

# Supplementary Materials for the Paper: “On Dropout and Nuclear Norm Regularization”

## A. Proofs of the Main Results

In this section, we provide the complete proofs of our main results. For notational simplicity, we define

$$\begin{aligned} W_{i \rightarrow j} &:= W_i W_{i-1} \cdots W_{j+1} W_j, \\ \bar{W}_{i \rightarrow j} &:= \frac{1}{\theta^{i-j}} W_i \text{diag}(\mathbf{b}_{i-1}) W_{i-1} \cdots \text{diag}(\mathbf{b}_{j+1}) W_{j+1} \text{diag}(\mathbf{b}_j) W_j. \end{aligned}$$

Since the Bernoulli random vectors are i.i.d., it holds that  $\mathbb{E}_{\mathbf{b}_i}[\bar{W}_{i \rightarrow j}] = W_{i \rightarrow j}$ . A quantity that shows up when analyzing dropout training under squared error is  $\mathbb{E}[\|\mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V} \mathbf{x}\|^2]$ , where the expectation is taken with respect to  $\mathbf{b}$ , which is a Bernoulli random vector with parameter  $\theta$ . The following lemma gives the closed form of this expectation.

**Lemma A.1.** *Let  $U \in \mathbb{R}^{d_2 \times r}$ ,  $V \in \mathbb{R}^{d_1 \times r}$ , and  $\mathbf{C} := \mathbb{E}[\mathbf{x} \mathbf{x}^\top]$ . It holds that*

$$\mathbb{E}[\|\mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V} \mathbf{x}\|^2] = \theta^2 \mathbb{E}[\|\mathbf{U} \mathbf{V} \mathbf{x}\|^2] + (\theta - \theta^2) \sum_{j=1}^r \|u_{\cdot j}\|^2 \|\mathbf{C}^{\frac{1}{2}} v_{j\cdot}\|^2.$$

The proof can be found in (Cavazza et al., 2018; Mianjy et al., 2018). Nonetheless, we provide a proof here for completeness.

*Proof of Lemma A.1.*

$$\begin{aligned} \mathbb{E}[\|\mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V} \mathbf{x}\|^2] &= \mathbb{E} \sum_{i=1}^{d_2} \mathbb{E}_{\mathbf{b}} \left( \sum_{j=1}^r u_{ij} b_j v_{j\cdot}^\top \mathbf{x} \right)^2 = \mathbb{E} \sum_{i=1}^{d_2} \mathbb{E}_{\mathbf{b}} \left[ \sum_{j,k=1}^r u_{ij} u_{ik} b_j b_k (v_{j\cdot}^\top \mathbf{x})(v_{k\cdot}^\top \mathbf{x}) \right] \\ &= \mathbb{E} \sum_{i=1}^{d_2} \sum_{j,k=1}^r u_{ij} u_{ik} (\theta^2 \mathbf{1}_{j \neq k} + \theta \mathbf{1}_{j=k}) (v_{j\cdot}^\top \mathbf{x})(v_{k\cdot}^\top \mathbf{x}) = \theta^2 \mathbb{E}[\|\mathbf{U} \mathbf{V} \mathbf{x}\|^2] + (\theta - \theta^2) \mathbb{E} \sum_{i=1}^{d_2} \sum_{j=1}^r u_{ij}^2 (v_{j\cdot}^\top \mathbf{x})^2 \\ &= \theta^2 \mathbb{E}[\|\mathbf{U} \mathbf{V} \mathbf{x}\|^2] + (\theta - \theta^2) \sum_{j=1}^r \mathbb{E}[v_{j\cdot}^\top \mathbf{x} \mathbf{x}^\top v_{j\cdot}] \sum_{i=1}^{d_2} u_{ij}^2 = \theta^2 \mathbb{E}[\|\mathbf{U} \mathbf{V} \mathbf{x}\|^2] + (\theta - \theta^2) \sum_{j=1}^r \|\mathbf{C}^{\frac{1}{2}} v_{j\cdot}\|^2 \|u_{\cdot j}\|^2. \end{aligned}$$

□

### A.1. Properties of the explicit regularizer

Recall that training a network with dropout aims at minimizing the following *dropout objective*

$$L_\theta(\{\mathbf{W}_i\}_{i=1}^{k+1}) = \mathbb{E}_{\substack{\mathbf{b}_i \sim \text{Bern}(\theta) \\ (\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}} \left[ \|\mathbf{y} - \frac{1}{\theta^k} \mathbf{W}_{k+1} \text{diag}(\mathbf{b}_k) \mathbf{W}_k \cdots \text{diag}(\mathbf{b}_1) \mathbf{W}_1 \mathbf{x}\|^2 \right].$$

In Proposition 2.1 we show that this objective can be decomposed into a summation of the population loss plus an *explicit regularizer*, i.e.  $L_\theta(\cdot) = L(\cdot) + R(\cdot)$ , and give the closed form expression for the explicit regularizer.

*Proof of Proposition 2.1.* We start by expanding the squared error:

$$\begin{aligned}
 L_\theta(\{\mathbf{W}_i\}_{i=1}^{k+1}) &= \mathbb{E}_{\substack{\mathbf{b}_i \sim \text{Bern}(\theta) \\ (x, y) \sim \mathcal{D}}} [\|y - \bar{\mathbf{W}}_{k+1 \rightarrow 1} \mathbf{x}\|^2] \\
 &= \mathbb{E}[\|y\|^2] - 2\mathbb{E}[\langle \bar{\mathbf{W}}_{k+1 \rightarrow 1} \mathbf{x}, y \rangle] + \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 1} \mathbf{x}\|^2] \\
 &= \mathbb{E}[\|y\|^2] - 2\mathbb{E}[\langle \mathbf{W}_{k+1 \rightarrow 1} \mathbf{x}, y \rangle] + \frac{1}{\theta^{2k}} \mathbb{E}[\|\mathbf{W}_{k+1} \text{diag}(\mathbf{b}_k) \mathbf{W}_k \dots \text{diag}(\mathbf{b}_1) \mathbf{W}_1 \mathbf{x}\|^2] \quad (10)
 \end{aligned}$$

We now focus on the last term in the right hand side of Equation (10).

$$\begin{aligned}
 \mathbb{E}[\|\mathbf{W}_{k+1} \text{diag}(\mathbf{b}_k) \mathbf{W}_k \dots \text{diag}(\mathbf{b}_1) \mathbf{W}_1 \mathbf{x}\|^2] &= \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 2} \text{diag}(\mathbf{b}_1) \mathbf{W}_1 \mathbf{x}\|^2] \\
 &= \theta^2 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 2} \mathbf{W}_1 \mathbf{x}\|^2] + (\theta - \theta^2) \sum_{j=1}^{d_1} \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 2}(:, j)\|^2] \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_1(j, :)\|^2 \quad (11)
 \end{aligned}$$

The second equality follows from Lemma A.1. Similarly, the first term on the right hand side of Equation (11) can be expressed as:

$$\begin{aligned}
 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 2} \mathbf{W}_1 \mathbf{x}\|^2] &= \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 3} \text{diag}(\mathbf{b}_2) \mathbf{W}_2 \rightarrow 1 \mathbf{x}\|^2] \\
 &= \theta^2 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 3} \mathbf{W}_2 \rightarrow 1 \mathbf{x}\|^2] + (\theta - \theta^2) \sum_{j=1}^{d_2} \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 3}(:, j)\|^2] \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_2 \rightarrow 1(j, :)\|^2
 \end{aligned}$$

By recursive application of the above identity and plugging the result into Equation (11), we obtain:

$$\mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow 1} \mathbf{x}\|^2] = \theta^{2k} \mathbb{E}[\|\mathbf{W}_{k+1 \rightarrow 1} \mathbf{x}\|^2] + (1 - \theta) \sum_{i=1}^k \sum_{j=1}^{d_i} \theta^{2i-1} \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+1}(:, j)\|^2] \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \quad (12)$$

Plugging back the above equality into Equation (10), we get

$$\begin{aligned}
 L_\theta(\{\mathbf{W}_i\}) &= \|y\|^2 - 2\mathbb{E}[\langle \mathbf{W}_{k+1 \rightarrow 1} \mathbf{x}, y \rangle] + \mathbb{E}[\|\mathbf{W}_{k+1 \rightarrow 1} \mathbf{x}\|^2] + \frac{1 - \theta}{\theta^{2k}} \sum_{i=1}^k \sum_{j=1}^{d_i} \theta^{2i-1} \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+1}(:, j)\|^2] \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \\
 &= \mathbb{E}_x[\|y - \mathbf{W}_{k+1 \rightarrow 1} \mathbf{x}\|^2] + (1 - \theta) \sum_{i=1}^k \sum_{j=1}^{d_i} \theta^{2(i-k)-1} \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+1}(:, j)\|^2] \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2. \quad (13)
 \end{aligned}$$

It remains to calculate the terms of the form  $\mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+1}(:, j)\|^2]$  in the right hand side of Equation (13). We introduce the variable  $x \sim \mathcal{N}(0, 1)$  so that we can use Lemma A.1 again:

$$\begin{aligned}
 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+1}(:, j)\|^2] &= \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+2} \text{diag}(\mathbf{b}_{i+1}) \mathbf{W}_{i+1}(:, j) x\|^2] \\
 &= \theta^2 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+2} \mathbf{W}_{i+1}(:, j)\|^2] + (\theta - \theta^2) \mathbb{E} \sum_{l=1}^{d_{i+1}} \|\bar{\mathbf{W}}_{k+1 \rightarrow i+2}(:, l)\|^2 \mathbf{W}_{i+1}(l, j)^2. \quad (14)
 \end{aligned}$$

The first term on the right hand side of Equation (14) can be expanded as:

$$\begin{aligned}
 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+2} \mathbf{W}_{i+1}(:, j)\|^2] &= \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+3} \text{diag}(\mathbf{b}_{i+2}) \mathbf{W}_{i+2 \rightarrow i+1}(:, j) x\|^2] \\
 &= \theta^2 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+3} \mathbf{W}_{i+2 \rightarrow i+1}(:, j)\|^2] + (\theta - \theta^2) \mathbb{E} \sum_{l=1}^{d_{i+2}} \|\bar{\mathbf{W}}_{k+1 \rightarrow i+3}(:, l)\|^2 \mathbf{W}_{i+2 \rightarrow i+1}(l, j)^2
 \end{aligned}$$

By recursive application of the above equality and plugging the results into Equation (14), we get

$$\mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+1}(:, j)\|^2] = \theta^{2(k-i)} \|\mathbf{W}_{k+1 \rightarrow i+1}(:, j)\|^2 + (1 - \theta) \sum_{m=1}^{k-i} \theta^{2m-1} \mathbb{E} \sum_{l=1}^{d_{i+m}} \|\bar{\mathbf{W}}_{k+1 \rightarrow i+1+m}(:, l)\|^2 \mathbf{W}_{i+m \rightarrow i+1}(l, j)^2$$

Plugging back the above identity into Equation (13) we get

$$\begin{aligned}
 R(\{\mathbf{W}_i\}) &= (1-\theta) \sum_{i=1}^k \sum_{j=1}^{d_i} \theta^{2(i-k)-1} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \mathbb{E}[\|\bar{\mathbf{W}}_{k+1 \rightarrow i+1}(:, j)\|^2] \\
 &= \frac{1-\theta}{\theta} \sum_{i=1}^k \sum_{j=1}^{d_i} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \|\mathbf{W}_{k+1 \rightarrow i+1}(:, j)\|^2 \\
 &\quad + \left(\frac{1-\theta}{\theta}\right)^2 \sum_{i=1}^k \sum_{j=1}^{d_i} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \sum_{m=1}^{k-i} \theta^{2(i+m-k)} \mathbb{E} \sum_{l=1}^{d_{i+m}} \mathbf{W}_{i+m \rightarrow i+1}(l, j)^2 \|\bar{\mathbf{W}}_{k+1 \rightarrow i+m+1}(:, l)\|^2 \\
 &= \frac{1-\theta}{\theta} \sum_{i=1}^k \sum_{j=1}^{d_i} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \|\mathbf{W}_{k+1 \rightarrow i+1}(:, j)\|^2 \\
 &\quad + \left(\frac{1-\theta}{\theta}\right)^2 \sum_{i=1}^k \sum_{j=1}^{d_i} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \sum_{m=1}^{k-i} \sum_{l=1}^{d_{i+m}} \mathbf{W}_{i+m \rightarrow i+1}(l, j)^2 \|\mathbf{W}_{k+1 \rightarrow i+m+1}(:, l)\|^2 \\
 &\quad + \left(\frac{1-\theta}{\theta}\right)^3 \sum_{i=1}^k \sum_{j=1}^{d_i} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{i \rightarrow 1}(j, :)\|^2 \sum_{m=1}^{k-i} \sum_{l=1}^{d_{i+m}} \mathbf{W}_{i+m \rightarrow i+1}(l, j)^2 \left( \sum_{mm=1}^{k-i-m} \theta^{2(i+m+mm)} \right. \\
 &\quad \left. \sum_{ll=1}^{d_{i+m+mm}} \mathbf{W}_{i+m+mm \rightarrow i+1}(ll, l)^2 \mathbb{E} \|\bar{\mathbf{W}}_{k+1 \rightarrow i+1+m+mm}(:, ll)\|^2 \right) \\
 &= \dots \\
 &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_1, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_1, \dots, i_1) \\ \in [d_{j_1}] \times \dots \times [d_{j_1}]}} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\|^2 \prod_{p=1 \dots l-1} \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p)^2 \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\|^2,
 \end{aligned}$$

which completes the proof.  $\square$

**Lemma A.2.** [Properties of  $R$  and  $\Theta$ ] The following statements hold true:

1. All sub-regularizers, and hence the explicit regularizer, are re-scaling invariant.
2. The infimum in Equation (2) is always attained.
3. If  $\mathbf{C} = \mathbf{I}$ , then  $\Theta(\mathbf{M})$  is a spectral function, i.e. if  $\mathbf{M}$  and  $\mathbf{M}'$  have the same singular values, then  $\Theta(\mathbf{M}) = \Theta(\mathbf{M}')$ .

*Proof of Lemma A.2.* First, it is easy to see that the explicit regularizer and the sub-regularizers are all *re-scaling invariant*. For any sequence of scalars  $\{\alpha_i\}$  such that  $\prod_{i=1}^{k+1} \alpha_i = 1$ , let  $\bar{\mathbf{W}}_i := \alpha_i \mathbf{W}_i$ . Then it holds that:

$$\begin{aligned}
 R_l(\{\bar{\mathbf{W}}_i\}) &= \sum_{\substack{(j_1, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_1, \dots, i_1) \\ \in [d_{j_1}] \times \dots \times [d_{j_1}]}} \left\| \prod_{q=1}^{j_1} \alpha_q \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\right\|^2 \prod_{p \in [l-1]} \prod_{q=j_p+1}^{j_{p+1}} \alpha_q^2 \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p)^2 \left\| \prod_{q=j_l+1}^{k+1} \alpha_q \mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\right\|^2 \\
 &= \prod_{q=1}^{k+1} \alpha_q^2 \sum_{\substack{(j_1, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_1, \dots, i_1) \\ \in [d_{j_1}] \times \dots \times [d_{j_1}]}} \|\mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\|^2 \prod_{p \in [l-1]} \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p)^2 \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\|^2 \\
 &= R_l(\{\mathbf{W}_i\})
 \end{aligned}$$

Therefore, without loss of generality, we can express the induced regularizer as follows:

$$\Theta(\mathbf{M}) := \inf_{\substack{\mathbf{W}_{k+1} \dots \mathbf{W}_1 = \mathbf{M} \\ \|\mathbf{W}_i\|_F \leq \|\mathbf{M}\|_F}} R(\{\mathbf{W}_i\}) \quad (15)$$

Note that  $R(\{\mathbf{W}_i\})$  is a continuous function and the feasible set  $\mathcal{F} := \{(\mathbf{W}_i)_{i=1}^{k+1} : \mathbf{W}_{k+1} \cdots \mathbf{W}_1 = \mathbf{M}, \|\mathbf{W}_i\|_F \leq \|\mathbf{M}\|_F\}$  is compact. Hence, by Weierstrass extreme value theorem, the infimum is attained.

Now let  $\mathbf{U} \in \mathbb{R}^{d_{k+1} \times d_{k+1}}$  and  $\mathbf{V} \in \mathbb{R}^{d_0 \times d_0}$  be a pair of rotation matrices, i.e.  $\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{I}$  and  $\mathbf{V}^\top \mathbf{V} = \mathbf{V} \mathbf{V}^\top = \mathbf{I}$ . When the data is isotropic, i.e.  $\mathbf{C} = \mathbf{I}$ , the following equalities hold

$$\begin{aligned} R(\{\mathbf{W}_i\}) &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_l, \dots, i_1) \\ \in [d_{j_l}] \times \cdots \times [d_{j_1}]}} \|\mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\|^2 \prod_{p=1 \cdots l-1} \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p)^2 \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\|^2 \\ &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_l, \dots, i_1) \\ \in [d_{j_l}] \times \cdots \times [d_{j_1}]}} \|\mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\|^\top \mathbf{V}\|^2 \prod_{p=1 \cdots l-1} \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p)^2 \|\mathbf{U}^\top \mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\|^2 \\ &= R(\mathbf{U}^\top \mathbf{W}_{k+1}, \mathbf{W}_k, \dots, \mathbf{W}_2, \mathbf{W}_1 \mathbf{V}) \end{aligned}$$

That is,  $R(\mathbf{U}^\top \mathbf{W}_{k+1}, \mathbf{W}_k, \dots, \mathbf{W}_2, \mathbf{W}_1 \mathbf{V}) = R(\mathbf{W}_{k+1}, \mathbf{W}_k, \dots, \mathbf{W}_2, \mathbf{W}_1)$  for all rotation matrices  $\mathbf{U}$  and  $\mathbf{V}$ . In particular, let  $\mathbf{U}, \mathbf{V}$  be the left and right singular vectors of  $\mathbf{M}$ , i.e.  $\mathbf{M} = \mathbf{U} \Sigma \mathbf{V}^\top$ . To prove that  $\Theta$  is a spectral function, we need to show that  $\Theta(\mathbf{M}) = \Theta(\Sigma)$ . Let  $\{\bar{\mathbf{W}}_i\}, \{\tilde{\mathbf{W}}_i\}$  be such that  $\Theta(\mathbf{M}) = R(\{\bar{\mathbf{W}}_i\}), \Theta(\Sigma) = R(\{\tilde{\mathbf{W}}_i\})$ . Note that such weight matrices always exist since the infimum is always attained. Then

$$\Theta(\Sigma) = \Theta(\mathbf{U}^\top \mathbf{M} \mathbf{V}) \leq R(\mathbf{U}^\top \bar{\mathbf{W}}_{k+1}, \bar{\mathbf{W}}_k, \dots, \bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1 \mathbf{V}) = R(\bar{\mathbf{W}}_{k+1}, \bar{\mathbf{W}}_k, \dots, \bar{\mathbf{W}}_2, \bar{\mathbf{W}}_1) = \Theta(\mathbf{M}).$$

Similarly, we have that  $\Theta(\mathbf{M}) \leq R(\mathbf{U}^\top \tilde{\mathbf{W}}_{k+1}, \tilde{\mathbf{W}}_k, \dots, \tilde{\mathbf{W}}_2, \tilde{\mathbf{W}}_1 \mathbf{V}) = R(\tilde{\mathbf{W}}_{k+1}, \tilde{\mathbf{W}}_k, \dots, \tilde{\mathbf{W}}_2, \tilde{\mathbf{W}}_1) = \Theta(\Sigma)$ , which completes the proof.  $\square$

## A.2. The induced regularizer and its convex envelope

*Proof of Theorem 2.6.* By Lemma 3.1, for any architecture, any dropout rate, and any set of weights  $\{\mathbf{W}_i\}$  that implements a network map  $\mathbf{W}_{k+1 \rightarrow 1}$ , the explicit regularizer is lower bounded by the effective regularization parameter times the product of the squared nuclear norm of the network map and the principal squared root of the second moment of  $\mathbf{x}$ , i.e.  $R(\{\mathbf{W}_i\}) \geq \nu_{\{d_i\}} \|\mathbf{W}_{k+1 \rightarrow 1} \mathbf{C}^{\frac{1}{2}}\|_*^2$ . Consequently, the induced regularizer can also be lower bounded as  $\Theta(\mathbf{M}) \geq \nu_{\{d_i\}} \|\mathbf{M} \mathbf{C}^{\frac{1}{2}}\|_*^2$ . On the other hand, Lemma 3.2 establishes that  $\Theta^{**}(\mathbf{M}) \leq \nu_{\{d_i\}} \|\mathbf{M} \mathbf{C}^{\frac{1}{2}}\|_*^2$  holds for any network map  $\mathbf{M}$ . Putting these two inequalities together, we arrive at

$$\Theta^{**}(\mathbf{M}) \leq \nu_{\{d_i\}} \|\mathbf{M} \mathbf{C}^{\frac{1}{2}}\|_*^2 \leq \Theta(\mathbf{M}).$$

Since  $\Theta^{**}(\mathbf{M})$  is the largest convex underestimator of  $\Theta(\mathbf{M})$ , and the squared nuclear norm is a convex function, we conclude that  $\Theta^{**}(\mathbf{M}) = \nu_{\{d_i\}} \|\mathbf{M}\|_*^2$ .  $\square$

Despite the complex form of the explicit regularizer given in Proposition 2.1, we can show that it is always lower bounded by *effective regularization parameter* times  $\|\mathbf{M} \mathbf{C}^{\frac{1}{2}}\|_*^2$ . This result is given by Lemma 3.1.

*Proof of Lemma 3.1.* Recall that the explicit regularizer  $R(\{\mathbf{W}_i\})$  is composed of  $k$  sub-regularizers

$$R(\{\mathbf{W}_i\}) = R_1(\{\mathbf{W}_i\}) + R_2(\{\mathbf{W}_i\}) + \cdots + R_k(\{\mathbf{W}_i\}).$$

The  $l$ -th sub-regularizer  $R_l(\{\mathbf{W}_i\})$  can be written in the form of:

$$R_l(\{\mathbf{W}_i\}) = \lambda^l \sum_{(j_l, \dots, j_1) \in \binom{[k]}{l}} R_{\{j_l, \dots, j_1\}}(\{\mathbf{W}_i\})$$

where

$$R_{\{j_l, \dots, j_1\}}(\{\mathbf{W}_i\}) := \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\|^2 \prod_{p=1 \cdots l-1} \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p)^2 \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\|^2.$$

The following set of equalities hold true:

$$\begin{aligned}
 R_{\{j_l, \dots, j_1\}}(\{\mathbf{W}_i\}) &= \sum_{i_l, \dots, i_1} \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\|^2 \mathbf{W}_{j_l \rightarrow j_{l-1}+1}(i_l, i_{l-1})^2 \cdots \mathbf{W}_{j_2 \rightarrow j_1+1}(i_2, i_1)^2 \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\|^2 \\
 &\geq \frac{1}{\prod_{i \in [l]} d_{j_i}} \left( \sum_{i_l, \dots, i_1} \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\| \|\mathbf{W}_{j_l \rightarrow j_{l-1}+1}(i_l, i_{l-1})\| \cdots \|\mathbf{W}_{j_2 \rightarrow j_1+1}(i_2, i_1)\| \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\| \right)^2 \\
 &= \frac{1}{\prod_{i \in [l]} d_{j_i}} \left( \sum_{i_l, \dots, i_1} \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l) \mathbf{W}_{j_l \rightarrow j_{l-1}+1}(i_l, i_{l-1}) \cdots \mathbf{W}_{j_2 \rightarrow j_1+1}(i_2, i_1) \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\| \mathbf{C}^{\frac{1}{2}} \right)^2 \\
 &\geq \frac{1}{\prod_{i \in [l]} d_{j_i}} \left\| \sum_{i_l, \dots, i_1} \mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l) \mathbf{W}_{j_l \rightarrow j_{l-1}+1}(i_l, i_{l-1}) \cdots \mathbf{W}_{j_2 \rightarrow j_1+1}(i_2, i_1) \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\right\| \mathbf{C}^{\frac{1}{2}} \right\|_*^2 \\
 &= \frac{1}{\prod_{i \in [l]} d_{j_i}} \left\| \sum_{i_l, i_1} \mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l) \left( \sum_{i_{l-1}, \dots, i_2} \mathbf{W}_{j_l \rightarrow j_{l-1}+1}(i_l, i_{l-1}) \cdots \mathbf{W}_{j_2 \rightarrow j_1+1}(i_2, i_1) \right) \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\right\| \mathbf{C}^{\frac{1}{2}} \right\|_*^2 \\
 &= \frac{1}{\prod_{i \in [l]} d_{j_i}} \left\| \sum_{i_l, i_1} \mathbf{W}_{k+1 \rightarrow j_l}(:, i_l) \mathbf{W}_{j_l-1 \rightarrow j_1}(i_l, i_1) \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\right\| \mathbf{C}^{\frac{1}{2}} \right\|_*^2 \\
 &= \frac{1}{\prod_{i \in [l]} d_{j_i}} \|\mathbf{W}_{k+1} \cdots \mathbf{W}_1 \mathbf{C}^{\frac{1}{2}}\|_*^2
 \end{aligned}$$

where the first inequality follows due to the Cauchy-Schwartz inequality, and the second inequality follows from the triangle inequality for the matrix norms. The inequality holds with equality if and only if all the summands inside the summation are equal to each other, and sum up to  $\frac{\|\mathbf{W}_{k+1 \rightarrow 1} \mathbf{C}^{\frac{1}{2}}\|_*}{\prod_{i \in [l]} d_{j_i}}$ , i.e. when

$$\|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\| \|\mathbf{W}_{j_l \rightarrow j_{l-1}+1}(i_l, i_{l-1})\| \cdots \|\mathbf{W}_{j_2 \rightarrow j_1+1}(i_2, i_1)\| \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\| = \frac{1}{\prod_{i \in [l]} d_{j_i}} \|\mathbf{W}_{k+1 \rightarrow 1} \mathbf{C}^{\frac{1}{2}}\|_*$$

for all  $(i_l, \dots, i_1) \in [d_{j_l}] \times \cdots \times [d_{j_1}]$ . This lowerbound holds for all  $l \in [k]$ , and for all  $(j_l, \dots, j_1) \in \binom{[k]}{l}$ . Thus, we get the following lowerbound on the regularizer:

$$R(\{\mathbf{W}_i\}) \geq \sum_{l \in [k]} \lambda^l \sum_{(j_l, \dots, j_1) \in \binom{[k]}{l}} \frac{1}{\prod_{i \in [l]} d_{j_i}} \|\mathbf{W}_{k+1 \rightarrow 1} \mathbf{C}^{\frac{1}{2}}\|_*^2 = \nu_{\{d_i\}} \|\mathbf{W}_{k+1 \rightarrow 1} \mathbf{C}^{\frac{1}{2}}\|_*^2$$

which completes the proof.  $\square$

Not only  $\nu_{\{d_i\}} \|\mathbf{M} \mathbf{C}^{\frac{1}{2}}\|_*^2$  is a lowerbound for the induced regularizer, but also is an upperbound for its convex envelope. We prove this result in Lemma 3.2.

*Proof of Lemma 3.2.* The induced regularizer is non-negative. Hence, the domain of the Fenchel dual of the induced regularizer is the whole  $\mathbb{R}^{d_{k+1} \times d_0}$ . The Fenchel dual of the induced regularizer  $\Theta(\cdot)$  is given by:

$$\begin{aligned}
 \Theta^*(\mathbf{M}) &= \max_{\mathbf{P}} \langle \mathbf{M}, \mathbf{P} \rangle - \Theta(\mathbf{P}) \\
 &= \max_{\mathbf{P}} \langle \mathbf{M}, \mathbf{P} \rangle - \min_{\substack{\{\mathbf{W}_i\} \\ \mathbf{W}_{k+1 \rightarrow 1} = \mathbf{P}}} R(\{\mathbf{W}_i\}) \\
 &= \max_{\{\mathbf{W}_i\}} \langle \mathbf{M}, \mathbf{W}_{k+1 \rightarrow 1} \rangle - R(\{\mathbf{W}_i\}). \tag{16}
 \end{aligned}$$

Define  $\Phi(\{\mathbf{W}_i\}) := \langle \mathbf{M}, \mathbf{W}_{k+1 \rightarrow 1} \rangle - R(\{\mathbf{W}_i\})$  as the objective in the right hand side of Equation (16). Due to the complicated products of the norms of the weights in the regularizer, maximizing  $\Phi$  with respect to  $\{\mathbf{W}_i\}$  is a daunting task. Here, we find a lower bound on this maximum value. Let  $\mathbf{W}_{k+1}^\alpha := \alpha \mathbf{u}_1 \mathbf{1}_{d_k}^\top$  and  $\mathbf{W}_1^\alpha := \mathbf{1}_{d_1} \mathbf{v}_1^\top \mathbf{C}^{-\frac{1}{2}}$ , where  $(\mathbf{u}_1, \mathbf{v}_1)$  is

the top singular vectors of  $\mathbf{MC}^{-\frac{1}{2}}$ , and  $\mathbf{1}_d$  is the  $d$ -dimensional vector of all 1s. Furthermore, let  $\mathbf{W}_i^\alpha := \mathbf{1}_{d_i} \mathbf{1}_{d_{i-1}}^\top$ , for all  $i \in \{2, \dots, k\}$ . Note that

$$\Theta^*(\mathbf{M}) = \max_{\{\mathbf{W}_i\}} \Phi(\{\mathbf{W}_i\}) \geq \max_{\alpha} \Phi(\{\mathbf{W}_i^\alpha\}).$$

We now simplify  $\Phi(\{\mathbf{W}_i^\alpha\})$ . First, the following equalities hold for the  $\langle \mathbf{M}, \mathbf{W}_{k+1 \rightarrow 1}^\alpha \rangle$ :

$$\begin{aligned} \langle \mathbf{M}, \mathbf{W}_{k+1 \rightarrow 1}^\alpha \rangle &= \sum_{(i_{k+1}, \dots, i_1) \in [d_{k+1}] \times \dots \times [d_1]} \langle \mathbf{M}, \mathbf{W}_{k+1}^\alpha(:, i_k) \prod_{j=\{k-1, \dots, 1\}} \mathbf{W}_{j+1}^\alpha(i_{j+1}, i_j) \mathbf{W}_1^\alpha(i_1, :)^T \rangle \\ &= \sum_{(i_{k+1}, \dots, i_1) \in [d_{k+1}] \times \dots \times [d_1]} \mathbf{W}_{k+1}^\alpha(:, i_k)^\top \mathbf{M} \mathbf{W}_1^\alpha(i_1, :) \\ &= \sum_{(i_{k+1}, \dots, i_1) \in [d_{k+1}] \times \dots \times [d_1]} \alpha \mathbf{u}_1^\top \mathbf{MC}^{-\frac{1}{2}} \mathbf{v}_1 \\ &= \sum_{(i_{k+1}, \dots, i_1) \in [d_{k+1}] \times \dots \times [d_1]} \alpha \|\mathbf{MC}^{-\frac{1}{2}}\|_2 \\ &= \alpha \|\mathbf{MC}^{-\frac{1}{2}}\|_2 \prod_{j \in [k]} d_j =: \alpha \|\mathbf{MC}^{-\frac{1}{2}}\|_2 D. \end{aligned}$$

The following terms show up in the expansion of the regularizer:

$$\begin{aligned} \mathbf{W}_{j_1 \rightarrow 1}^\alpha(i_1, :)^T &= \mathbf{W}_{j_1}^\alpha(i_1, :) \mathbf{W}_{j_1-1}^\alpha \cdots \mathbf{W}_2^\alpha \mathbf{W}_1^\alpha = \mathbf{1}_{d_{j_1-1}}^\top \mathbf{1}_{d_{j_1-1}} \mathbf{1}_{d_{j_1-2}}^\top \cdots \mathbf{1}_{d_2} \mathbf{1}_{d_1}^\top \mathbf{1}_{d_1} \mathbf{v}_1^\top = \prod_{i \in [j_1-1]} d_i \mathbf{v}_1^\top \mathbf{C}^{-\frac{1}{2}} \\ \mathbf{W}_{j_{p+1} \rightarrow j_p+1}^\alpha(i_{p+1}, i_p) &= \mathbf{W}_{j_{p+1}}^\alpha(i_{p+1}, :) \mathbf{W}_{j_{p+1}-1}^\alpha \cdots \mathbf{W}_{j_p+2}^\alpha \mathbf{W}_{j_p+1}^\alpha(:, i_p) \\ &= \mathbf{1}_{d_{j_{p+1}-1}}^\top \mathbf{1}_{d_{j_{p+1}-1}} \mathbf{1}_{d_{j_{p+1}-2}}^\top \cdots \mathbf{1}_{d_{j_p+2}} \mathbf{1}_{d_{j_p+1}}^\top \mathbf{1}_{d_{j_p+1}} = \prod_{i \in \{j_p+1, \dots, j_{p+1}-1\}} d_i \\ \mathbf{W}_{k+1 \rightarrow j_l+1}^\alpha(:, i_l) &= \alpha \mathbf{W}_{k+1}^\alpha \mathbf{W}_k^\alpha \cdots \mathbf{W}_{j_l+2}^\alpha \mathbf{W}_{j_l+1}^\alpha(:, i_l) = \alpha \mathbf{u}_1 \mathbf{1}_{d_k}^\top \mathbf{1}_{d_k} \mathbf{1}_{d_{k-1}}^\top \cdots \mathbf{1}_{d_{j_l+2}} \mathbf{1}_{d_{j_l+1}}^\top \mathbf{1}_{d_{j_l+1}} = \alpha \prod_{i \in \{j_l+1, \dots, k\}} d_i \mathbf{u}_1 \end{aligned}$$

With the above equalities, the explicit regularizer reduces to:

$$\begin{aligned} R(\{\mathbf{W}_i^\alpha\}) &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{(i_l, \dots, i_1) \in [d_{j_l}] \times \dots \times [d_{j_1}]} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{j_1 \rightarrow 1}^\alpha(i_1, :)\|^2 \prod_{p=1 \dots l-1} \mathbf{W}_{j_{p+1} \rightarrow j_p+1}^\alpha(i_{p+1}, i_p)^2 \|\mathbf{W}_{k+1 \rightarrow j_l+1}^\alpha(:, i_l)\|^2 \\ &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{(i_l, \dots, i_1) \in [d_{j_l}] \times \dots \times [d_{j_1}]} \|\mathbf{C}^{\frac{1}{2}} \mathbf{C}^{-\frac{1}{2}} \mathbf{v}_1 \prod_{i \in [j_1-1]} d_i\|^2 \prod_{p=1 \dots l-1} \prod_{i \in \{j_p+1, \dots, j_{p+1}-1\}} d_i^2 \|\alpha \mathbf{u}_1 \prod_{i \in \{j_l+1, \dots, k\}} d_i\|^2 \\ &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{(i_l, \dots, i_1) \in [d_{j_l}] \times \dots \times [d_{j_1}]} \prod_{i \in [j_1-1]} d_i^2 \prod_{p=1 \dots l-1} \prod_{i \in \{j_p+1, \dots, j_{p+1}-1\}} d_i^2 \alpha^2 \prod_{i \in \{j_l+1, \dots, k\}} d_i^2 \\ &= \alpha^2 \sum_{l=1}^k \lambda^l \sum_{(j_l, \dots, j_1) \in \binom{[k]}{l}} \sum_{(i_l, \dots, i_1) \in [d_{j_l}] \times \dots \times [d_{j_1}]} \frac{\prod_{i \in [k]} d_i^2}{\prod_{i \in [l]} d_{j_i}^2} =: \alpha^2 \rho \end{aligned}$$

Plugging back the above equalities into the definition of  $\Phi$ , we arrive at  $\Phi(\{\mathbf{W}_i^\alpha\}) = \alpha \|\mathbf{MC}^{-\frac{1}{2}}\|_2 D - \alpha^2 \rho$ . The maximum of  $\Phi(\{\mathbf{W}_i^\alpha\})$  with respect to  $\alpha$  is achieved when  $\alpha^* = \frac{\|\mathbf{MC}^{-\frac{1}{2}}\|_2 D}{2\rho}$ , in which case we have

$$\Theta^*(\mathbf{M}) \geq \Phi(\{\mathbf{W}_i^{\alpha^*}\}) = \frac{D^2}{4\rho} \|\mathbf{MC}^{-\frac{1}{2}}\|_2^2 =: \Psi(\mathbf{M}).$$

Since Fenchel dual is order reversing, we get

$$\begin{aligned}
 \Theta^{**}(\mathbf{M}) &\leq \Psi^*(\mathbf{M}) \\
 &= \frac{\rho}{D^2} \|\mathbf{MC}^{\frac{1}{2}}\|_*^2 \\
 &= \frac{\sum_{l=1}^k \lambda^l \sum_{(j_1, \dots, j_l) \in \binom{[k]}{l}} \sum_{(i_1, \dots, i_l) \in [d_{j_1}] \times \dots \times [d_{j_l}]} \frac{\prod_{i \in [k]} d_i^2}{\prod_{i \in [l]} d_{j_i}^2}}{\prod_{j \in [k]} d_j^2} \|\mathbf{MC}^{\frac{1}{2}}\|_*^2 \\
 &= \sum_{l=1}^k \lambda^l \sum_{(j_1, \dots, j_l) \in \binom{[k]}{l}} \sum_{(i_1, \dots, i_l) \in [d_{j_1}] \times \dots \times [d_{j_l}]} \frac{1}{\prod_{i \in [l]} d_{j_i}^2} \|\mathbf{MC}^{\frac{1}{2}}\|_*^2 \\
 &= \sum_{l=1}^k \lambda^l \sum_{(j_1, \dots, j_l) \in \binom{[k]}{l}} \frac{1}{\prod_{i \in [l]} d_{j_i}} \|\mathbf{MC}^{\frac{1}{2}}\|_*^2 \\
 &= \nu_{\{d_i\}} \|\mathbf{MC}^{\frac{1}{2}}\|_*^2
 \end{aligned}$$

where the first equality follows from the fact that if  $f(\mathbf{M}) = \beta \|\mathbf{MA}\|^2$  and  $\mathbf{A} \succ 0$  then  $f^*(\mathbf{M}) = \frac{1}{4\beta} \|\mathbf{MA}^{-1}\|_*^2$ . This result is standard in the literature, but we prove it here for completeness. Note that

$$\begin{aligned}
 \langle \mathbf{Y}, \mathbf{M} \rangle - \beta \|\mathbf{YA}\|^2 &= \langle \mathbf{YA}, \mathbf{MA}^{-1} \rangle - \beta \|\mathbf{YA}\|^2 \\
 &\leq \|\mathbf{YA}\| \|\mathbf{MA}^{-1}\|_* - \beta \|\mathbf{YA}\|^2
 \end{aligned}$$

where the inequality is due to Holder's identity. The right hand side above is a quadratic in  $\|\mathbf{YA}\|$  and is maximized when  $\|\mathbf{YA}\| = \frac{1}{2\beta} \|\mathbf{MA}^{-1}\|_*$ , in which case we have

$$f^*(\mathbf{M}) = \sup_{\mathbf{Y}} \langle \mathbf{Y}, \mathbf{M} \rangle - \beta \|\mathbf{YA}\|^2 = \frac{1}{2\beta} \|\mathbf{MA}^{-1}\|_* \|\mathbf{MA}^{-1}\|_* - \beta \left( \frac{1}{2\beta} \|\mathbf{MA}^{-1}\|_* \right)^2 = \frac{1}{4\beta} \|\mathbf{MA}^{-1}\|_*^2.$$

□

### A.3. Characterization of the global optima of the dropout objective

*Proof of Proposition 3.3.* When the network map has rank equal to one, it can be expressed as  $\mathbf{uv}^\top$ , where  $\mathbf{u} \in \mathbb{R}^{d_{k+1}}$  and  $\mathbf{v} \in \mathbb{R}^{d_0}$ . We show that for any architecture  $\{d_i\}$  and any network mapping  $\mathbf{uv}^\top \in \mathbb{R}^{d_{k+1} \times d_0}$ , it is always possible to represent  $\mathbf{uv}^\top = \mathbf{W}_{k+1} \cdots \mathbf{W}_1$  such that the resulting network is equalized. One such factorization is when  $\mathbf{W}_1 = \frac{1_{d_1} \mathbf{v}^\top}{\sqrt{d_1}}$ ,  $\mathbf{W}_{k+1} = \frac{\mathbf{u} 1_{d_k}^\top}{\sqrt{d_k}}$ , and  $\mathbf{W}_i = \frac{1_{d_i} 1_{d_{i-1}}^\top}{\sqrt{d_i d_{i-1}}}$  for  $i \in \{2, \dots, k\}$ . For these weight parameters, we have that

$$\begin{aligned}
 \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)^\top &= \mathbf{W}_{j_1}(i_1, :)^\top \mathbf{W}_{j_1-1} \cdots \mathbf{W}_2 \mathbf{W}_1 = \frac{1_{d_{j_1-1}}^\top}{\sqrt{d_{j_1} d_{j_1-1}}} \frac{1_{d_{j_1-1}} 1_{d_{j_1-2}}^\top}{\sqrt{d_{j_1-1} d_{j_1-2}}} \cdots \frac{1_{d_2} 1_{d_1}^\top}{\sqrt{d_2 d_1}} \frac{1_{d_1} \mathbf{v}^\top}{\sqrt{d_1}} = \frac{\mathbf{v}^\top}{\sqrt{d_{j_1}}} \\
 \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p) &= \mathbf{W}_{j_{p+1}}(i_{p+1}, :)^\top \mathbf{W}_{j_{p+1}-1} \cdots \mathbf{W}_{j_p+2} \mathbf{W}_{j_p+1}(:, i_p) \\
 &= \frac{1_{d_{j_{p+1}-1}}^\top}{\sqrt{d_{j_{p+1}} d_{j_{p+1}-1}}} \frac{1_{d_{j_{p+1}-1}} 1_{d_{j_{p+1}-2}}^\top}{\sqrt{d_{j_{p+1}-1} d_{j_{p+1}-2}}} \cdots \frac{1_{d_{j_p+2}} 1_{d_{j_p+1}}^\top}{\sqrt{d_{j_p+2} d_{j_p+1}}} \frac{1_{d_{j_p+1}}}{\sqrt{d_{j_p+1} d_{j_p}}} = \frac{1}{\sqrt{d_{j_{p+1}} d_{j_p}}} \\
 \mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l) &= \mathbf{W}_{k+1} \mathbf{W}_k \cdots \mathbf{W}_{j_l+2} \mathbf{W}_{j_l+1}(:, i_l) \\
 &= \frac{\mathbf{u} 1_{d_k}^\top}{\sqrt{d_k}} \frac{1_{d_k} 1_{d_{k-1}}^\top}{\sqrt{d_k d_{k-1}}} \cdots \frac{1_{d_{j_l+2}} 1_{d_{j_l+1}}^\top}{\sqrt{d_{j_l+2} d_{j_l+1}}} \frac{1_{d_{j_l+1}}}{\sqrt{d_{j_l+1} d_{j_l}}} = \frac{\mathbf{u}}{\sqrt{d_{j_l}}}
 \end{aligned}$$

With the above equalities, the regularizer reduces to:

$$\begin{aligned}
 R(\{\mathbf{W}_i\}) &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_l, \dots, i_1) \\ \in [d_{j_l}] \times \dots \times [d_{j_1}]}} \|\mathbf{C}^{\frac{1}{2}} \mathbf{W}_{j_1 \rightarrow 1}(i_1, :)\|^2 \prod_{p=1 \dots l-1} \mathbf{W}_{j_{p+1} \rightarrow j_p+1}(i_{p+1}, i_p)^2 \|\mathbf{W}_{k+1 \rightarrow j_l+1}(:, i_l)\|^2 \\
 &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_l, \dots, i_1) \\ \in [d_{j_l}] \times \dots \times [d_{j_1}]}} \|\mathbf{C}^{\frac{1}{2}} \frac{\mathbf{v}}{\sqrt{d_{j_1}}}\|^2 \prod_{p=1 \dots l-1} \frac{1}{d_{j_{p+1}} d_{j_p}} \|\frac{\mathbf{u}}{\sqrt{d_{j_l}}}\|^2 \\
 &= \sum_{l=1}^k \lambda^l \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \sum_{\substack{(i_l, \dots, i_1) \\ \in [d_{j_l}] \times \dots \times [d_{j_1}]}} \frac{\|\mathbf{C}^{\frac{1}{2}} \mathbf{v}\|^2 \|\mathbf{u}\|^2}{\prod_{p \in [l]} d_{j_p}^2} \\
 &= \sum_{l=1}^k \sum_{\substack{(j_l, \dots, j_1) \\ \in \binom{[k]}{l}}} \frac{\lambda^l}{\prod_{p \in [l]} d_{j_p}} \|\mathbf{u} \mathbf{v}^\top \mathbf{C}^{\frac{1}{2}}\|_*^2 = \nu_{\{d_i\}} \|\mathbf{u} \mathbf{v}^\top \mathbf{C}^{\frac{1}{2}}\|_*^2
 \end{aligned}$$

where we used the fact that  $\|\mathbf{u}\| \|\mathbf{C}^{\frac{1}{2}} \mathbf{v}\| = \|\mathbf{u} \mathbf{v}^\top \mathbf{C}^{\frac{1}{2}}\|_*$ . Moreover, note that the network specified by the above weight matrices is equalized, since

$$|\alpha_{j_1, i_1}| \prod_{p=1 \dots l-1} |\beta_p| |\gamma_{j_p, i_p}| = \sqrt{\frac{\|\mathbf{u} \mathbf{v}^\top \mathbf{C}^{\frac{1}{2}}\|_*^2}{\prod_{p \in [l]} d_{j_p}^2}} = \frac{1}{\prod_{p \in [l]} d_{j_p}} \|\mathbf{u} \mathbf{v}^\top \mathbf{C}^{\frac{1}{2}}\|_*.$$

□

**Lemma A.3.** For any integer  $r$ , and for any  $\nu \in \mathbb{R}_+$ , it holds that

$$(\mathbf{I}_r + \nu \mathbf{I}_r \mathbf{I}_r^\top)^{-1} = \mathbf{I}_r - \frac{\nu}{1 + r\nu} \mathbf{I}_r \mathbf{I}_r^\top.$$

Lemma A.3 is an instance of the Woodbury's matrix identity. Here, we include a proof for completeness.

*Proof of Lemma A.3.* The proof simply follows from the following set of equations.

$$\begin{aligned}
 (\mathbf{I}_r + \nu \mathbf{1}_r \mathbf{1}_r^\top)(\mathbf{I}_r - \frac{\nu}{1 + r\nu} \mathbf{1}_r \mathbf{1}_r^\top) &= \mathbf{I}_r + \nu \mathbf{1}_r \mathbf{1}_r^\top - \frac{\nu}{1 + r\nu} \mathbf{1}_r \mathbf{1}_r^\top - \frac{\nu^2}{1 + r\nu} \mathbf{1}_r \mathbf{1}_r^\top \mathbf{1}_r \mathbf{1}_r^\top \\
 &= \mathbf{I}_r + \left( \nu - \frac{\nu}{1 + r\nu} - \frac{\nu^2 r}{1 + r\nu} \right) \mathbf{1}_r \mathbf{1}_r^\top = \mathbf{I}_r
 \end{aligned}$$

□

**Lemma A.4.** Consider the following optimization problem where the induced regularizer in Problem 3 is replaced with its convex envelope:

$$\min_{\mathbf{W} \in \mathbb{R}^{d_{k+1} \times d_0}} \mathbb{E}[\|\mathbf{y} - \mathbf{W}\mathbf{x}\|^2] + \Theta^{**}(\mathbf{W}), \quad \text{rank}(\mathbf{W}) \leq \min_{i \in [k+1]} d_i =: r \quad (17)$$

Define the ‘‘model’’  $\bar{\mathbf{M}} := \mathbf{C}_{\mathbf{y}\mathbf{x}} \mathbf{C}^{-\frac{1}{2}}$ . The global optimum of problem 17 is given as  $\mathbf{M}^* = \mathcal{S}_{\alpha_\rho}(\bar{\mathbf{M}}) \mathbf{C}^{-\frac{1}{2}}$ , where  $\alpha_\rho := \frac{\nu_{\{d_i\}} \sum_{j=1}^r \sigma_j(\bar{\mathbf{M}})}{1 + \rho \nu_{\{d_i\}}}$ , and  $\rho \in [\min\{r, \text{rank}(\bar{\mathbf{M}})\}]$  is the largest integer such that for all  $i \in [\rho]$ , it holds that  $\sigma_i(\bar{\mathbf{M}}) > \alpha_\rho$ .

*Proof of Lemma A.4.* Denote the objective in the optimization problem (17) as  $\mathcal{E}_{\nu_{\{d_i\}}}(\mathbf{W}) := \mathbb{E}[\|\mathbf{y} - \mathbf{W}\mathbf{x}\|^2] + \nu_{\{d_i\}} \|\mathbf{W} \mathbf{C}^{\frac{1}{2}}\|_*^2$ .



Let  $C_y := \mathbb{E}[yy^\top]$  and  $C_{xy} := \mathbb{E}[xy^\top]$ . Note that

$$\begin{aligned} \min_{\text{rank}(\mathbf{W}) \leq r} \mathcal{E}_{\nu_{\{d_i\}}}(\mathbf{W}) &= \min_{\text{rank}(\mathbf{W}) \leq r} \mathbb{E}[\|y\|^2] + \mathbb{E}[\|\mathbf{W}x\|^2] - 2\mathbb{E}[\langle y, \mathbf{W}x \rangle] + \nu_{\{d_i\}} \|\mathbf{W}C^{\frac{1}{2}}\|_*^2 \\ &\equiv \min_{\text{rank}(\mathbf{W}) \leq r} \text{Tr}(\mathbb{E}[\mathbf{W}x x^\top \mathbf{W}^\top]) - 2\text{Tr}(\mathbb{E}[\mathbf{W}xy^\top]) + \nu_{\{d_i\}} \|\mathbf{W}C^{\frac{1}{2}}\|_*^2 \\ &= \min_{\text{rank}(\mathbf{W}) \leq r} \text{Tr}(\mathbf{W}C\mathbf{W}^\top) - 2\text{Tr}(\mathbf{W}C_{xy}) + \nu_{\{d_i\}} \|\mathbf{W}C^{\frac{1}{2}}\|_*^2 \end{aligned}$$

Make the change of variable  $\bar{\mathbf{W}} \leftarrow \mathbf{W}C^{\frac{1}{2}}$  and denote  $\bar{\mathbf{M}} := C_{yx}C^{-\frac{1}{2}}$ , the goal is to solve the following problem

$$\min_{\text{rank}(\bar{\mathbf{W}}) \leq r} \text{Tr}(\bar{\mathbf{W}}\bar{\mathbf{W}}^\top) - 2\langle \bar{\mathbf{W}}, \bar{\mathbf{M}} \rangle + \nu_{\{d_i\}} \|\bar{\mathbf{W}}\|_*^2 \equiv \min_{\text{rank}(\bar{\mathbf{W}}) \leq r} \|\bar{\mathbf{M}} - \bar{\mathbf{W}}\|_F^2 + \nu_{\{d_i\}} \|\bar{\mathbf{W}}\|_*^2 \quad (18)$$

If  $\bar{\mathbf{W}}$  is a solution to the above problem, then a solution to the original problem in Equation (17) is given as  $\bar{\mathbf{W}}C^{-\frac{1}{2}}$ . Following (Cavazza et al., 2018; Mianjy et al., 2018), we show that the global optimum of Problem 18 is given in terms of an appropriate shrinkage-thresholding on the spectrum of  $\bar{\mathbf{M}}$ . Define  $r' := \max\{\text{rank}(\bar{\mathbf{M}}), r\}$ . Let  $\bar{\mathbf{M}} = U_{\bar{\mathbf{M}}} \Sigma_{\bar{\mathbf{M}}} V_{\bar{\mathbf{M}}}^\top$  and  $\bar{\mathbf{W}} = U_{\bar{\mathbf{W}}} \Sigma_{\bar{\mathbf{W}}} V_{\bar{\mathbf{W}}}^\top$  be rank- $r'$  SVDs of  $\bar{\mathbf{M}}$  and  $\bar{\mathbf{W}}$  respectively, such that  $\sigma_i(\bar{\mathbf{M}}) \geq \sigma_{i+1}(\bar{\mathbf{M}})$  and  $\sigma_i(\bar{\mathbf{W}}) \geq \sigma_{i+1}(\bar{\mathbf{W}})$  for all  $i \in [r' - 1]$ . Rewriting objective of Problem 18 in terms of these decompositions gives:

$$\begin{aligned} \|\bar{\mathbf{M}} - \bar{\mathbf{W}}\|_F^2 + \nu_{\{d_i\}} \|\bar{\mathbf{W}}\|_*^2 &= \|U_{\bar{\mathbf{M}}} \Sigma_{\bar{\mathbf{M}}} V_{\bar{\mathbf{M}}}^\top - U_{\bar{\mathbf{W}}} \Sigma_{\bar{\mathbf{W}}} V_{\bar{\mathbf{W}}}^\top\|_F^2 + \nu_{\{d_i\}} \|U_{\bar{\mathbf{W}}} \Sigma_{\bar{\mathbf{W}}} V_{\bar{\mathbf{W}}}^\top\|_*^2 \\ &= \|\Sigma_{\bar{\mathbf{M}}} - U' \Sigma_{\bar{\mathbf{W}}} V'^\top\|_F^2 + \nu_{\{d_i\}} \|\Sigma_{\bar{\mathbf{W}}}\|_*^2 \\ &= \|\Sigma_{\bar{\mathbf{M}}}\|_F^2 + \|\Sigma_{\bar{\mathbf{W}}}\|_F^2 - 2\langle \Sigma_{\bar{\mathbf{M}}}, U' \Sigma_{\bar{\mathbf{W}}} V'^\top \rangle + \nu_{\{d_i\}} \|\Sigma_{\bar{\mathbf{W}}}\|_*^2 \end{aligned}$$

where  $U' = U_{\bar{\mathbf{M}}}^\top U_{\bar{\mathbf{W}}}$  and  $V' = V_{\bar{\mathbf{M}}}^\top V_{\bar{\mathbf{W}}}$ . By Von Neumann's trace inequality, for a fixed  $\Sigma_{\bar{\mathbf{W}}}$  we have that

$$\langle \Sigma_{\bar{\mathbf{M}}}, U' \Sigma_{\bar{\mathbf{W}}} V'^\top \rangle \leq \sum_{i=1}^{r'} \sigma_i(\bar{\mathbf{M}}) \sigma_i(\bar{\mathbf{W}}),$$

where the equality is achieved when  $U_{\bar{\mathbf{M}}} = U_{\bar{\mathbf{W}}}$  and  $V_{\bar{\mathbf{M}}} = V_{\bar{\mathbf{W}}}$ . Hence, problem 18 is reduced to

$$\min_{\substack{\|\Sigma_{\bar{\mathbf{W}}}\|_0 \leq r, \\ \Sigma_{\bar{\mathbf{W}}} \geq 0}} \|\Sigma_{\bar{\mathbf{M}}} - \Sigma_{\bar{\mathbf{W}}}\|_F^2 + \nu_{\{d_i\}} (\text{Trace}(\Sigma_{\bar{\mathbf{W}}}))^2 = \min_{\bar{\sigma} \in \mathbb{R}_+^r} \sum_{i=1}^r (\lambda_i(\bar{\mathbf{M}}) - \bar{\sigma}_i)^2 + \nu_{\{d_i\}} \left( \sum_{i=1}^r \bar{\sigma}_i \right)^2$$

The Lagrangian is given by

$$L(\bar{\sigma}, \alpha) = \sum_{i=1}^r (\lambda_i(\bar{\mathbf{M}}) - \bar{\sigma}_i)^2 + \nu_{\{d_i\}} \left( \sum_{i=1}^r \bar{\sigma}_i \right)^2 - \sum_{i=1}^r \alpha_i \bar{\sigma}_i$$

The KKT conditions ensures that at the optima it holds for all  $i \in [r]$  that

$$\bar{\sigma}_i \geq 0, \alpha_i \geq 0, \bar{\sigma}_i \alpha_i = 0, \quad 2(\bar{\sigma}_i - \lambda_i(\bar{\mathbf{M}})) + 2\nu_{\{d_i\}} \sum_{j=1}^r \bar{\sigma}_j - \alpha_i = 0$$

Let  $\rho = |\{i : \bar{\sigma}_i > 0\}| \leq r$  be the number of nonzero  $\bar{\sigma}_i$ , i.e. rank of the global optimum  $\bar{\mathbf{W}}$ . For  $i \in [\rho]$ , we have  $\alpha_i = 0$ . Therefore, we have that:

$$\begin{aligned} \bar{\sigma}_i + \nu_{\{d_i\}} \sum_{j=1}^r \bar{\sigma}_j = \lambda_i(\bar{\mathbf{M}}) &\implies (\mathbf{I}_\rho + \nu_{\{d_i\}} \mathbf{1}_\rho \mathbf{1}_\rho^\top) \bar{\sigma}_{1:\rho} = \sigma_{1:\rho}(\bar{\mathbf{M}}) \\ &\implies \bar{\sigma}_{1:\rho} = \left( \mathbf{I}_\rho - \frac{\nu_{\{d_i\}}}{1 + \rho \nu_{\{d_i\}}} \mathbf{1}_\rho \mathbf{1}_\rho^\top \right) \sigma_{1:\rho}(\bar{\mathbf{M}}) = \sigma_{1:\rho}(\bar{\mathbf{M}}) - \frac{\nu_{\{d_i\}} \rho \kappa_\rho}{1 + \rho \nu_{\{d_i\}}} \mathbf{1}_\rho \end{aligned}$$

where  $\kappa_j := \frac{1}{j} \sum_{i=1}^j \sigma_i(\bar{\mathbf{M}})$ . The equation above tell us that for  $i \in [\rho]$ , the singular values of  $\bar{\mathbf{W}}$  are just a shrinkage of the singular values of  $\bar{\mathbf{M}}$ . In particular, it means that  $\rho \leq \text{rank}(\bar{\mathbf{M}})$ . Therefore, without loss of generality, we assume that  $r \leq \text{rank}(\bar{\mathbf{M}})$ . Also, since  $\bar{\sigma}_i > 0$  for all  $i \in [\rho]$ , it holds that  $\sigma_i(\bar{\mathbf{M}}) > \frac{\nu_{\{d_i\}} \rho \kappa_\rho}{1 + \rho \nu_{\{d_i\}}}$  for all  $i \in [\rho]$ . For  $i \in \{\rho + 1, \dots, r\}$ , on the other hand,  $\bar{\sigma}_i = 0$  and we have

$$\frac{1}{2} \alpha_i = \bar{\sigma}_i + \nu_{\{d_i\}} \sum_{j=1}^r \bar{\sigma}_j - \sigma_i(\bar{\mathbf{M}}) = 0 + \frac{\nu_{\{d_i\}}}{1 + \rho \nu_{\{d_i\}}} \sum_{j=1}^{\rho} \sigma_j(\bar{\mathbf{M}}) - \sigma_i(\bar{\mathbf{M}}) = -\sigma_i(\bar{\mathbf{M}}) + \frac{\nu_{\{d_i\}} \rho \kappa_\rho}{1 + \rho \nu_{\{d_i\}}},$$

where we used the fact that

$$\sum_{i=1}^r \bar{\sigma}_i = \mathbf{1}_\rho^\top \bar{\sigma}_{1:\rho} = \sum_{i=1}^{\rho} \sigma_i(\bar{\mathbf{M}}) - \frac{\nu_{\{d_i\}} \rho^2 \kappa_\rho}{1 + \rho \nu_{\{d_i\}}} = \left(1 - \frac{\nu_{\{d_i\}} \rho}{1 + \rho \nu_{\{d_i\}}}\right) \kappa_\rho = \frac{\rho \kappa_\rho}{1 + \rho \nu_{\{d_i\}}}.$$

By dual feasibility, we conclude that  $\sigma_i(\bar{\mathbf{M}}) \leq \frac{\nu_{\{d_i\}} \rho \kappa_\rho}{1 + \rho \nu_{\{d_i\}}}$  for all  $i \in \{\rho + 1, \dots, r\}$ , which completes the proof.  $\square$

*Proof of Theorem 2.7.* Consider  $\mathbf{W}^*$ , a global optimum of problem 3. If all such global optima can be implemented by equalized networks, then by Theorem 2.6 it holds that  $\Theta(\mathbf{W}^*) = \Theta^{**}(\mathbf{W}^*) = \nu_{\{d_i\}} \|\mathbf{W}^* \mathbf{C}^{\frac{1}{2}}\|_*^2$ . In this case, the lifted problem in Equation 3 boils down to the following convex problem

$$\min_{\mathbf{W} \in \mathbb{R}^{d_{k+1} \times d_0}} \mathbb{E}[\|\mathbf{y} - \mathbf{W}\mathbf{x}\|^2] + \nu_{\{d_i\}} \|\mathbf{W}\mathbf{C}^{\frac{1}{2}}\|_*^2, \quad \text{rank}(\mathbf{W}) \leq \min_{i \in [k+1]} d_i =: r. \quad (19)$$

Proposition 3.3, on the other hand, states that any rank-1 network map can be implemented by an equalized network. Therefore, the key idea of the proof is to make sure that the global optimum of problem 19 has rank equal to one. It suffices to notice that under the assumption  $\sigma_1(\bar{\mathbf{M}}) - \sigma_2(\bar{\mathbf{M}}) \geq \frac{1}{\nu_{\{d_i\}}} \sigma_2(\bar{\mathbf{M}})$ , it holds that  $\sigma_1(\bar{\mathbf{M}}) > \frac{\nu_{\{d_i\}} \sigma_1(\bar{\mathbf{M}})}{1 + \nu_{\{d_i\}}}$  and  $\sigma_j(\bar{\mathbf{M}}) \leq \frac{\nu_{\{d_i\}} \sigma_1(\bar{\mathbf{M}})}{1 + \nu_{\{d_i\}}}$  for all  $j > 1$ . In this case, using Lemma A.4, the solution  $\mathcal{S}_{\alpha_1}(\bar{\mathbf{M}}) \mathbf{C}^{-\frac{1}{2}}$  has rank equal to one, which completes the proof.  $\square$