambda},\text{rep}}^{\text{test}}$ .

**Proposition 5.3.**  $\mathcal{L}_{\lambda,\text{rep}}^{\text{tr-val}}(\cdot; (n_1, n_2))$  and  $\mathcal{L}_{\lambda,\text{rep}}^{\text{test}}(\cdot; \bar{n}_1)$  are equivalent if  $\bar{n}_1 = n_1$  and  $\bar{\lambda} = \lambda$

*Proof.* We again note using the central limit theorem and Equation (9) that

$$\begin{aligned}
 \mathcal{L}_{\lambda,\text{rep}}^{\text{tr-val}}(\mathbf{A}) &= \lim_{T \rightarrow \infty} \widehat{\mathcal{L}}_{\lambda,\text{rep}}^{\text{tr-val}}(\mathbf{A}) = \mathbb{E}_{\rho_{\mathbf{v}} \sim \bar{\mu}} \left[ \mathbb{E}_{(S^{\text{tr}}, S^{\text{val}}) \sim \rho_{\mathbf{v}}^n} \left[ \frac{1}{n_2} \left\| \mathbf{X}^{\text{val}} \mathbf{A} \mathbf{w}_\lambda(\mathbf{A}; S^{\text{tr}}) - \mathbf{Y}^{\text{val}} \right\|^2 \right] \right] \\
 &\stackrel{(a)}{=} \mathbb{E}_{\rho_{\mathbf{v}} \sim \bar{\mu}} \left[ \mathbb{E}_{S^{\text{tr}} \sim \rho_{\mathbf{v}}^{n_1}} \left[ \mathbb{E}_{S^{\text{val}} \sim \rho_{\mathbf{v}}^{n_2}} \left[ \frac{1}{n_2} \left\| \mathbf{X}^{\text{val}} \mathbf{A} \mathbf{w}_\lambda(\mathbf{A}; S^{\text{tr}}) - \mathbf{Y}^{\text{val}} \right\|^2 \right] \right] \right] \\
 &= \mathbb{E}_{\rho_{\mathbf{v}} \sim \bar{\mu}} \left[ \mathbb{E}_{S^{\text{tr}} \sim \rho_{\mathbf{v}}^{n_1}} \left[ \mathbb{E}_{(\mathbf{x}, y) \sim \rho_{\mathbf{v}}} \left( \mathbf{x}^\top \mathbf{A} \mathbf{w}_\lambda(\mathbf{A}; S^{\text{tr}}) - y \right)^2 \right] \right] \\
 &= \mathbb{E}_{\rho_{\mathbf{v}} \sim \bar{\mu}} \left[ \mathbb{E}_{S^{\text{tr}} \sim \rho_{\mathbf{v}}^{n_1}} \left[ \left\| \mathbf{A} \mathbf{w}_\lambda(\mathbf{A}; S^{\text{tr}}) - \mathbf{v} \right\|^2 \right] \right] + \sigma^2 \\
 &\stackrel{(b)}{=} \mathcal{L}_{\lambda,\text{rep}}^{\text{test}}(\mathbf{A}; n_1)
 \end{aligned}$$

where (a) follows by noticing that sample  $S \sim \rho_{\mathbf{v}}^n$  and splitting randomly into  $S^{\text{tr}}$  and  $S^{\text{val}}$  is equivalent to independently sampling  $S^{\text{tr}} \sim \rho_{\mathbf{v}}^{n_1}$  and  $S^{\text{val}} \sim \rho_{\mathbf{v}}^{n_2}$  and (b) follows from the definition of  $\mathcal{L}^{\text{test}}$  from Equation (11).  $\square$

We now prove the main benefit of the tr-val split: learning of low rank linear representations. For that, we use this general result that computes closed form expression for  $\mathcal{L}_{\lambda,\text{rep}}^{\text{tr-val}}$ .

**Theorem C.1.** Let  $\lambda = 0$ . For a first layer  $\mathbf{A} \in \mathbb{R}^{d \times d}$  with  $r = \text{rank}(\mathbf{A})$ , let  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$  be the SVD, where  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$  and  $\mathbf{S} \in \mathbb{R}^{r \times r}$ . Furthermore let  $P_{\mathbf{U}} = \mathbf{U}\mathbf{U}^\top$  denote the span of columns of  $\mathbf{U}$  (and thus  $\mathbf{A}$ ) and let  $P_{\mathbf{U}}^\perp = I_d - P_{\mathbf{U}}$ . Then the tr-val objective has the following form

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 = \begin{cases} (1 + \alpha(n_1, r)) \|P_{\mathbf{U}}^\perp \mathbf{A}^*\|^2 + \alpha(n_1, r)\sigma^2 & \text{if } r < n_1 - 1 \\ \left(\frac{n_1}{r} + \alpha(r, n_1)\right) \|P_{\mathbf{U}}^\perp \mathbf{A}^*\|^2 + \frac{r-n_1}{r} + \beta(\mathbf{S}) + \alpha(r, n_1)\sigma^2 & \text{if } r > n_1 + 1 \\ \infty & \text{otherwise} \end{cases}$$

where  $\beta(\mathbf{S}) \geq 0$  and  $\beta(\mathbf{S}) = 0$  when  $\mathbf{S} = \kappa \mathbf{I}_r$  for some  $\kappa > 0$ . Also  $\alpha$  is defined as

$$\alpha(a, b) = \frac{b}{a - b - 1} \quad (21)$$

We prove this theorem in Section C.2 First we prove the main result by using this theorem. We note that similar results can be shown for different regimes of  $k, d, n_1$  and  $\sigma$ . The result below is for one reasonable regime where  $n_1 = \Omega(k)$  and  $\sigma = \mathcal{O}(1)$ . Similar results can be obtained if we further assume that  $k \ll d$  with weaker conditions on  $\sigma$ , we leave that for future work. Furthermore, this result can be extended to  $\tau$ -optimal solutions to the tr-val objective rather than just the optimal solution.

**Theorem 5.4.** Let  $\lambda = 0$ . Suppose  $n_1 \geq 2k + 2$  and  $\sigma^2 \in (0, \frac{n_1 - k - 1}{3k})$ , then any optimal solution  $\mathbf{A}_{\text{best}} \in \arg \min_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2))$  will satisfy

$$\text{rank}(\mathbf{A}_{\text{best}}) = k, \quad P_{\mathbf{A}_{\text{best}}} \mathbf{A}^* = \mathbf{A}^*, \quad \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}_{\text{best}}; (n_1, n_2)) - \sigma^2 = \sigma^2 \frac{k}{n_1 - k - 1}$$

where  $P_{\mathbf{A}} = \mathbf{A}\mathbf{A}^\dagger$  is the projection matrix onto the columnspace of  $\mathbf{A}$ . For any matrix  $\mathbf{A}_{\text{good}} \in \mathbb{R}^{d \times d}$  that is  $\tau$ -optimal, i.e. it satisfies  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}_{\text{good}}; (n_1, n_2)) \leq \inf_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) + \tau$  for  $\tau \in (0, \frac{\sigma^2}{n_1 - k - 1})$ , then we have

$$\text{rank}(\mathbf{A}_{\text{good}}) = k, \quad \|P_{\mathbf{A}_{\text{good}}} \mathbf{A}^* - \mathbf{A}^*\|^2 \leq \tau$$

The meta-testing performance of  $\mathbf{A}_{\text{good}}$  on a new task with  $\bar{n}_1 > 2k + 2$  samples satisfies

$$\inf_{\bar{\lambda} \geq 0} \mathcal{L}_{\lambda, \text{rep}}^{\text{test}}(\mathbf{A}_{\text{good}}; \bar{n}_1) - \sigma^2 \leq 2\tau + \sigma^2 \frac{2k}{\bar{n}_1}$$

*Proof.* Suppose the optimal value of  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2))$  is  $\mathcal{L}^*$ , i.e.

$$\mathcal{L}^* = \inf_{\mathbf{A} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) \quad (22)$$

Let  $\mathbf{A}_{\text{good}} \in \mathbb{R}^{d \times d}$  be the “good” matrix that is  $\tau$ -optimal, i.e.  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}_{\text{good}}; (n_1, n_2)) - \mathcal{L}^* \leq \tau$ . Note that the result for  $\mathbf{A}_{\text{best}}$  follows from the result for  $\tau = 0$ .

We use the expression for  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\cdot; (n_1, n_2))$  from Theorem C.1 to find properties for  $\mathbf{A}_{\text{good}}$  that can ensure that it is  $\tau$ -optimal. Consider a representation  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and let  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$  be its SVD with  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$  and  $\mathbf{S} \in \mathbb{R}^{r \times r}$  where  $r = \text{rank}(\mathbf{A})$ . We consider 3 cases for  $r$ :  $r < k$ ,  $k \leq r \leq n_1$  and  $r > n_1$ , and find properties that can result in  $\mathbf{A}$  being  $\tau$ -optimal in each of the 3 cases. For the ranges of  $n_1, \sigma, \tau$  in the theorem statement, it will turn out that when  $r < k$  or  $r > n_1$ ,  $\mathbf{A}$  must satisfy  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) > \mathcal{L}^* + \tau$ , thus concluding that  $\text{rank}(\mathbf{A}_{\text{good}})$  cannot be in these ranges.

To do so we analyze the optimal value of  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\cdot; (n_1, n_2))$  that can be achieved for all three cases. We first start analyzing the most promising case of  $k \leq r \leq n_1$ .

**Case 1:**  $k \leq r \leq n_1$  For this case we can use the  $r < n_1 - 1$  regime from Theorem C.1. Note that for  $r \in \{n_1 - 1, n_1\}$ ,  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2))$  is unbounded. Thus we get

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 = (1 + \alpha(n_1, r)) \|P_{\mathcal{U}}^{\perp} \mathbf{A}^*\|^2 + \alpha(n_1, r) \sigma^2 \quad (23)$$

where  $\alpha(n_1, r)$  is defined in Equation (34). Note that  $\|P_{\mathcal{U}}^{\perp} \mathbf{A}^*\| \geq 0$  and in fact equality can be achieved for every  $r$  in this case since  $r = \text{rank}(\mathbf{U}) \geq k = \text{rank}(\mathbf{A}^*)$ . Furthermore  $\alpha(n_1, r) = \frac{r}{n_1 - r - 1}$  is an increasing function of  $r$ , and thus  $\alpha(n_1, r) \geq \alpha(n_1, k)$ , which can also be achieved by picking  $r = k$ . Thus for the range of  $k \leq r \leq n_1$ , we get

$$\mathcal{L}^* - \sigma^2 \leq \min_{\substack{\mathbf{A} \text{ s.t.} \\ k \leq r \leq n_1}} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 = \sigma^2 \alpha(n_1, k) = \sigma^2 \frac{k}{n_1 - k - 1} \quad (24)$$

with the minimum being achieved when  $\text{rank}(\mathbf{A}) = k$  and  $\|P_{\mathcal{U}}^{\perp} \mathbf{A}^*\|$ , which is the same as  $P_{\mathcal{U}} \mathbf{A}^* = \mathbf{A}^*$ .

Suppose  $r_{\text{good}} = \text{rank}(\mathbf{A}_{\text{good}})$  lies in this range, i.e.  $k \leq r_{\text{good}} \leq n_1$ . Using the fact that  $\mathbf{A}_{\text{good}}$  is  $\tau$ -optimal, we can show an upper bound the  $r_{\text{good}}$  and  $\|P_{\mathbf{A}_{\text{good}}} \mathbf{A}^* - \mathbf{A}^*\|$ . From the  $\tau$ -optimality condition and Equation (24), we get

$$\tau \geq \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}_{\text{good}}; (n_1, n_2)) - \mathcal{L}^* \geq \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}_{\text{good}}; (n_1, n_2)) - \sigma^2 \alpha(n_1, k) \quad (25)$$

$$\stackrel{(a)}{\geq} (1 + \alpha(n_1, r)) \|P_{\mathcal{U}}^{\perp} \mathbf{A}^*\|^2 + \sigma^2 \alpha(n_1, r_{\text{good}}) - \sigma^2 \alpha(n_1, k) \quad (26)$$

$$= (1 + \alpha(n_1, r)) \|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|^2 + \sigma^2 (\alpha(n_1, r_{\text{good}}) - \alpha(n_1, k)) \quad (27)$$

$$\stackrel{(b)}{=} \|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|^2 + \sigma^2 \alpha'(n_1, r')(r_{\text{good}} - k) = \|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|^2 + \sigma^2 \frac{r'(n_1 - 1)}{(n_1 - r' - 1)^2} (r_{\text{good}} - k) \quad (28)$$

$$\stackrel{(c)}{\geq} \|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|^2 + \sigma^2 \frac{k(n_1 - 1)}{(n_1 - k - 1)^2} (r_{\text{good}} - k) \geq \|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|^2 + \sigma^2 \frac{k}{(n_1 - k - 1)} (r_{\text{good}} - k) \quad (29)$$

where (a) follows from Equation (23) instantiated for  $\mathbf{A}_{\text{good}}$ , (b) follows from the mean value theorem applied to the continuous function  $\alpha(n_1, \cdot)$  with  $r' \in [k, r_{\text{good}}]$ , (c) follows by  $k > r'$  and the monotonicity of  $\alpha'(n_1, \cdot)$ . Thus from the inequality in Equation (29), we can get an upper bound on  $r_{\text{good}}$  as follows

$$\tau \geq \sigma^2 \frac{k}{(n_1 - k - 1)} (r_{\text{good}} - k) \implies r_{\text{good}} \leq k + \frac{\tau(n_1 - k - 1)}{\sigma^2 k} \stackrel{(a)}{\leq} k + \frac{1}{k} \quad (30)$$

$$\implies r_{\text{good}} = k \quad (31)$$

where (a) follows from the upper bound on  $\tau$  in the theorem statement. We can also get an upper bound on  $\|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|$  using Equation (29) as  $\|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|^2 \leq \tau$ . Thus any  $\mathbf{A}_{\text{good}}$  with  $k \leq r_{\text{good}} \leq n_1$  that is  $\tau$ -optimal will satisfy  $r_{\text{good}} = k$  and  $\|P_{\mathbf{A}_{\text{good}}}^{\perp} \mathbf{A}^*\|^2 \leq \tau$ .

Additionally the minimum achievable value for this range of rank is  $\sigma^2 \frac{k}{n_1 - k - 1}$  from Equation (24). We now analyze the other two cases, where we will show that no  $\mathbf{A}$  can even be  $\tau$ -optimal.

**Case 2:**  $r < k$  In this case, since  $k < n_1$ , we are still in the  $r < n_1$  regime. The key point here is that when  $r < k$ , it is impossible for  $\mathbf{A}$  to span all of  $\mathbf{A}^*$ . In fact for rank  $r$ , it is clear that  $\mathbf{A}$  can cover only at most  $r$  out of  $k$  directions from  $\mathbf{A}^*$ . Thus the inexpressiveness term  $\|P_{\mathcal{U}}^{\perp} \mathbf{A}^*\|^2$  will be at least  $\frac{k-r}{k} \|\mathbf{A}^*\|^2 = \frac{k-r}{k}$ . Using the  $r < n_1$  expression from Theorem C.1, we get

$$\begin{aligned} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 &= (1 + \alpha(n_1, r)) \|P_{\mathcal{U}}^{\perp} \mathbf{A}^*\|^2 + \alpha(n_1, r) \sigma^2 \\ &\stackrel{(a)}{\geq} \frac{k-r}{k} + \sigma^2 \frac{r}{n_1 - r - 1} \\ &= \frac{k-r}{k} + \sigma^2 \frac{r}{n_1 - r - 1} - \sigma^2 \frac{k}{n_1 - r - 1} + \sigma^2 \frac{k}{n_1 - r - 1} - \sigma^2 \frac{k}{n_1 - k - 1} + \sigma^2 \frac{k}{n_1 - k - 1} \\ &= \frac{k-r}{k} - \sigma^2 \frac{k-r}{n_1 - r - 1} - \sigma^2 k \left( \frac{1}{n_1 - k - 1} - \frac{1}{n_1 - r - 1} \right) + \sigma^2 \frac{k}{n_1 - k - 1} \\ &= \frac{k-r}{k} - \sigma^2 \frac{k-r}{n_1 - r - 1} - \sigma^2 \frac{k(k-r)}{(n_1 - r - 1)(n_1 - k - 1)} + \sigma^2 \frac{k}{n_1 - k - 1} \end{aligned}$$

$$\begin{aligned}
&\geq^{(b)} \frac{k-r}{k} - \sigma^2 \frac{k-r}{n_1-k-1} - \sigma^2 \frac{(k-r)}{(n_1-k-1)} + \sigma^2 \frac{k}{n_1-k-1} \\
&= \frac{k-r}{k} - \sigma^2 \frac{2(k-r)}{n_1-k-1} + \sigma^2 \frac{k}{n_1-k-1}
\end{aligned} \tag{32}$$

where for (a) we use  $\alpha(n_1, r) \geq 0$  and for (b), we use that  $k < n_1 - k - 1 < n_1 - r - 1$ . Since we assume  $\sigma^2 < (n_1 - k - 1)/2k$ , the difference between the first 2 terms is at least 0. Thus the error when  $r < k$  is at least  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 > \sigma^2 \frac{k}{n_1 - k - 1}$ , which is larger than the previous case. So the optimal solution cannot have  $r < k$ .

We now check if such an  $\mathbf{A}$  can be a  $\tau$ -optimal solution instead. The answer is negative due to the following calculation

$$\begin{aligned}
\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \mathcal{L}^* &\geq^{(a)} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 \frac{k}{n_1 - k - 1} \\
&\geq^{(b)} \frac{k-r}{k} - \sigma^2 \frac{2(k-r)}{n_1 - k - 1} \\
&\geq^{(c)} \frac{3k\sigma^2}{n_1 - k - 1} \frac{k-r}{k} - \sigma^2 \frac{2(k-r)}{n_1 - k - 1} \\
&= \frac{\sigma^2(k-r)}{n_1 - k - 1} \geq \frac{\sigma^2}{n_1 - k - 1} > \tau
\end{aligned}$$

where (a) follows from Equation (24), (b) follows from Equation (32), (c) follows from  $1 \geq \frac{3k\sigma^2}{n_1 - k - 1}$  from the condition on  $\sigma$  in the theorem statement and the last inequalities follow from  $r < k$  condition and the range of  $\tau$  in the theorem statement respectively. Thus we cannot even have a  $\tau$ -optimal solution with  $r < k$ . Next we show the same for the case of  $r > n_1$ .

**Case 3:  $r > n_1$**  For this case we can use the  $r > n_1 + 1$  regime from Theorem C.1. Note that for  $r \in \{n_1 + 1, n_1\}$ ,  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2))$  is unbounded. Thus we have  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 = (\frac{n_1}{r} + \alpha(r, n_1)) \|P_{\mathbf{U}}^\perp \mathbf{A}^*\|^2 + \frac{r-n_1}{r} + \beta(\mathbf{S}) + \sigma^2 \alpha(r, n_1)$ . We again note that since  $r > n_1 > k$ , we can make  $\|P_{\mathbf{U}}^\perp \mathbf{A}^*\|^2 = 0$  by simply picking an “expressive” rank- $r$  subspace  $\mathbf{U}$ . Further more  $\beta(\mathbf{S}) = 0$  is easy to achieve by making  $\mathbf{S} = \mathbf{I}_r$  which can be done independently of  $\mathbf{U}$ . Thus we can lower bound the error as follows

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 \geq \frac{r-n_1}{r} + \sigma^2 \alpha(r, n_1) = 1 - \frac{n_1}{r} + \sigma^2 \frac{n_1}{r-n_1-1} \tag{33}$$

where equality can be achieved by letting columns of  $\mathbf{A}$  span the subspace  $\mathbf{A}^*$  and by picking  $\mathbf{S} = \mathbf{I}_r$ . We now lower bound this value further to show that we cannot have a  $\tau$ -optimal solution.

Firstly we note that if  $\sigma \geq 1$ , then Equation (33) gives us

$$\begin{aligned}
\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 &\geq 1 - \frac{n_1}{r} + \frac{n_1}{r-n_1-1} \geq 1 \geq^{(a)} \sigma^2 \frac{3k}{n_1-k-1} \\
&\geq^{(b)} (\mathcal{L}^* - \sigma^2) + \sigma^2 \frac{2k}{n_1-k-1} \geq^{(c)} (\mathcal{L}^* - \sigma^2) + \tau(n_1-k-1) \frac{2k}{n_1-k-1} \\
&> (\mathcal{L}^* - \sigma^2) + \tau
\end{aligned}$$

where (a) follows from the upper bound on  $\sigma^2$  from the theorem statement, (b) follows from Equation (24) and (c) follows from the upper bound on  $\tau$  from the theorem statement. Thus we cannot have a  $\tau$ -optimal solution in this case.

If  $\sigma < 1$  on the other hand, we can lower bound the error as

$$\begin{aligned}
\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 &= 1 - \frac{n_1}{r} + \sigma^2 \frac{n_1}{r-n_1-1} > 1 - \frac{n_1}{r} + \sigma^2 \frac{n_1}{r-n_1} \\
&\geq \min_{r > n_1} 1 - \frac{n_1}{r} + \sigma^2 \frac{n_1}{r-n_1}
\end{aligned}$$

Setting the derivate w.r.t.  $r$  to 0 for the above expression, we get that this is minimized at  $r = n_1/(1-\sigma)$ . Plugging in this value and simplifying, we get

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 \geq 1 - (1-\sigma)^2 = 2\sigma - \sigma^2 \geq \sigma^2$$

$$\begin{aligned}
 &= \sigma^2 \frac{n_1 - 2k - 1}{n_1 - k - 1} + \sigma^2 \frac{k}{n_1 - k - 1} \stackrel{(a)}{>} \frac{\sigma^2}{n_1 - k - 1} + (\mathcal{L}^* - \sigma^2) \\
 &\stackrel{(b)}{>} \tau + (\mathcal{L}^* - \sigma^2)
 \end{aligned}$$

where (a) follows from the condition that  $n_1 > 2k + 2$  and (b) follows again from Equation (24). Thus the  $r > n_1$  case can never give us a  $\tau$ -optimal solution, let alone the optimal solution and so  $\mathbf{A}_{\text{good}}$  cannot have rank  $r_{\text{good}} > n_1$  either. Combining all 3 cases, we can conclude that any  $\tau$ -optimal solution  $\mathbf{A}_{\text{good}}$  will satisfy  $\text{rank}(\mathbf{A}_{\text{good}}) = k$  and  $\|P_{\mathbf{A}_{\text{good}}}^\perp \mathbf{A}^*\|^2 = \|P_{\mathbf{A}_{\text{good}}} \mathbf{A}^* - \mathbf{A}^*\|^2 \leq \tau$ .

For the second part, we use Proposition 5.3 to first get that for  $\lambda = \bar{\lambda} = 0$ ,  $\mathcal{L}_{\bar{\lambda}, \text{rep}}^{\text{test}}(\mathbf{A}_{\text{good}}; \bar{n}_1) = \mathcal{L}_{\bar{\lambda}, \text{rep}}^{\text{tr-val}}(\mathbf{A}_{\text{good}}; (\bar{n}_1, \cdot))$ . Since  $\bar{n}_1 > 2k + 2$  and  $\text{rank}(\mathbf{A}_{\text{good}}) = k < \bar{n}_1$ , we can use Theorem C.1 to get

$$\begin{aligned}
 \mathcal{L}_{\bar{\lambda}, \text{rep}}^{\text{test}}(\mathbf{A}_{\text{good}}; \bar{n}_1) - \sigma^2 &= \mathcal{L}_{\bar{\lambda}, \text{rep}}^{\text{tr-val}}(\mathbf{A}_{\text{good}}; (\bar{n}_1, \cdot)) - \sigma^2 = (1 + \alpha(\bar{n}_1, k)) \|P_{\mathbf{A}_{\text{good}}}^\perp \mathbf{A}^*\|^2 + \alpha(\bar{n}_1, k) \sigma^2 \\
 &\leq^{(a)} 2 \|P_{\mathbf{A}_{\text{good}}}^\perp \mathbf{A}^*\|^2 + \sigma^2 \frac{k}{\bar{n}_1 - k - 1} \leq^{(b)} 2\tau + \sigma^2 \frac{k}{\bar{n}_1 - \bar{n}_1/2} = 2\tau + \sigma^2 \frac{2k}{\bar{n}_1}
 \end{aligned}$$

where (a) follows by noting that  $\alpha(\bar{n}_1, k) < 1$  and (b) follows from  $k + 1 < \bar{n}_1/2$ . This completes the proof of result for  $\tau$ -optimal solution  $\mathbf{A}_{\text{good}}$ . Setting  $\tau = 0$  gives results for optimal solution  $\mathbf{A}_{\text{best}}$  of  $\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\cdot, (n_1, n_2))$ .  $\square$

## C.2. Proof of Theorem C.1

We now prove the crucial result that gives a closed form solution for the tr-val objective

**Theorem C.1.** *Let  $\lambda = 0$ . For a first layer  $\mathbf{A} \in \mathbb{R}^{d \times d}$  with  $r = \text{rank}(\mathbf{A})$ , let  $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$  be the SVD, where  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$  and  $\mathbf{S} \in \mathbb{R}^{r \times r}$ . Furthermore let  $P_{\mathbf{U}} = \mathbf{U} \mathbf{U}^\top$  denote the span of columns of  $\mathbf{U}$  (and thus  $\mathbf{A}$ ) and let  $P_{\mathbf{U}}^\perp = I_d - P_{\mathbf{U}}$ . Then the tr-val objective has the following form*

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 = \begin{cases} (1 + \alpha(n_1, r)) \|P_{\mathbf{U}}^\perp \mathbf{A}^*\|^2 + \alpha(n_1, r) \sigma^2 & \text{if } r < n_1 - 1 \\ \left(\frac{n_1}{r} + \alpha(r, n_1)\right) \|P_{\mathbf{U}}^\perp \mathbf{A}^*\|^2 + \frac{r - n_1}{r} + \beta(\mathbf{S}) + \alpha(r, n_1) \sigma^2 & \text{if } r > n_1 + 1 \\ \infty & \text{otherwise} \end{cases}$$

where  $\beta(\mathbf{S}) \geq 0$  and  $\beta(\mathbf{S}) = 0$  when  $\mathbf{S} = \kappa \mathbf{I}_r$  for some  $\kappa > 0$ . Also  $\alpha$  is defined as

$$\alpha(a, b) = \frac{b}{a - b - 1} \quad (34)$$

*Proof.* We will prove the result for the cases  $r < n_1$  and  $r \geq n_1$ . Firstly, using Proposition 5.3, we already get that

$$\begin{aligned}
 \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) &= \mathcal{L}_{\lambda, \text{rep}}^{\text{test}}(\mathbf{A}; n_1) = \sigma^2 + \mathbb{E}_{\rho_v \sim \bar{\mu}} \left[ \mathbb{E}_{(\mathbf{X}, \mathbf{Y}) \sim \rho_v^{n_1}} [\|\mathbf{A} \mathbf{w}_\lambda(\mathbf{A}; (\mathbf{X}, \mathbf{Y})) - \mathbf{v}\|^2] \right] \\
 &= \sigma^2 + \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)} \left[ \mathbb{E}_{\substack{\mathbf{X} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)^{n_1}, \\ \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}_d, \sigma^2 \mathbf{I}_n)}} [\|\mathbf{A} \mathbf{w}_\lambda(\mathbf{A}; (\mathbf{X}, \mathbf{X} \mathbf{v} + \boldsymbol{\eta})) - \mathbf{v}\|^2] \right] \quad (35)
 \end{aligned}$$

We will use this expression in the rest of the proof.

**Case 1:  $r \leq n_1$**  In this case, the rank of the representations  $\mathbf{X} \mathbf{A} \in \mathbb{R}^{n_1 \times d}$  for training data is higher than the number of samples. Thus the unique minimizer for the dataset  $(\mathbf{X}, \mathbf{Y})$  is

$$\lim_{\lambda \rightarrow 0} \mathbf{w}_\lambda(\mathbf{A}; (\mathbf{X}, \mathbf{Y})) = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \ell_{\lambda, \text{rep}}(\mathbf{w}; \mathbf{A}, (\mathbf{X}, \mathbf{Y})) = (\mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{X}^\top \mathbf{Y}$$



Plugging this into Equation (35), we get

$$\begin{aligned} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 &= \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)} \left[ \mathbb{E}_{\substack{\mathbf{X} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)^{n_1}, \\ \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}_d, \sigma^2 \mathbf{I}_{n_1})}} \left[ \|\mathbf{A} (\mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{v} + \boldsymbol{\eta}) - \mathbf{v}\|^2 \right] \right] \\ &\stackrel{(a)}{=} \underbrace{\mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{A} (\mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right]}_{\text{bias}(\mathbf{A})} + \underbrace{\mathbb{E}_{\boldsymbol{\eta}, \mathbf{X}} \left[ \|\mathbf{A} (\mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{X}^\top \boldsymbol{\eta}\|^2 \right]}_{\text{variance}(\mathbf{A})} \end{aligned}$$

where (a) follows from the independence of  $\boldsymbol{\eta}$  and  $(\mathbf{X}, \mathbf{v})$  and that  $\mathbb{E}[\boldsymbol{\eta}] = \mathbf{0}_n$ . We will analyze the bias and variance terms separately below

**Bias:** Recall that  $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$  is the singular value decomposition, with  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$  and  $\mathbf{S} \in \mathbb{R}^{r \times r}$ . We set  $\mathbf{X}_U = \mathbf{X} \mathbf{U} \in \mathbb{R}^{n_1 \times r}$ . The two key ideas that we will exploit are that  $\mathbf{X}_U \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_r)^{n_1}$  and that  $\mathbf{X}_U$  is independent of  $\mathbf{X} P_U^\perp$  since  $\mathbf{X}$  is Gaussian-distributed.

$$\begin{aligned} \text{bias}(\mathbf{A}) &= \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{A} (\mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right] \\ &\stackrel{(a)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} \mathbf{S} \mathbf{V}^\top (\mathbf{V} \mathbf{S} \mathbf{U}^\top \mathbf{X}^\top \mathbf{X} \mathbf{U} \mathbf{S} \mathbf{V}^\top)^\dagger \mathbf{V} \mathbf{S} \mathbf{U}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right] \\ &\stackrel{(b)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} \mathbf{S} (\mathbf{S} (\mathbf{X}_U)^\top \mathbf{X}_U \mathbf{S})^{-1} \mathbf{S} (\mathbf{X}_U)^\top \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right] \\ &\stackrel{(c)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} \mathbf{S} (\mathbf{S} \mathbf{X}_U^\top \mathbf{X}_U \mathbf{S})^{-1} \mathbf{S} \mathbf{X}_U^\top \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right] \stackrel{(d)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} \mathbf{S} \mathbf{S}^{-1} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{S}^{-1} \mathbf{S} \mathbf{X}_U^\top \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right] \\ &\stackrel{(e)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{X}_U^\top \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right] = \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{X}_U^\top \mathbf{X} (\mathbf{U} \mathbf{U}^\top + P_U^\perp) \mathbf{v} - (\mathbf{U} \mathbf{U}^\top + P_U^\perp) \mathbf{v}\|^2 \right] \\ &\stackrel{(f)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{X}_U^\top \mathbf{X} \mathbf{U} \mathbf{U}^\top \mathbf{v} + \mathbf{U} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{X}_U^\top \mathbf{X} P_U^\perp \mathbf{v} - \mathbf{U} \mathbf{U}^\top \mathbf{v} - P_U^\perp \mathbf{v}\|^2 \right] \\ &\stackrel{(g)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U} \mathbf{U}^\top \mathbf{v} - \mathbf{U} \mathbf{U}^\top \mathbf{v} + \mathbf{U} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{X}_U^\top \mathbf{X} P_U^\perp \mathbf{v}\|^2 + \|P_U^\perp \mathbf{v}\|^2 \right] \\ &\stackrel{(h)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \text{tr} \left( (P_U^\perp \mathbf{v})^\top P_U^\perp \mathbf{X}^\top \mathbf{X}_U (\mathbf{X}_U^\top \mathbf{X}_U)^{-2} \mathbf{X}_U^\top \mathbf{X} P_U^\perp (P_U^\perp \mathbf{v}) \right) \right] + \|P_U^\perp \mathbf{A}^*\|^2 \\ &\stackrel{(i)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \text{tr} \left( \mathbf{X}_U (\mathbf{X}_U^\top \mathbf{X}_U)^{-2} \mathbf{X}_U^\top \mathbb{E}_{P_U^\perp \mathbf{X}} \left[ \mathbf{X} P_U^\perp (P_U^\perp \mathbf{v}) (P_U^\perp \mathbf{v})^\top P_U^\perp \mathbf{X}^\top \right] \right) \right] + \|P_U^\perp \mathbf{A}^*\|^2 \end{aligned}$$

while (a) just uses the SVD of  $\mathbf{A}$ , (b) uses the simple fact that  $\mathbf{V}^\top (\mathbf{V} \mathbf{B} \mathbf{V}^\top)^\dagger \mathbf{V} = \mathbf{B}^\dagger$  for an orthogonal matrix  $\mathbf{V}$ . (d) follows from the fact that  $\mathbf{X}_U^\top \mathbf{X}_U \in \mathbb{R}^{r \times r}$  is full rank with probability 1, and thus invertible. (f) simply decomposes  $\mathbf{v} = \mathbf{U} \mathbf{U}^\top \mathbf{v} + P_U^\perp \mathbf{v}$ , while (g) follows from the orthogonality of  $P_U^\perp \mathbf{v}$  and any vector in the span of  $\mathbf{U}$ . (h) uses the simple facts that  $\|\mathbf{a}\|^2 = \text{tr}(\mathbf{a} \mathbf{a}^\top)$  and  $P_U^\perp P_U^\perp = P_U^\perp$  and (i) uses the crucial observation that  $\mathbf{X}_U$  is independent of  $\mathbf{X} P_U^\perp$ , since for Gaussian distribution, all subspaces are independent of its orthogonal subspace, and that  $\text{tr}$  is a linear operator.

We now look closer at the term  $\mathbf{M} = \mathbb{E} \left[ \mathbf{X} P_U^\perp (P_U^\perp \mathbf{v}) (P_U^\perp \mathbf{v})^\top P_U^\perp \mathbf{X}^\top \right]$  which is a matrix in  $\mathbb{R}^{n_1 \times n_1}$ .

$$\begin{aligned} \mathbf{M}_{i,j} &= \mathbb{E}_{P_U^\perp \mathbf{X}} \left[ \mathbf{x}_i^\top (P_U^\perp \mathbf{v}) (P_U^\perp \mathbf{v})^\top \mathbf{x}_j \right] \\ &= \text{tr} \left( (P_U^\perp \mathbf{v}) (P_U^\perp \mathbf{v})^\top \mathbb{E} [\mathbf{x}_j \mathbf{x}_i^\top] \right) = \begin{cases} 0 & \text{if } i \neq j \\ \|P_U^\perp \mathbf{v}\|^2 & \text{if } i = j \end{cases} \end{aligned} \quad (36)$$

This gives us that  $\mathbf{M} = \|P_U^\perp \mathbf{v}\|^2 \mathbf{I}_{n_1}$ . We can now complete the computation for  $\text{bias}(\mathbf{A})$ .

$$\begin{aligned} \text{bias}(\mathbf{A}) &= \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \text{tr} \left( \mathbf{X}_U (\mathbf{X}_U^\top \mathbf{X}_U)^{-2} \mathbf{X}_U^\top \|P_U^\perp \mathbf{v}\|^2 \right) \right] + \|P_U^\perp \mathbf{A}^*\|^2 \\ &= \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( \mathbf{X}_U (\mathbf{X}_U^\top \mathbf{X}_U)^{-2} \mathbf{X}_U^\top \right) \right] \mathbb{E}_{\mathbf{v}} \left[ \|P_U^\perp \mathbf{v}\|^2 \right] + \|P_U^\perp \mathbf{A}^*\|^2 \end{aligned}$$

$$= \mathbb{E}_{\mathbf{X}_U} [\text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1})] \|P_U^\perp \mathbf{A}^*\|^2 + \|P_U^\perp \mathbf{A}^*\|^2$$

We will deal with the  $\mathbb{E}_{\mathbf{X}_U} [\text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1})]$  term later and show that it is equal to  $\alpha(n_1, r)$ .

**Variance:** We will use many ideas that were used for the bias term. Again using the SVD, we get

$$\begin{aligned} \text{variance}(\mathbf{A}) &= \mathbb{E}_{\eta, \mathbf{X}} \left[ \|\mathbf{A} (\mathbf{A}^\top \mathbf{X}^\top \mathbf{X} \mathbf{A})^\dagger \mathbf{A}^\top \mathbf{X}^\top \eta\|^2 \right] \\ &=^{(a)} \mathbb{E}_{\eta, \mathbf{X}_U} \left[ \|\mathbf{U} \mathbf{S} \mathbf{V}^\top (\mathbf{V} \mathbf{S} \mathbf{X}_U^\top \mathbf{X}_U \mathbf{S} \mathbf{V}^\top)^\dagger \mathbf{V} \mathbf{S} \mathbf{X}_U^\top \eta\|^2 \right] \\ &=^{(b)} \mathbb{E}_{\eta, \mathbf{X}_U} \left[ \|\mathbf{U} \mathbf{S} \mathbf{S}^{-1} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{S}^{-1} \mathbf{S} \mathbf{X}_U^\top \eta\|^2 \right] =^{(c)} \mathbb{E}_{\eta, \mathbf{X}_U} \left[ \|\mathbf{X}_U^\top \mathbf{X}_U\|^{-1} \|\mathbf{X}_U^\top \eta\|^2 \right] \\ &=^{(d)} \mathbb{E}_{\eta, \mathbf{X}_U} \left[ \text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{X}_U^\top \eta \eta^\top \mathbf{X}_U (\mathbf{X}_U^\top \mathbf{X}_U)^{-1}) \right] \\ &=^{(e)} \sigma^2 \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1} \mathbf{X}_U^\top \mathbf{X}_U (\mathbf{X}_U^\top \mathbf{X}_U)^{-1}) \right] \\ &=^{(f)} \sigma^2 \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1}) \right] \end{aligned}$$

Here (a) uses SVD, (b) uses the fact that as before  $\mathbf{X}_U^\top \mathbf{X}_U$  is full rank with probability 1, (d) follows from the norm and trace relationship, (e) follows from the noise distribution  $\eta \sim \mathcal{N}(\mathbf{0}_{n_1}, \sigma^2 \mathbf{I}_{n_1})$  and its independence from  $\mathbf{X}$ .

Thus combining the bias and variance terms, we get

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, \mathbf{a})) - \sigma^2 = \text{bias}(\mathbf{A}) + \text{variance}(\mathbf{A}) = \left( 1 + \mathbb{E}_{\mathbf{X}_U} [\text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1})] \right) \|P_U^\perp\|^2 + \mathbb{E}_{\mathbf{X}_U} [\text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1})] \sigma^2$$

Thus the only thing remaining to show is that  $\mathbb{E}_{\mathbf{X}_U} [\text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1})] = \alpha(n_1, r)$ . To show this, we will use the fact that  $(\mathbf{X}_U^\top \mathbf{X}_U)^{-1}$  is from the inverse Wishart distribution, and that  $\mathbb{E}_{\mathbf{X}_U} (\mathbf{X}_U^\top \mathbf{X}_U)^{-1} = \frac{\mathbf{I}_{n_1}}{r - n_1 - 1}$  when  $r > n_1 + 1$  and unbounded when  $r \in \{n_1, n_1 + 1\}$  (Mardia et al., 1979; Belkin et al., 2020). For  $r < n_1 - 1$ , this gives us  $\mathbb{E}_{\mathbf{X}_U} [\text{tr}((\mathbf{X}_U^\top \mathbf{X}_U)^{-1})] = \text{tr} \left( \mathbb{E}_{\mathbf{X}_U} [(\mathbf{X}_U^\top \mathbf{X}_U)^{-1}] \right) = \text{tr} \left( \frac{\mathbf{I}_{n_1}}{r - n_1 - 1} \right) = \frac{n_1}{r - n_1 - 1}$ , which completes the proof. We now prove the result for  $r > n_1$ .

**Case 2:  $r > n_1$**  In this case, the rank of the representations  $\mathbf{X} \mathbf{A} \in \mathbb{R}^{n_1 \times d}$  for training data is lower than the number of samples. Thus we can use the dual formulation to get the minimum  $\ell_2$  norm solution for dataset  $(\mathbf{X}, \mathbf{Y})$

$$\lim_{\lambda \rightarrow 0} \mathbf{w}_\lambda(\mathbf{A}; (\mathbf{X}, \mathbf{Y})) = \arg \min_{\mathbf{X} \mathbf{A} \mathbf{w} = \mathbf{Y}} \|\mathbf{w}\|_2 = \mathbf{A}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top)^{-1} \mathbf{Y}$$

Plugging this into Equation (35), we get

$$\begin{aligned} \mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, n_2)) - \sigma^2 &= \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)} \left[ \mathbb{E}_{\substack{\mathbf{X} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)^{n_1} \\ \eta \sim \mathcal{N}(\mathbf{0}_d, \sigma^2 \mathbf{I}_{n_1})}} \left[ \|\mathbf{A} \mathbf{A}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top)^{-1} (\mathbf{X} \mathbf{v} + \eta) - \mathbf{v}\|^2 \right] \right] \\ &=^{(a)} \underbrace{\mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{A} \mathbf{A}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{v} - \mathbf{v}\|^2 \right]}_{\text{bias}(\mathbf{A})} + \underbrace{\mathbb{E}_{\eta, \mathbf{X}} \left[ \|\mathbf{A} \mathbf{A}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top)^{-1} \eta\|^2 \right]}_{\text{variance}(\mathbf{A})} \end{aligned}$$

We again handle the bias and variance terms separately. Let  $\mathbf{X}_{U^\perp} = \mathbf{X} P_U^\perp$  and we will again use the fact that  $\mathbf{X}_U = \mathbf{X} \mathbf{U} \mathbf{U}^\top$  and  $\mathbf{X}_{U^\perp}$  are independent

**Bias:**

$$\begin{aligned}
 \text{bias}(\mathbf{A}) &= \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{A}\mathbf{A}^\top \mathbf{X}^\top (\mathbf{X}\mathbf{A}\mathbf{A}^\top \mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{v} - \mathbf{v}\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U}\mathbf{S}^2 \mathbf{U}^\top \mathbf{X}^\top (\mathbf{X}\mathbf{U}\mathbf{S}^2 \mathbf{U}^\top \mathbf{X}^\top)^{-1} \mathbf{X}(\mathbf{U}\mathbf{U}^\top + \mathbf{P}_U^\perp) \mathbf{v} - (\mathbf{U}\mathbf{U}^\top + \mathbf{P}_U^\perp) \mathbf{v}\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U}(\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbf{U}^\top \mathbf{v} + \mathbf{U} \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X} \mathbf{P}_U^\perp \mathbf{v} - \mathbf{P}_U^\perp \mathbf{v}\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{U}(\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbf{U}^\top \mathbf{v} + \mathbf{U} \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v}\|^2 + \|\mathbf{P}_U^\perp \mathbf{v}\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v} + \mathbf{C}_{\mathbf{X}_U} \mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v}\|^2 \right] + \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2
 \end{aligned}$$

where  $\mathbf{B}_{\mathbf{X}_U} = \mathbf{U}(\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r)$  and  $\mathbf{C}_{\mathbf{X}_U} = \mathbf{U} \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1}$  only depends on the components of  $\mathbf{X}$  in the direction of  $\mathbf{U}$ , i.e.  $\mathbf{X}_U$ . The main difference from the  $r < n_1$  case is that here  $\mathbf{B}_{\mathbf{X}_U} \in \mathbb{R}^{r \times r}$  is not zero, since the rank of  $\mathbf{X}_U$  is  $\min\{n_1, r\} = n_1$ . We can expand the bias term further

$$\begin{aligned}
 \text{bias}(\mathbf{A}) &= \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2 + \|\mathbf{C}_{\mathbf{X}_U} \mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v}\|^2 + 2\text{tr}(\mathbf{v}^\top \mathbf{U} \mathbf{B}_{\mathbf{X}_U}^\top \mathbf{C}_{\mathbf{X}_U} \mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v}) \right] + \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2 \\
 &\stackrel{(a)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2 \right] + \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{C}_{\mathbf{X}_U} \mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v}\|^2 \right] + 2 \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \text{tr} \left( \mathbf{v}^\top \mathbf{U} \mathbf{B}_{\mathbf{X}_U}^\top \mathbf{C}_{\mathbf{X}_U} \mathbb{E}_{\mathbf{X}_{U^\perp}} [\mathbf{X}_{U^\perp}] \mathbf{P}_U^\perp \mathbf{v} \right) \right] + \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2 \\
 &\stackrel{(b)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2 \right] + \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \|\mathbf{C}_{\mathbf{X}_U} \mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v}\|^2 \right] + \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2 \\
 &\stackrel{(c)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2 \right] + \mathbb{E}_{\mathbf{v}, \mathbf{X}} \left[ \text{tr}(\mathbf{v}^\top \mathbf{P}_U^\perp \mathbf{X}_{U^\perp}^\top \mathbf{C}_{\mathbf{X}_U}^\top \mathbf{C}_{\mathbf{X}_U} \mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v}) \right] + \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2 \\
 &\stackrel{(d)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2 \right] + \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \text{tr} \left( \mathbf{C}_{\mathbf{X}_U}^\top \mathbf{C}_{\mathbf{X}_U} \mathbb{E}_{\mathbf{X}_{U^\perp}} [\mathbf{X}_{U^\perp} \mathbf{P}_U^\perp \mathbf{v} \mathbf{v}^\top \mathbf{P}_U^\perp \mathbf{X}_{U^\perp}^\top] \right) \right] + \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2 \\
 &\stackrel{(e)}{=} \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2 \right] + \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr}(\mathbf{C}_{\mathbf{X}_U}^\top \mathbf{C}_{\mathbf{X}_U}) \right] \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2 + \|\mathbf{P}_U^\perp \mathbf{A}^*\|^2 \tag{37}
 \end{aligned}$$

where (a) uses the fact that  $\mathbf{X}_{U^\perp}$  is independent of  $\mathbf{X}_U$  and  $\mathbf{v}$ , (b) uses the Gaussianity of  $\mathbf{X}_{U^\perp}$  and that it has mean  $\mathbf{0}$ , (c) uses  $\|\mathbf{a}\|^2 = \text{tr}(\mathbf{a}\mathbf{a}^\top)$ , (d) uses the independence of  $\mathbf{X}_U$  and  $\mathbf{X}_{U^\perp}$  and (e) uses the calculation from Equation (36). We first tackle the  $\mathbf{C}_{\mathbf{X}_U}$  term

$$\begin{aligned}
 \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr}(\mathbf{C}_{\mathbf{X}_U}^\top \mathbf{C}_{\mathbf{X}_U}) \right] &= \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^2 \mathbf{U}^\top \mathbf{U} \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \\
 &= \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \tag{38}
 \end{aligned}$$

We will encounter this function again in the variance term and we will tackle this later. At a high level, we will show that this term is going to be at least as large as the term for  $\mathbf{S} = \mathbf{I}_r$ , which reduces to  $\mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \right) \right]$  which has a closed form expression, again using inverse Wishart distribution. For now we will first deal with the  $\mathbf{B}_{\mathbf{X}_U}$  term.

$$\begin{aligned}
 \mathbb{E}_{\mathbf{v}, \mathbf{X}_U} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2 \right] &= \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{U}^\top \mathbf{A}^* \mathbf{A}^* \mathbf{U}), \mathbf{X}_U} \left[ \|\mathbf{B}_{\mathbf{X}_U} \mathbf{u}\|^2 \right] = \mathbb{E}_{\mathbf{u}, \mathbf{X}_U} \left[ \|\mathbf{U}(\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbf{u}\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U} \left[ \|\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r\| \mathbf{u}\|^2 \right]
 \end{aligned}$$

We are now going to exploit the symmetry in Gaussian distribution once again. Recall that  $\mathbf{X}_U \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)$ . We notice that  $\mathbf{X}_U \equiv_D \mathbf{X}_U \mathbf{P} \mathbf{D}$ , where  $\mathbf{P}$  is a random permutation matrix and  $\mathbf{D}$  is a diagonal matrix with random entries in  $\pm 1$ . Essentially this is saying that randomly shuffling the coordinates and multiplying each coordinate by a random

sign results in the same isotropic Gaussian distribution. We observe that  $PDS^2DP^\top$  for diagonal matrix  $S$  and that  $PP^\top = P^\top P = D^2 = I_r$  rewrite the above expectation.

$$\begin{aligned}
 \mathbb{E}_{\mathbf{u}, \mathbf{X}_U} [\|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2] &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U} \left[ \left\| (\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbf{u} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U, \mathbf{P}, \mathbf{D}} \left[ \left\| (\mathbf{S}^2 \mathbf{D} \mathbf{P}^\top \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{P} \mathbf{D} \mathbf{S}^2 \mathbf{D} \mathbf{P}^\top \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{P} \mathbf{D} - \mathbf{I}_r) \mathbf{u} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U, \mathbf{P}, \mathbf{D}} \left[ \left\| (\mathbf{S}^2 \mathbf{D} \mathbf{P}^\top \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{P} \mathbf{D} - \mathbf{D} \mathbf{P}^\top \mathbf{P} \mathbf{D}) \mathbf{u} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U, \mathbf{P}, \mathbf{D}} \left[ \left\| \mathbf{D} \mathbf{P}^\top (\mathbf{P} \mathbf{D} \mathbf{S}^2 \mathbf{D} \mathbf{P}^\top \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbf{P} \mathbf{D} \mathbf{u} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U, \mathbf{P}, \mathbf{D}} \left[ \left\| \mathbf{D} \mathbf{P}^\top (\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbf{P} \mathbf{D} \mathbf{u} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U, \mathbf{P}, \mathbf{D}} \left[ \left\| (\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbf{P} \mathbf{D} \mathbf{u} \right\|^2 \right] \\
 &= \mathbb{E}_{\mathbf{u}, \mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r) \mathbb{E}_{\mathbf{P}, \mathbf{D}} [\mathbf{P} \mathbf{D} \mathbf{u} \mathbf{u}^\top \mathbf{D} \mathbf{P}^\top] (\mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r)^\top \right) \right]
 \end{aligned}$$

It is not hard to see that randomly multiplying coordinates of  $\mathbf{u}$  by  $\pm 1$  and then shuffling the coordinates will lead to  $\mathbb{E}_{\mathbf{P}, \mathbf{D}} [\mathbf{P} \mathbf{D} \mathbf{u} \mathbf{u}^\top \mathbf{D} \mathbf{P}^\top] = \frac{\|\mathbf{u}\|^2}{r} \mathbf{I}_r$ .

$$\begin{aligned}
 \mathbb{E}_{\mathbf{u}, \mathbf{X}_U} [\|\mathbf{B}_{\mathbf{X}_U} \mathbf{U}^\top \mathbf{v}\|^2] &= \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}_r, \mathbf{U}^\top \mathbf{A}^* \mathbf{A}^* \mathbf{U})} \left[ \frac{\|\mathbf{u}\|^2}{r} \right] \mathbb{E}_{\mathbf{X}_U} \left[ \left\| \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - \mathbf{I}_r \right\|^2 \right] \\
 &= \frac{\|P_U \mathbf{A}^*\|^2}{r} \mathbb{E}_{\mathbf{X}_U} \left[ \left\| \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U P_{\mathbf{X}_U} - P_{\mathbf{X}_U} - P_{\mathbf{X}_U}^\perp \right\|^2 \right] \\
 &= \frac{\|P_U \mathbf{A}^*\|^2}{r} \mathbb{E}_{\mathbf{X}_U} \left[ \left\| \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - P_{\mathbf{X}_U} \right\|^2 + \|P_{\mathbf{X}_U}^\perp\|^2 \right] \\
 &= \beta_1(\mathbf{S}) + \frac{\|P_U \mathbf{A}^*\|^2}{r} \mathbb{E}_{\mathbf{X}_U} [\text{rank}(P_{\mathbf{X}_U}^\perp)] \\
 &= \beta_1(\mathbf{S}) + \frac{r - n_1}{r} \|P_U \mathbf{A}^*\|^2 \tag{39}
 \end{aligned}$$

where  $\beta_1(\mathbf{S}) = \frac{\|P_U \mathbf{A}^*\|^2}{r} \mathbb{E}_{\mathbf{X}_U} \left[ \left\| \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U - P_{\mathbf{X}_U} \right\|^2 \right]$  is a function that satisfies  $\beta_1(\mathbf{S}) \geq 0$  and  $\beta_1(\kappa \mathbf{I}_r) = 0$  for any  $\kappa$ . Also we used that  $\text{rank}(P_{\mathbf{X}_U}^\perp) = r - \text{rank}(\mathbf{X}_U) = r - n_1$  with probability 1.

Plugging Equations (38) and (39) into Equation (37), we get the following final expression for the bias

$$\begin{aligned}
 \text{bias}(\mathbf{A}) &= \left( 1 + \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \right) \|P_U^\perp \mathbf{A}^*\|^2 + \frac{r - n_1}{r} \|P_U \mathbf{A}^*\|^2 + \beta_1(\mathbf{S}) \\
 &= \left( 1 + \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \right) \|P_U^\perp \mathbf{A}^*\|^2 + \frac{r - n_1}{r} (1 - \|P_U^\perp \mathbf{A}^*\|^2) + \beta_1(\mathbf{S}) \\
 &= \left( 1 + \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] - \frac{r - n_1}{r} \right) \|P_U^\perp \mathbf{A}^*\|^2 + \frac{r - n_1}{r} + \beta_1(\mathbf{S}) \\
 &= \left( \frac{n_1}{r} + \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \right) \|P_U^\perp \mathbf{A}^*\|^2 + \frac{r - n_1}{r} + \beta_1(\mathbf{S}) \tag{40}
 \end{aligned}$$

**Variance:** We now move to the variance term

$$\begin{aligned}
 \text{variance}(\mathbf{A}) &= \mathbb{E}_{\boldsymbol{\eta}, \mathbf{X}} \left[ \left\| \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top)^{-1} \boldsymbol{\eta} \right\|^2 \right] \\
 &= \sigma^2 \mathbb{E}_{\mathbf{X}} \left[ \text{tr} \left( \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{A} \mathbf{A}^\top \mathbf{X}^\top)^{-2} \mathbf{X} \mathbf{A} \mathbf{A}^\top \right) \right]
 \end{aligned}$$

$$\begin{aligned}
&= \sigma^2 \mathbb{E}_{\mathbf{X}} \left[ \text{tr} \left( \mathbf{U} \mathbf{S}^2 \mathbf{U}^\top \mathbf{X}^\top (\mathbf{X} \mathbf{U} \mathbf{S}^2 \mathbf{U}^\top \mathbf{X}^\top)^{-2} \mathbf{X} \mathbf{U} \mathbf{S}^2 \mathbf{U}^\top \right) \right] \\
&= \sigma^2 \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-2} \mathbf{X}_U \mathbf{S}^2 \right) \right] \\
&= \sigma^2 \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \tag{41}
\end{aligned}$$

To further simplify both, the bias and variance terms, we need the following result

**Lemma C.2.** For  $\mathbf{X}_U \sim \mathcal{N}(\mathbf{0}_r, \mathbf{I}_r)^{n_1}$  and diagonal matrix  $\mathbf{S} \in \mathbb{R}^{r \times r}$

$$\mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] = \begin{cases} \alpha(r, n_1) + \beta_2(\mathbf{S}) & \text{if } r > n_1 + 1 \\ \infty & \text{if } r \in \{n_1, n_1 + 1\} \end{cases}$$

where  $\alpha(r, n_1)$  is defined in Equation (34) and  $\beta_2(\mathbf{S}) \geq 0$  and  $\beta_2(\kappa \mathbf{I}_r) = 0$  for any  $\kappa > 0$ .

Using Lemma C.2 and plugging it into Equations (40) and (41), we get

$$\mathcal{L}_{\lambda, \text{rep}}^{\text{tr-val}}(\mathbf{A}; (n_1, \mathbf{a})) - \sigma^2 = \text{bias}(\mathbf{A}) + \text{variance}(\mathbf{A}) = \left( \frac{n_1}{r} + \alpha(r, n_1) \right) \|P_U^\perp \mathbf{A}^*\|^2 + \frac{r - n_1}{r} + \sigma^2 \alpha(r, n_1) + \beta(\mathbf{S})$$

where  $\beta(\mathbf{S}) = \beta_1(\mathbf{S}) + \beta_2(\mathbf{S}) (\|P_U^\perp \mathbf{A}^*\|^2 + \sigma^2)$  is a non-negative function that is 0 at  $\kappa \mathbf{I}_r$  for all  $\kappa > 0$ .

We now complete the proof by proving Lemma C.2

*Proof of Lemma C.2.* Let the L.H.S. be  $\gamma(\mathbf{S}) = \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right]$ . We first show that equality holds for  $\mathbf{S} = \kappa \mathbf{I}_d$ . In this case, we have  $\gamma(\kappa \mathbf{I}_r) = \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \right) \right] = \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \right) \right]$ . Using the closed-form expression for inverse Wishart distribution once again, we conclude

$$\gamma(\kappa \mathbf{I}_r) = \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \right) \right] = \text{tr} \left( \mathbb{E}_{\mathbf{X}_U} \left[ (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \right] \right) = \begin{cases} \text{tr} \left( \frac{\mathbf{I}_{n_1}}{r - n_1 - 1} \right) = \frac{n_1}{r - n_1 - 1} & \text{if } r > n_1 + 1 \\ \infty & \text{otherwise} \end{cases}$$

Thus we get  $\gamma(\kappa \mathbf{I}_r) = \alpha(r, n_1)$  when  $r > n_1 + 1$  and unbounded otherwise.

We will show for an arbitrary diagonal matrix  $\mathbf{S}$  that the value is at least as large as  $\mathbf{I}_r$ , i.e.  $\gamma(\mathbf{S}) \geq \gamma(\kappa \mathbf{I}_r)$ , which will complete the proof. We first observe that  $\mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top = \mathbf{X}_U \mathbf{S}^2 \mathbf{S}^2 \mathbf{X}_U^\top \succeq \mathbf{X}_U \mathbf{S}^2 P_{\mathbf{X}_U} \mathbf{S}^2 \mathbf{X}_U^\top = \mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top$ . Using this, we get

$$\begin{aligned}
\gamma(\mathbf{S}) &= \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^4 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \\
&\geq \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top (\mathbf{X}_U \mathbf{S}^2 \mathbf{X}_U^\top)^{-1} \right) \right] \\
&\geq \mathbb{E}_{\mathbf{X}_U} \left[ \text{tr} \left( (\mathbf{X}_U \mathbf{X}_U^\top)^{-1} \right) \right] \\
&= \gamma(\mathbf{I}_r)
\end{aligned}$$

□

□

## D. Experiments

This section contains experimental details and also additional experiments on more datasets and settings. The code for all experiments will be made public.

**Algorithm 1** RepLearn( $\theta^{\text{rep}}$ ,  $\text{TaskLoader}$ ,  $\text{variant}$ )

---

**Parameters:** Regularization ( $\lambda$ ), Outer steps ( $T_{\text{out}}$ ), outer lr ( $\eta_{\text{out}}$ ), batch size ( $b$ ), Inner steps ( $T_{\text{in}}$ ), inner lr ( $\eta_{\text{in}}$ )

**Input:** Representation model  $\theta^{\text{rep}}$ ,  $\text{TaskLoader}$ ,  $\text{variant}$  {tr-tr or tr-val}

$\theta^{\text{rep}}(0) \leftarrow \text{RandInit}$

**for**  $t = 0$  **to**  $T_{\text{out}} - 1$  **do**

$\text{TaskBatch} \leftarrow \text{TaskLoader}(\text{batch\_size} = b)$

**for**  $i = 1$  **to**  $b$  **do**

$S \leftarrow \text{TaskBatch}[i]$  {Dataset of size  $n$ }

**if**  $\text{variant}$  is ‘tr-val’ **then**

$(S^{\text{tr}}, S^{\text{val}}) \leftarrow_{\text{split}} S$  {Split into tr-val sets of sizes  $n_1 + n_2 = n$ }

$\widetilde{\mathbf{W}} \leftarrow_{\text{stop gradient}} \text{InnerLoop}(\theta^{\text{rep}}(t), S^{\text{tr}}, \lambda)$  {Ignore the dependence of  $\widetilde{\mathbf{W}}$  on  $\theta^{\text{rep}}(t)$ }

$\nabla_i \leftarrow \nabla_{\theta^{\text{rep}}} \ell_{\text{rep}}(\widetilde{\mathbf{W}}; \theta^{\text{rep}}, S^{\text{val}})|_{\theta^{\text{rep}}(t)}$

**else if**  $\text{variant}$  is ‘tr-tr’ **then**

$\widetilde{\mathbf{W}} \leftarrow_{\text{stop gradient}} \text{InnerLoop}(\theta^{\text{rep}}(t), S, \lambda)$

$\nabla_i \leftarrow \nabla_{\theta^{\text{rep}}} \ell_{\text{rep}}(\widetilde{\mathbf{W}}; \theta^{\text{rep}}, S)|_{\theta^{\text{rep}}(t)}$  {Defined in Equation (42)}

**end if**

$\nabla = \frac{1}{b} \sum_{i=1}^b [\nabla_i]$

$\theta^{\text{rep}}(t+1) \leftarrow \text{Adam}(\theta^{\text{rep}}(t), \nabla, \eta_{\text{in}})$

**end for**

**end for**

**return**  $\theta^{\text{rep}}(T_{\text{out}})$

---

**Algorithm 2** InnerLoop( $\theta^{\text{rep}}$ ,  $S$ )

---

**Parameters:** Regularization ( $\lambda$ ), Inner steps ( $T_{\text{in}}$ ), inner lr ( $\eta_{\text{in}}$ )

**Input:** Representation model  $\theta^{\text{rep}}$ ,  $S$

$\mathbf{W}(0) \leftarrow \mathbf{0}_{d \times k}$

**for**  $t = 0$  **to**  $T_{\text{in}} - 1$  **do**

$\nabla \leftarrow \nabla_{\mathbf{W}} \{ \ell_{\text{rep}}(\mathbf{W}; \theta^{\text{rep}}, S) + \frac{\lambda}{2} \|\mathbf{W}\|_F^2 \} |_{\mathbf{W}(t)}$

$\mathbf{W}(t+1) \leftarrow \text{SGD}(\mathbf{W}(t), \nabla, \eta_{\text{in}}, \text{momentum} = 0.9)$

**end for**

**return**  $\mathbf{W}(T_{\text{in}})$

---

### D.1. RepLearn Algorithm and Datasets

**RepLearn** We describe the inner and outer loops for the RepLearn algorithm that we use for experiments in Algorithms 2 and 1 respectively. Recall the definitions for inner and outer losses.

$$\ell_{\text{rep}}(\mathbf{W}; \theta^{\text{rep}}, S) = \mathbb{E}_{(\mathbf{x}, y) \sim S} [\ell(\mathbf{W}^\top f_{\theta^{\text{rep}}}(\mathbf{x}), y)] \quad (42)$$

$$\ell_{\lambda, \text{rep}}(\mathbf{W}; \theta^{\text{rep}}, S) = \mathbb{E}_{(\mathbf{x}, y) \sim S} [\ell(\mathbf{W}^\top f_{\theta^{\text{rep}}}(\mathbf{x}), y)] + \frac{\lambda}{2} \|\mathbf{W}\|_2^2 \quad (43)$$

$$\mathcal{A}_{\lambda, \text{rep}}(\theta^{\text{rep}}; S) = \arg \min_{\mathbf{W}} \ell_{\lambda, \text{rep}}(\mathbf{W}; \theta^{\text{rep}}, S) \quad (44)$$

tr-tr Inner loop: For dataset  $S$  and current initialization  $\theta_i^{\text{rep}}$ , run  $T_{\text{in}}$  of gradient descent (with momentum 0.9) with

learning rate  $\eta_{\text{in}}$  on  $\ell_{\lambda, \text{rep}}(\mathbf{W}; \theta^{\text{rep}}, S)$  to get an approximation  $\widetilde{\mathbf{W}}_{\lambda, \text{rep}}(\theta^{\text{rep}}; S) \approx \mathcal{A}_{\lambda, \text{rep}}(\theta^{\text{rep}}; S)$ . Compute the gradient:  $\nabla^{\text{tr-val}}(\theta^{\text{rep}}) = \nabla \ell_{\lambda, \text{rep}}(\widetilde{\mathbf{W}}_{\lambda, \text{rep}}(\theta_i^{\text{rep}}, S); \cdot, S) |_{\theta_i^{\text{rep}}}$

Outer loop: Run Adam with learning rate  $\eta_{\text{out}}$  (other parameters at default value) with batch size  $b$  for  $T_{\text{out}}$  steps by using the gradient

Meta-testing: Tune  $T_{\text{in}}$  and  $\bar{\lambda}$  using validation tasks

**Datasets** We conduct experiments on the Omniglot (Lake et al., 2015) and MiniImageNet (Vinyals et al., 2016) datasets. The Omniglot dataset consists of 1623 different handwritten characters from 50 different alphabets. Each character was hand drawn by 20 different people. The original Omniglot dataset was split into a background set comprised of 30 alphabets and an evaluation set of 20 alphabets. We use the split recommended by Vinyals et al. (2016), which contains of a training split of 1028 characters, a validation split of 172 characters, and test split of 423 characters. Vinyals et al. (2016) construct the MiniImageNet dataset by sampling 100 random classes from ImageNet. We use 64 classes for training, 16 for validation, and 20 for testing. We use torchmeta (Deleu et al., 2019) to load datasets. All our evaluations in meta-test time are conducted in the transductive setting.

## D.2. Omniglot Experiments

**Relearn on Omniglot** We use a batch size  $b = 32$ . We use the standard 4-layer convolutional backbone with each layer having 64 output filters followed by batch normalization and ReLU activations. We resize the images to be 28x28 and apply 90, 180, and 270 degree rotations to augment the data, as in prior work, during training and evaluation. We train for  $T_{\text{out}} = 30000$  meta-steps and use  $T_{\text{in}} = 100$  inner steps regardless of model, and use an inner learning rate  $\eta_{\text{in}} = 0.05$ , unless there is a failure of optimization, in which case we reduce the learning rate to 0.01. We use an outer learning rate  $\eta_{\text{out}} = 0.001$ . We evaluate on 600 tasks at meta-test time. At meta-test time, for each model, we pick the best  $\bar{\lambda}$  and inner step size based on the validation set, where we explore  $\bar{\lambda} \in [0, .1, .3, 1.0, 2.0, 3.0, 10.0, 100.0]$  and inner step size in  $[50, 100, 200]$  and evaluate on the test set.

**Increasing width of a FC network on Omniglot** We examine the performance gap between tr-val versus tr-tr as we increase the expressive power of the representation network. We use a baseline 4 hidden layer fully connected network (FCN) inspired by Finn et al. (2017), but with  $64\ell$  nodes at all 4 layers for different values of  $\ell \in \{1, 4, 8, 16, 32\}$ . We set inner learning rates to be .05, except for  $\ell \in \{16, 32\}$  where we found a smaller inner learning rate of .01 was needed for convergence.

**iMAML** We use the original author code<sup>10</sup> as a starting point, which creates a convolutional neural network with four convolutional layers, followed by Batch Normalization and ReLU activations. We modify their code to add a tr-tr variant by using the combined data for the inner loop and the outer loop updates. We apply 90, 180, and 270 degree rotations to Omniglot data, resizing each image to 28x28 pixels. We use 5 conjugate gradient steps. For meta-testing we pick the best  $\bar{\lambda} \in \{0, .1, .3, 1.0, 2.0, 3.0, 10.0, 100.0\}$  and  $n_{\text{in}} \in \{8, 16, 32, 64\}$  that maximizes accuracy on the validation set. All other hyper-parameters at meta-test time are equal to the values used during training. We use an outer learning rate of  $1e-3$ . We train for 30000 outer steps for all models tested, and set the number of inner steps  $n_{\text{steps}} = 16$  for 5-way 1-shot and 25 for 20-way 1-shot. We investigate the performance of tr-val versus tr-tr by examining different settings of the regularization parameter  $\lambda$ . We report our results for tr-val versus tr-tr for Omniglot 5-way 1-shot and 20-way 1-shot in Table 6. We find that tr-val significantly outperforms tr-tr in all settings, and the gap is much larger than for RepLearn.

**iMAML representations and t-SNE** We train a model on Omniglot 5-way 1-shot using iMAML with batch size 32, inner learning rate .05, and meta steps 30000, again using the original author code convolutional neural network. We use  $\lambda = 2.0$  for tr-val and  $\lambda = 1.0$  for tr-tr. We use the tr-val and tr-tr CNN models to get image representations for input images, and then perform t-SNE on 10 randomly selected classes. We report t-SNE and singular value decay results for iMAML in Figure 4. We also plot the t-SNE representations for varying values of the perplexity in Figure 2. As in the case of RepLearn, the tr-val representations are much better clustered in the tr-tr representations. Furthermore the tr-val representations have a sharper drop in singular values, suggesting that they have lower effective rank than tr-tr representations.

<sup>10</sup>[https://github.com/aravindr93/imaml\\_dev](https://github.com/aravindr93/imaml_dev)

Table 6. meta-test accuracies in % for a CNN trained with iMAML for tr-val versus tr-tr on Omniglot 5-way 1-shot and Omniglot 20-way 1-shot. We find a huge gap in performance between tr-val and tr-tr, even larger than that for RepLearn.

	5-way 1-shot	20-way 1-shot
$\lambda = 2.0$ tr-val	$97.90 \pm 0.58$	$91.0 \pm 0.54$
$\lambda = 10.0$ tr-tr	$49.22 \pm 1.83$	$14.45 \pm 0.61$
$\lambda = 2.0$ tr-tr	$43.71 \pm 1.92$	$16.18 \pm 0.65$
$\lambda = 1.0$ tr-tr	$47.18 \pm 1.90$	$16.96 \pm 0.66$
$\lambda = 0.3$ tr-tr	$48.61 \pm 1.90$	$16.50 \pm 0.64$
$\lambda = 0.1$ tr-tr	$48.60 \pm 1.93$	$17.70 \pm 0.69$
$\lambda = 0.0$ tr-tr	$49.21 \pm 1.92$	$18.30 \pm 0.67$

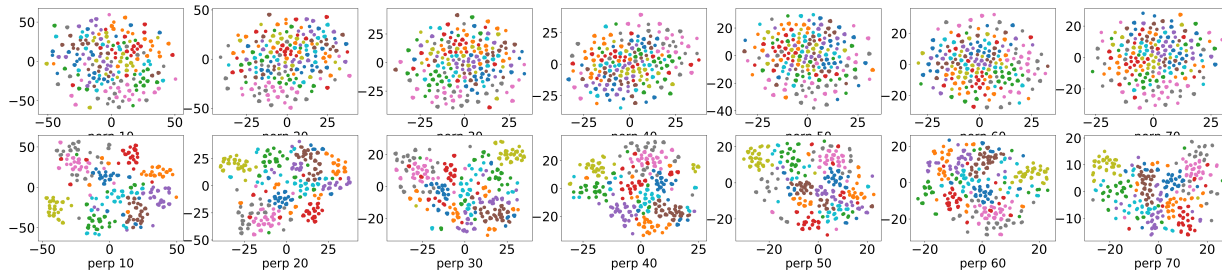


Figure 2. t-SNE plots for tr-val versus tr-tr for varying values of the perplexity for a CNN model trained with iMAML on Omniglot 5-way 1-shot. Plots shown for 10 randomly selected classes from the test split of Omniglot. Top row depicts tr-tr and bottom row depicts tr-val. We find that the tr-val representations appear to be more clustered for all values of perplexity of t-SNE.

**Rank and Expressivity** For a fully connected model of width factor  $\ell = 32$  trained with RepLearn on Omniglot 5-way 1-shot, we conduct linear regression twice. The first regression predicts the tr-tr representations given the tr-val representations, and the second predicts the tr-val representations given the tr-tr representations. We find that the  $R^2$  scores were 0.0973 and 0.0967, respectively. Thus the two sets of representations can express each other well enough, even though tr-val representations have lower effective rank.

We conduct the same experiment for a CNN model trained with iMAML on Omniglot 5-way 1-shot, and find that the  $R^2$  scores were 0.0103 and 0.171, respectively, thus suggesting that tr-val representations are a bit more expressive than tr-tr representations, in addition to having lower effective rank.

**Adding explicit regularization to RepLearn** For a CNN model trained with RepLearn on Omniglot 5-way 1-shot, we add explicit regularization to tr-tr by adding the Frobenius norm of the representation to the loss. We report the accuracy in percentage evaluated on 1200 tasks in Table 8 in the top row. We find that this significantly improves the performance of the tr-tr method, compared to the tr-tr models without explicit regularization in Table 10, or the bottom row of Table 8. This fits our intuition that the tr-tr method requires some form of regularization to learn low rank representations and to have guaranteed good performance on new tasks.

### D.3. MiniImageNet Experiments

**RepLearn on MiniImageNet** As standard (Rajeswaran et al., 2019), we resize the data to 84x84 pixel images and apply 90, 180, and 270 degree rotations and use a batch size of 16 during training. We use a convolutional neural network with four convolutional layers, with output filter sizes of 32, 64, 128, and 128 each followed by batch normalization and ReLU activations. We investigate the performance of tr-val versus tr-tr for RepLearn on the MiniImageNet 5-way 1-shot and 5-way 5-shot setting and report our results in Table 7. The findings are very similar to those from the Omniglot dataset, suggesting that our insights hold across multiple benchmark datasets.



## A Representation Learning Perspective on Train-Validation Splitting

Table 7. Tr-val versus tr-tr meta-test accuracies in % for a CNN model trained with RepLearn on MiniImageNet 5-way 1-shot and 5-way 5-shot for varying values of the regularization parameter,  $\lambda$ . The final value of the tr-tr objective is depicted in the last two columns for MiniImageNet 5-way 1-shot and MiniImageNet 5-way 5-shot, respectively. The tr-tr models make the tr-tr loss very small, which is what they were trained to minimize. Thus their failure on few-shot learning is not due to failure of optimization.

	5-way 1-shot	5-way 5-shot	tr-tr loss 5-way 1-shot	tr-tr loss 5-way 5-shot
$\lambda = 0.0$ tr-val	$46.16 \pm 1.67$	$65.36 \pm 0.91$	0.01	0.01
$\lambda = 0.0$ tr-tr	$25.53 \pm 1.43$	$33.49 \pm 0.82$	$1.1e-8$	$2.1e-6$
$\lambda = 0.1$ tr-tr	$24.69 \pm 1.32$	$34.91 \pm 0.85$	$3.5e-8$	$5.5e-7$
$\lambda = 1.0$ tr-tr	$25.88 \pm 1.45$	$40.19 \pm 1.12$	$1.9e-6$	$9.3e-5$

Table 8. Meta-test accuracies in % for tr-tr RepLearn trained on Omniglot 5-way 1-shot with explicit Frobenius norm regularization added to the representations.

	$\lambda = 0.0$ (tr-tr)	$\lambda = 0.1$ (tr-tr)	$\lambda = 3.0$ (tr-tr)
with representation regularization	$94.73 \pm 0.55$	$95.05 \pm 0.55$	$94.74 \pm 0.55$
no representation regularization	$67.78 \pm 1.60$	$67.53 \pm 1.66$	$89.00 \pm 1.08$

Table 9. Accuracies in % of representations parameterized by CNN networks of varying number of filters on MiniImageNet. Representations trained using tr-val objective consistently outperforms those learned using tr-tr objective, and the gap increases as width increases.

capacity = $num\_filters * \ell$	MiniImageNet 5-way 1-shot		Supervised 20-way	
	tr-val	tr-tr	tr-val	tr-tr
$\ell = 0.5$	$46.66 \pm 1.69$	$26.25 \pm 1.45$	1.	1.
$\ell = 1$	$48.44 \pm 1.62$	$26.81 \pm 1.44$	1.	1.
$\ell = 4$	$52.22 \pm 1.68$	$24.66 \pm 1.26$	1.	1.
$\ell = 8$	$52.25 \pm 1.71$	$25.28 \pm 1.37$	1.	1.

Table 10. Tr-val versus tr-tr meta-test accuracies in % for a CNN model trained with RepLearn on Omniglot 5-way 1-shot for varying values of the regularization parameter,  $\lambda$ .

	5-way 1-shot tr-val	5-way 1-shot tr-tr
$\lambda = 0.0$	$97.25 \pm 0.57$	$67.78 \pm 1.60$
$\lambda = 0.1$	$97.34 \pm 0.59$	$67.53 \pm 1.64$
$\lambda = 0.3$	$97.59 \pm 0.55$	$66.06 \pm 1.67$
$\lambda = 1.0$	$97.66 \pm 0.52$	$87.25 \pm 1.13$
$\lambda = 3.0$	$97.19 \pm 0.59$	$89.00 \pm 1.08$
$\lambda = 10.0$	$96.50 \pm 0.61$	$85.41 \pm 1.22$

**Increasing capacity of a CNN model on MiniImageNet** We start with a CNN model trained on MiniImageNet with four convolutional layers with 32, 64, 128, and 128 filters, respectively. Each convolution is followed by batch normalization and ReLU activations. We train with RepLearn. We increase the capacity by increasing the number of filters by a capacity factor,  $\ell$ , so that the convolutional layers contain  $32\ell$ ,  $64\ell$ ,  $128\ell$ , and  $128\ell$  output filters, respectively. We depict our results in Table 9. We find that increasing the network capacity improves the performance of tr-val representations, but slightly hurts tr-tr performance, just like the findings for Omniglot dataset with fully-connected networks. Thus the tr-val method is more robust to architecture choice/capacity and datasets.

**t-SNE on MiniImageNet** We take the baseline convolutional neural network with capacity factor  $\ell = 1$  from the previous section, and conduction t-SNE on the representations produced by the tr-val versus the tr-tr model. We report our results in Figure 3.

Table 11. Tr-val versus tr-tr meta-test accuracies in % for a CNN model trained with RepLearn on Omniglot 20-way 1-shot for varying values of the regularization parameter,  $\lambda$ .

	20-way 1-shot trval	20-way 1-shot tr-tr
$\lambda = 0.0$	$92.26 \pm 0.45$	$49.00 \pm 0.88$
$\lambda = 0.1$	$92.38 \pm 0.44$	$50.84 \pm 0.85$
$\lambda = 0.3$	$92.21 \pm 0.47$	$55.38 \pm 0.92$
$\lambda = 1.0$	$92.44 \pm 0.47$	$84.14 \pm 0.63$
$\lambda = 3.0$	$92.70 \pm 0.48$	$88.20 \pm 0.55$
$\lambda = 10.0$	$91.50 \pm 0.48$	$85.85 \pm 0.58$

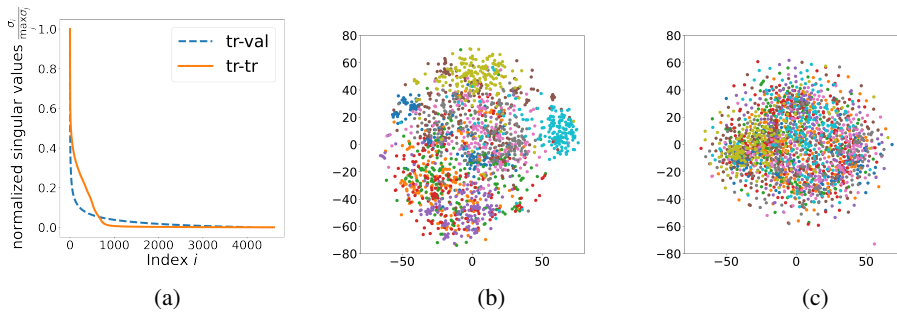


Figure 3. Plot of singular values and t-SNE plots for MiniImageNet tr-val versus tr-tr. (a): Singular values for representations produced by tr-val versus tr-tr of a CNN model trained with RepLearn on MiniImageNet. (b): t-SNE plot of representations produced by tr-val CNN model trained with RepLearn on MiniImageNet for 10 randomly selected classes of test split. (c): t-SNE plot of representations produced by tr-tr CNN model trained with RepLearn on MiniImageNet for the same 10 randomly selected classes of test split. We find that the tr-val representations appear to be more clustered than the tr-tr representations.

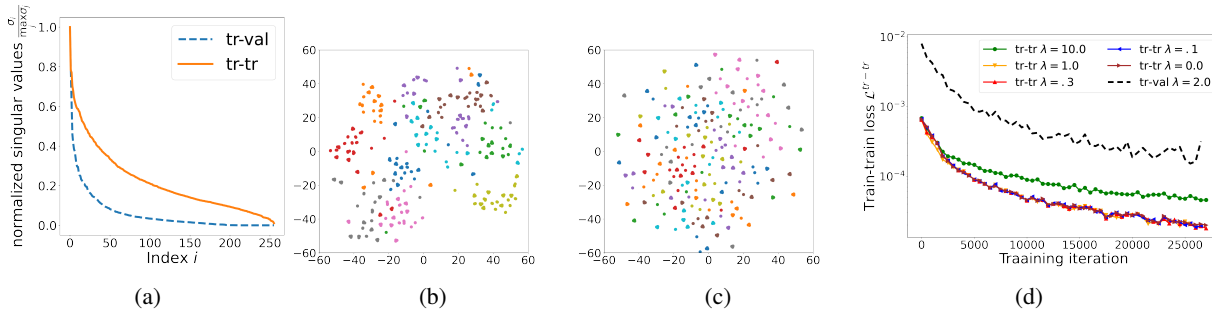


Figure 4. Plot of singular values and t-SNE plots for iMAML tr-val versus tr-tr. (a): Singular values for representations produced by tr-val versus tr-tr of a CNN model trained with iMAML on Omniglot 5-way 1-shot. (b): t-SNE plot of representations produced by tr-val CNN model trained with iMAML on Omniglot 5-way 1-shot for 10 randomly selected classes of test split. (c): t-SNE plot of representations produced by tr-tr CNN trained with iMAML model on Omniglot 5-way 1-shot for the same 10 randomly selected classes of test split. We find that the tr-val representations appear to be more clustered than the tr-tr representations. (d): Train-train loss values for tr-tr versus tr-val models.