

A PROOF OF THEOREMS AND TECHNICAL LEMMAS

A.1 PROOF OF LEMMA 2.1

Recall the shorthand $\bar{y}_a = (\sum_{i \in B_a} y_i)/k$, for $a \in [m]$. We have,

$$\begin{aligned} \mathcal{L}_{\text{ins}}(\boldsymbol{\theta}) &= \frac{1}{mk} \sum_{a=1}^m \sum_{i \in B_a} (\bar{y}_a - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2 \\ &= \frac{1}{mk} \sum_{a=1}^m \sum_{i \in B_a} \left(\bar{y}_a - \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) + \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right)^2 \\ &= \frac{1}{mk} \sum_{a=1}^m \sum_{i \in B_a} \left(\bar{y}_a - \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) \right)^2 + \frac{1}{mk} \sum_{a=1}^m \sum_{i \in B_a} \left(\frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right)^2 \\ &\quad + \frac{2}{mk} \sum_{a=1}^m \sum_{i \in B_a} \left(\bar{y}_a - \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) \right) \left(\frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right) \end{aligned}$$

Note that the first term can be written as

$$\frac{1}{mk} \sum_{a=1}^m \sum_{i \in B_a} \left(\bar{y}_a - \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) \right)^2 = \frac{1}{k} \sum_{a=1}^m \left(\bar{y}_a - \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) \right)^2 = \mathcal{L}_{\text{bag}}(\boldsymbol{\theta}).$$

For the second term, we have

$$\sum_{i \in B_a} \left(\frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right)^2 = \frac{1}{k} \sum_{i,j \in B_a} (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - f_{\boldsymbol{\theta}}(\mathbf{x}_j))^2 = \mathcal{R}(\boldsymbol{\theta}),$$

where we used the following identity for a_1, \dots, a_k and $\bar{a} = (\sum_{i=1}^k a_i)/k$:

$$\sum_{i=1}^k (a_i - \bar{a})^2 = \frac{1}{k} \sum_{i,j=1}^k (a_i - a_j)^2. \quad (11)$$

Finally, the third term works out at zero because

$$\begin{aligned} &\sum_{a=1}^m \sum_{i \in B_a} \left(\bar{y}_a - \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) \right) \left(\frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right) \\ &= \sum_{a=1}^m \left(\bar{y}_a - \frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) \right) \left[\sum_{i \in B_a} \left(\frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right) \right] = 0, \end{aligned}$$

since

$$\sum_{i \in B_a} \left(\frac{1}{k} \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right) = \sum_{j \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_j) - \sum_{i \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_i) = 0.$$

Combining the three terms together we arrive at $\mathcal{L}_{\text{ins}}(\boldsymbol{\theta}) = \mathcal{L}_{\text{bag}}(\boldsymbol{\theta}) + \mathcal{R}(\boldsymbol{\theta})$.

A.2 PROOF OF LEMMA 2.2

We use the shorthand $\bar{f}_a = (\sum_{i \in B_a} f_{\boldsymbol{\theta}}(\mathbf{x}_i))/k$, for $a \in [m]$. By Taylor's expansion of the loss ℓ on its second argument we have

$$\ell(\bar{y}_a, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) = \ell(\bar{y}_a, \bar{f}_a) + \frac{\partial}{\partial b} \ell(\bar{y}_a, \bar{f}_a) (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - \bar{f}_a) + \frac{1}{2} \frac{\partial^2}{\partial b^2} \ell(\bar{y}_a, \bar{f}_a) (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - \bar{f}_a)^2,$$

for some f between \bar{f}_a and $f_{\boldsymbol{\theta}}(\mathbf{x}_i)$ and $\partial/\partial b$, $\partial^2/\partial b^2$ indicate the first and second derivative of $\ell(a, b)$ with respect to the second input b .

Summing both sides of the above equation over $i \in B_a$, the second term works out at zero since $\sum_{i \in B_a} (f_{\theta}(\mathbf{x}_i) - \bar{f}_a) = 0$. Using the bound on the second derivative we arrive at

$$\sum_{i \in B_a} \ell(\bar{y}_a, f_{\theta}(\mathbf{x}_i)) \leq k\ell(\bar{y}_a, \bar{f}_a) + \sum_{i \in B_a} C(f_{\theta}(\mathbf{x}_i) - \bar{f}_a)^2.$$

Next, summing both sides of the above equation over bags $a \in [m]$, and dividing by mk , we get

$$\mathcal{L}_{\text{ins}}(\theta) \leq \mathcal{L}_{\text{bag}}(\theta) + \frac{1}{k} \sum_{a=1}^m \sum_{i \in B_a} C(f_{\theta}(\mathbf{x}_i) - \bar{f}_a)^2.$$

By invoking identity (11) in the above we arrive at (5).

We are now ready to prove the second part of the statement. If the loss $\ell(\cdot, \cdot)$ is convex in the second input, by the Jensen's inequality we have

$$\frac{1}{k} \sum_{i \in B_a} \ell(\bar{y}_a, f_{\theta}(\mathbf{x}_i)) \geq \ell\left(\bar{y}_a, \frac{1}{k} \sum_{i \in B_a} f_{\theta}(\mathbf{x}_i)\right)$$

Taking the average of both side over the bags $a \in [m]$, we obtain that $\mathcal{L}_{\text{ins}}(\theta) \geq \mathcal{L}_{\text{bag}}(\theta)$, which completes the proof of lemma.

B PROOF OF THEOREM 2.5

Recall m as the number of bags, and n as the number of samples. Since the bags are non-overlapping and each of size k , we have $m = n/k$. Define $\mathbf{S} \in \mathbb{R}^{m \times n}$, as a matrix the encodes the bagging structure, with $S_{ia} = 1/\sqrt{k} \mathbf{1}_{\{j \in B_a\}}$ where B_a indicates a -th bag, for $a \in [m]$.

We next write the bag-level loss function and the instance-level loss function in terms of \mathbf{S} as follows:

$$\begin{aligned} \mathcal{L}_{\text{bag}}(\theta) &= \frac{1}{km} \|\mathbf{S}(\mathbf{y} - \mathbf{X}\theta)\|_2^2, \\ \mathcal{L}_{\text{ins}}(\theta) &= \frac{1}{km} \|\mathbf{S}^{\top} \mathbf{S} \mathbf{y} - \mathbf{X}\theta\|_2^2. \end{aligned}$$

The interpolating loss function (7) then reads as

$$\mathcal{L}_{\text{int}}(\theta) = \frac{1}{mk} \left((1 - \rho) \|\mathbf{S} \mathbf{X} \theta - \mathbf{S} \mathbf{y}\|_2^2 + \rho \|\mathbf{X} \theta - \mathbf{S}^{\top} \mathbf{S} \mathbf{y}\|_2^2 \right).$$

This can equivalently be written as

$$\mathcal{L}_{\text{int}}(\theta) = \frac{1}{mk} \left\| \left(\frac{\sqrt{\rho} \mathbf{I}}{\sqrt{1 - \rho}} \right) \mathbf{X} \theta - \left(\frac{\sqrt{\rho} \mathbf{S}^{\top} \mathbf{S} \mathbf{y}}{\sqrt{1 - \rho}} \right) \right\|_2^2.$$

The minimizer of the above loss admits a closed-form solution given by $\hat{\theta}_{\text{int}} = \mathbf{B} \mathbf{y}$, with

$$\mathbf{B} = (\mathbf{X}^{\top} (\rho \mathbf{I} + (1 - \rho) \mathbf{S}^{\top} \mathbf{S}) \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{S}^{\top} \mathbf{S}.$$

We define the shorthand $\mathbf{E} := \rho \mathbf{I} + (1 - \rho) \mathbf{S}^{\top} \mathbf{S} \in \mathbb{R}^{n \times n}$, which is non-singular for $\rho > 0$, and $\mathbf{M} = (\mathbf{X}^{\top} \mathbf{E} \mathbf{X})^{-1} \mathbf{X}^{\top} \in \mathbb{R}^{d \times n}$. We then have $\mathbf{B} = \mathbf{M} \mathbf{S}^{\top} \mathbf{S}$.

We next recall the bias-variance decomposition (6), where the bias and variance are given by

$$\begin{aligned} \text{Bias}(\hat{\theta}_{\text{int}}) &= \|(\mathbf{B} \mathbf{X} - \mathbf{I}) \theta_0\|_2^2 \\ &= \|(\mathbf{M} \mathbf{S}^{\top} \mathbf{S} \mathbf{X} - \mathbf{M} \mathbf{E} \mathbf{X}) \theta_0\|_2^2 \\ &= \|\mathbf{M} (\mathbf{S}^{\top} \mathbf{S} - \mathbf{E}) \mathbf{X} \theta_0\|_2^2, \end{aligned} \tag{12}$$

$$\text{Var}(\hat{\theta}_{\text{int}}) = \sigma^2 \|\mathbf{M} \mathbf{S}^{\top} \mathbf{S}\|_F^2, \tag{13}$$

with $\|\cdot\|_F$ indicating the matrix Frobenius norm.

We continue by treating the bias and the variance separately.

B.1 CALCULATING THE BIAS

Since the distribution of the features matrix \mathbf{X} is invariant under rotation, we can assume that $\boldsymbol{\theta}_0 = \|\boldsymbol{\theta}_0\| \mathbf{e}_i$, where $\mathbf{e}_i \in \mathbb{R}^d$ is the vector with one at i -th entry and zero everywhere else. By taking average on $i \in [d]$ we obtain

$$\begin{aligned} \text{Bias}(\widehat{\boldsymbol{\theta}}_{\text{int}}) &\stackrel{(d)}{=} \frac{\|\boldsymbol{\theta}_0\|_2^2}{d} \sum_{i \in [d]} \|\mathbf{M}(\mathbf{S}^\top \mathbf{S} - \mathbf{E})\mathbf{X}\mathbf{e}_i\|_2^2 \\ &= \frac{\|\boldsymbol{\theta}_0\|_2^2}{d} \text{tr} \left(\mathbf{M}(\mathbf{S}^\top \mathbf{S} - \mathbf{E})\mathbf{X} \left(\sum_{i \in [p]} \mathbf{e}_i \mathbf{e}_i^\top \right) \mathbf{X}^\top (\mathbf{S}^\top \mathbf{S} - \mathbf{E})\mathbf{M}^\top \right) \\ &= \frac{\|\boldsymbol{\theta}_0\|_2^2}{d} \|\mathbf{M}(\mathbf{S}^\top \mathbf{S} - \mathbf{E})\mathbf{X}\|_F^2. \end{aligned}$$

Let us define $\boldsymbol{\Lambda} \in \mathbb{R}^{n \times n}$ as follows:

$$\begin{aligned} \boldsymbol{\Lambda} &:= -(\mathbf{S}^\top \mathbf{S} - \mathbf{E}) \\ &= -(\mathbf{S}^\top \mathbf{S} - (\rho \mathbf{I} + (1 - \rho)\mathbf{S}^\top \mathbf{S})) \\ &= \rho(\mathbf{I} - \mathbf{S}^\top \mathbf{S}). \end{aligned} \tag{14}$$

The bias can then be written in terms of $\boldsymbol{\Lambda}$ as $\text{Bias}(\boldsymbol{\theta}) = \frac{1}{d} \|\mathbf{M}\boldsymbol{\Lambda}\mathbf{X}\|_F^2$. In our next lemma, we characterize the asymptotic behavior of the bias.

Lemma B.1. *Under the asymptotic regime of Assumption 2.3 we have*

$$\frac{1}{d} \|\mathbf{M}\boldsymbol{\Lambda}\mathbf{X}\|_F^2 \stackrel{(p)}{\rightarrow} \alpha_*^2 + \frac{\alpha_*^2}{\frac{(k-1)\psi}{k^2(1-\alpha_*)^2} - \left(\frac{\alpha_*}{1-\alpha_*}\right)^2 \frac{1}{k} - \frac{k-1}{k}},$$

where α_* is the nonnegative fixed point of the following equation:

$$\rho + \frac{\psi}{k(1-\alpha_*)} - 1 = \frac{\psi}{k\alpha_*} \rho(k-1).$$

Since $\|\boldsymbol{\theta}_0\| \rightarrow 1$, the result (8) follows from Lemma B.1.

We refer to the supplementary D.1 for the proof of Lemma B.1.

B.2 CALCULATING THE VARIANCE

Since the bags are non-overlapping we have $\mathbf{S}\mathbf{S}^\top = \mathbf{I}_m$. Therefore $\mathbf{S}^\top \mathbf{S}$ is a projection matrix and can be written as $\mathbf{S}^\top \mathbf{S} = \mathbf{U}\mathbf{U}^\top$, with $\mathbf{U} \in \mathbb{R}^{n \times m}$ an orthogonal matrix. Recall that the variance is given by $\text{Var}(\widehat{\boldsymbol{\theta}}_{\text{int}}) = \sigma^2 \|\mathbf{M}\mathbf{S}^\top \mathbf{S}\|_F^2$. We use the next lemma to characterize the asymptotic behavior of the variance.

Lemma B.2. *Under the asymptotic regime of Assumption 2.3 for any vector $\mathbf{a} \in \mathbb{R}^m$ we have*

$$\frac{n}{\|\mathbf{a}\|^2} \|\mathbf{M}\mathbf{U}\mathbf{a}\|_2^2 \stackrel{(p)}{\rightarrow} \frac{k}{v_*},$$

where v_* is given as the fixed point of the following system of equations in (v, u) :

$$\begin{cases} \frac{\psi}{1+u} + \frac{\rho\psi(k-1)}{\rho+u} = k, \\ \frac{\psi(1+v)}{(1+u)^2} + \frac{\rho^2\psi(k-1)}{(\rho+u)^2} = k. \end{cases}$$

Proof of Lemma B.2 is given in the supplementary D.2.

We next use the above lemma for each row of \mathbf{U} separately (as the vector \mathbf{a}) and add them together.

Using the fact that $\|\mathbf{U}\|_F^2 = m$ and $m/n = k$, we get that $\|\mathbf{M}\mathbf{U}\mathbf{U}^\top\|_F^2 \stackrel{(p)}{\rightarrow} 1/v_*$, which completes the variance calculation.

B.3 PROOF OF LEMMA 4.2

We use the idea of (Dwork et al., 2014, Theorem 3.6) to prove this lemma. Given a database $D = (y_1, y_2, \dots, y_n)$, Algorithm 1 (we denote this mapping by $\mathcal{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$) outputs m real numbers $(\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_m)$. Given the database D , we define the map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ by $D \mapsto (\bar{y}_1, \bar{y}_2, \dots, \bar{y}_m)$, which computes the mean of labels in each bag. Fix any pair of neighboring databases D, D' that differ in the label of a single example. We have $\|f(D) - f(D')\|_1 := \sum_{a \in [m]} |f(D)_a - f(D')_a| \leq \Delta f := \frac{C\sqrt{\log n}}{k}$. In this argument, we used Assumption 2.4 that assumes non-overlapping bags, and therefore, changing a certain y_i in D leads to a change in only one of \bar{y}_a by at most Δf . Let $p_{\mathcal{A}(D)}(z)$ and $p_{\mathcal{A}(D')}(z)$ denote the probability density function of $\mathcal{A}(D)$ and $\mathcal{A}(D')$. We have

$$\begin{aligned} \frac{p_{\mathcal{A}(D)}(z)}{p_{\mathcal{A}(D')}(z)} &= \prod_{a \in [m]} \frac{\exp\left(-\frac{\varepsilon |f(D)_a - z_a|}{\Delta f}\right)}{\exp\left(-\frac{\varepsilon |f(D')_a - z_a|}{\Delta f}\right)} \\ &= \prod_{a \in [m]} \exp\left(\frac{\varepsilon (|f(D')_a - z_a| - |f(D)_a - z_a|)}{\Delta f}\right) \\ &\leq \prod_{a \in [m]} \exp\left(\frac{\varepsilon |f(D')_a - f(D)_a|}{\Delta f}\right) \\ &= \exp\left(\frac{\varepsilon \|f(D) - f(D')\|_1}{\Delta f}\right) \\ &\leq e^\varepsilon, \end{aligned}$$

which completes the proof.

C PROOF OF THEOREM 4.3

Recall that in Algorithm 1, the individual responses are first truncated by $C\sqrt{\log n}$ and then after the aggregate responses are computed, a Laplace noise is added to them to ensure label DP. We define \mathcal{E} is the event that no truncation happens, namely:

$$\mathcal{E} := \mathbf{1}_{\{|y_i| \leq C\sqrt{\log n}, \forall i \in [n]\}}. \quad (15)$$

Since $y_i \sim N(0, \|\theta_0\|^2 + \sigma^2)$, $\|\theta_0\| = 1$, by using Gaussian tail bound along with union bounding we arrive at

$$\mathbb{P}(\mathcal{E}) \geq 1 - n \exp\left(-\frac{C^2}{2(1+\sigma^2)} \log n\right) = 1 - n^{-c}, \quad (16)$$

with $c = \frac{C^2}{2(1+\sigma^2)} - 1 > 0$.

We next bound the risk of estimator $\hat{\theta}_{\text{int}}$ as follows:

$$\text{Risk}(\hat{\theta}_{\text{int}}) = \mathbb{E}[\|\hat{\theta}_{\text{int}} - \theta_0\|^2 \mathbf{1}_{\{\mathcal{E}\}} | \mathbf{X}] + \mathbb{E}[\|\hat{\theta}_{\text{int}} - \theta_0\|^2 \mathbf{1}_{\{\mathcal{E}^c\}} | \mathbf{X}]. \quad (17)$$

For the first term, note that on the instance \mathcal{E} (no truncation), the privatized aggregate responses are just the aggregate responses with an additive zero mean noise with variance $2C^2 \log n / (k\varepsilon)^2$. So we can use the analysis in the proof of Theorem 2.5 with the inflated noise variance. Let $\hat{\theta}_{\text{int}}^{\text{nt}}$ be the estimator using untruncated responses in Algorithm 1. We then have This gives us

$$\begin{aligned} \frac{1}{\log n} \mathbb{E}[\|\hat{\theta}_{\text{int}} - \theta_0\|^2 \mathbf{1}_{\{\mathcal{E}\}} | \mathbf{X}] &= \frac{1}{\log n} \mathbb{E}[\|\hat{\theta}_{\text{int}}^{\text{nt}} - \theta_0\|^2 \mathbf{1}_{\{\mathcal{E}\}} | \mathbf{X}] \\ &= \frac{1}{\log n} \mathbb{E}[\|\hat{\theta}_{\text{int}}^{\text{nt}} - \theta_0\|^2 | \mathbf{X}] - \frac{1}{\log n} \mathbb{E}[\|\hat{\theta}_{\text{int}}^{\text{nt}} - \theta_0\|^2 \mathbf{1}_{\{\mathcal{E}^c\}} | \mathbf{X}] \\ &= \frac{1}{\log n} \text{Bias}(\hat{\theta}_{\text{int}}^{\text{nt}}) + \frac{1}{\log n} \text{Var}(\hat{\theta}_{\text{int}}^{\text{nt}}) - \frac{1}{\log n} \mathbb{E}[\|\hat{\theta}_{\text{int}}^{\text{nt}} - \theta_0\|^2 \mathbf{1}_{\{\mathcal{E}^c\}} | \mathbf{X}], \end{aligned} \quad (18)$$

where the bias is given by (8) and variance is given by (9), where σ^2/k is replaced with the inflated variance $\sigma^2/k + 2C^2 \log n/(k\varepsilon)^2$. Since $\text{Bias}(\hat{\boldsymbol{\theta}}_{\text{int}}^{\text{nt}})$ has a finite limit, the first term above vanishes as $n \rightarrow \infty$. For the second term we have

$$\frac{1}{\log n} \text{Var}(\hat{\boldsymbol{\theta}}_{\text{int}}^{\text{nt}}) \xrightarrow{(p)} \frac{2C^2}{k\varepsilon^2} \frac{1}{v_*}.$$

since $\sigma^2/\log n \rightarrow 0$. For the third term, by Cauchy–Schwarz inequality we have

$$\mathbb{E}[\|\hat{\boldsymbol{\theta}}_{\text{int}}^{\text{nt}} - \boldsymbol{\theta}_0\|^2 \mathbf{1}_{\{\mathcal{E}^c\}} | \mathbf{X}] \leq \mathbb{E}[\|\hat{\boldsymbol{\theta}}_{\text{int}}^{\text{nt}} - \boldsymbol{\theta}_0\|^4 | \mathbf{X}]^{1/2} \mathbb{P}(\mathcal{E}^c). \quad (19)$$

Using the high probability bound on the minimum singular value of the Gaussian matrix \mathbf{X} (Ver-shynin, 2018, Theorem 4.6.1), we can show that $\mathbb{E}[\|\hat{\boldsymbol{\theta}}_{\text{int}}^{\text{nt}} - \boldsymbol{\theta}_0\|^4 | \mathbf{X}]$ is bounded in probability and since $\mathbb{P}(\mathcal{E}^c) \leq n^{-c}$, we conclude that the third term in (18) also vanishes as $n \rightarrow \infty$, in probability. Combining these together we arrive at

$$\frac{1}{\log n} \mathbb{E}[\|\hat{\boldsymbol{\theta}}_{\text{int}} - \boldsymbol{\theta}_0\|^2 \mathbf{1}_{\{\mathcal{E}\}} | \mathbf{X}] \xrightarrow{(p)} \frac{2C^2}{k\varepsilon^2} \frac{1}{v_*}. \quad (20)$$

Similar to (19) we can also show that

$$\frac{1}{\log n} \mathbb{E}[\|\hat{\boldsymbol{\theta}}_{\text{int}} - \boldsymbol{\theta}_0\|^2 \mathbf{1}_{\{\mathcal{E}^c\}} | \mathbf{X}] \xrightarrow{(p)} 0,$$

which along with (20) and (17) implies that

$$\frac{1}{\log n} \text{Risk}(\hat{\boldsymbol{\theta}}_{\text{int}}) \xrightarrow{(p)} \frac{2C^2}{k\varepsilon^2} \frac{1}{v_*},$$

completing the proof.

D PROOF OF INTERMEDIATE LEMMAS

D.1 PROOF OF LEMMA B.1

Write $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_d]$ with \mathbf{x}_i representing the i -th column. We then have $\|\mathbf{M}\boldsymbol{\Lambda}\mathbf{X}\|^2 = \sum_{i=1}^d \|\mathbf{M}\boldsymbol{\Lambda}\mathbf{x}_i\|^2$. We compute the asymptotic behavior of each of the summand separately. Indeed, by symmetry of the distributions of \mathbf{x}_i , we will see that all summands converge to the same limit.

Recall that $\mathbf{M} = (\mathbf{X}^\top \mathbf{E} \mathbf{X})^{-1} \mathbf{X}^\top$. Consider the following optimization problem:

$$\boldsymbol{\alpha}_i = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \frac{1}{d} \|\mathbf{E}^{1/2} \mathbf{X} \boldsymbol{\alpha} - \mathbf{E}^{-1/2} \boldsymbol{\Lambda} \mathbf{x}_i\|_2^2. \quad (21)$$

It is easy to see that by the KKT condition, $\boldsymbol{\alpha}_i = (\mathbf{X}^\top \mathbf{E} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Lambda} \mathbf{x}_i = \mathbf{M} \boldsymbol{\Lambda} \mathbf{x}_i$. Therefore, we are interested in characterizing $\|\boldsymbol{\alpha}_i\|$ in the asymptotic regime, described in Assumption 2.3.

We write $\boldsymbol{\alpha}$ as $(\alpha_i, \boldsymbol{\alpha}_{\sim i})$ to separate its i -th entry from the rest. Likewise we write $\mathbf{X} = [\mathbf{x}_i \mathbf{X}_{\sim i}]$ to separate the i -th columns from the rest. We then have

$$\begin{aligned} & \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \frac{1}{d} \|\mathbf{E}^{1/2} \mathbf{X} \boldsymbol{\alpha} - \mathbf{E}^{-1/2} \boldsymbol{\Lambda} \mathbf{x}_i\|_2^2 \\ &= \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \frac{1}{d} \|\mathbf{E}^{1/2} \mathbf{x}_i \alpha_i + \mathbf{E}^{1/2} \mathbf{X}_{\sim i} \boldsymbol{\alpha}_{\sim i} - \mathbf{E}^{-1/2} \boldsymbol{\Lambda} \mathbf{x}_i\|_2^2 \\ &= \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \frac{1}{d} \|\mathbf{E}^{1/2} \mathbf{X}_{\sim i} \boldsymbol{\alpha}_{\sim i} + (\alpha_i \mathbf{E}^{1/2} - \mathbf{E}^{-1/2} \boldsymbol{\Lambda}) \mathbf{x}_i\|_2^2 \\ &= \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \max_{\mathbf{v} \in \mathbb{R}^n} \frac{2}{d} \left(\mathbf{v}^\top (\alpha_i \mathbf{E}^{1/2} - \mathbf{E}^{-1/2} \boldsymbol{\Lambda}) \mathbf{x}_i + \mathbf{v}^\top \mathbf{E}^{1/2} \mathbf{X}_{\sim i} \boldsymbol{\alpha}_{\sim i} - \frac{1}{2} \|\mathbf{v}\|_2^2 \right), \end{aligned} \quad (22)$$

where in the last step we used the identity $\max_{\mathbf{v}} (\mathbf{v}^\top \mathbf{x} - \|\mathbf{v}\|^2/2) = \|\mathbf{x}\|^2/2$ for any vector \mathbf{x} .

We next note that $\mathbf{S}\mathbf{S}^\top = \mathbf{I}$ since the bags are non-overlapping. Therefore we can write $\mathbf{S}^\top \mathbf{S} = \mathbf{U}\mathbf{U}^\top$ for an orthogonal matrix $\mathbf{U} \in \mathbb{R}^{n \times m}$. We then have

$$\mathbf{E} := \rho \mathbf{I} + (1 - \rho) \mathbf{S}^\top \mathbf{S} = \mathbf{U}\mathbf{U}^\top + \rho \mathbf{U}_\perp \mathbf{U}_\perp^\top, \quad \mathbf{\Lambda} = \rho(\mathbf{I} - \mathbf{S}^\top \mathbf{S}) = \rho \mathbf{U}_\perp \mathbf{U}_\perp^\top.$$

where \mathbf{U}_\perp is an orthogonal matrix representing the orthogonal space to the column space of \mathbf{U} . We next decompose the vector \mathbf{v} in the above optimization as $\mathbf{v} = \mathbf{U}\mathbf{v}_1 + \mathbf{U}_\perp \mathbf{v}_2$ and therefore $\|\mathbf{v}\|^2 = \|\mathbf{v}_1\|^2 + \|\mathbf{v}_2\|^2$.

We introduce the change of variable $\tilde{\mathbf{v}} = \mathbf{E}^{1/2} \mathbf{v}$ in optimization (22). Note that $\tilde{\mathbf{v}} = \mathbf{U}\mathbf{v}_1 + \sqrt{\rho} \mathbf{U}_\perp \mathbf{v}_2$. Continuing with (22) in terms of $\tilde{\mathbf{v}}$ we have

$$\min_{\alpha \in \mathbb{R}^d} \max_{\tilde{\mathbf{v}} \in \mathbb{R}^n} \frac{2}{d} \left(\tilde{\mathbf{v}}^\top (\alpha_i \mathbf{I} - \mathbf{E}^{-1} \mathbf{\Lambda}) \mathbf{x}_i + \tilde{\mathbf{v}}^\top \mathbf{X}_{\sim i} \alpha_{\sim i} - \frac{1}{2} \|\mathbf{E}^{-1/2} \tilde{\mathbf{v}}\|_2^2 \right). \quad (23)$$

To analyze the asymptotic behavior of the solution to the above minimax optimization, we use the Convex-Gaussian-Minimax-Theorem (CGMT) (Thrapoulidis et al., 2015, Theorem 3), which is a power extension of the classical Gordon's Gaussian min-max theorem (Gordon (1988)), under additional convexity assumptions. According to CGMT, the above optimization is equivalent to the following auxiliary optimization problem:

$$\min_{\alpha \in \mathbb{R}^d} \max_{\tilde{\mathbf{v}} \in \mathbb{R}^n} \frac{2}{d} \left(\tilde{\mathbf{v}}^\top (\alpha_i \mathbf{I} - \mathbf{E}^{-1} \mathbf{\Lambda}) \mathbf{x}_i + \|\alpha_{\sim i}\| \tilde{\mathbf{v}}^\top \mathbf{g} + \|\tilde{\mathbf{v}}\| \mathbf{h}^\top \alpha_{\sim i} - \frac{1}{2} \|\mathbf{E}^{-1/2} \tilde{\mathbf{v}}\|_2^2 \right), \quad (24)$$

with $\mathbf{g} \sim N(0, \mathbf{I}_n)$ and $\mathbf{h} \sim N(0, \mathbf{I}_{d-1})$ independent Gaussian vectors. We next write the above optimization in terms of the components \mathbf{v}_1 and \mathbf{v}_2 , noting that $\mathbf{E}^{-1} \mathbf{\Lambda} = \mathbf{U}_\perp \mathbf{U}_\perp^\top$, as follows:

$$\min_{\alpha \in \mathbb{R}^d} \max_{\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^n} \frac{2}{d} \left(\alpha_i \mathbf{v}_1^\top \mathbf{U}^\top \mathbf{x}_i + \sqrt{\rho} \mathbf{v}_2^\top \mathbf{U}_\perp^\top (\alpha_i \mathbf{I} - \mathbf{U}_\perp \mathbf{U}_\perp^\top) \mathbf{x}_i + \|\alpha_{\sim i}\| (\mathbf{v}_1^\top \mathbf{U}^\top \mathbf{g} + \sqrt{\rho} \mathbf{v}_2^\top \mathbf{U}_\perp^\top \mathbf{g}) + \sqrt{\|\mathbf{v}_1\|^2 + \rho \|\mathbf{v}_2\|^2} \mathbf{h}^\top \alpha_{\sim i} - \frac{1}{2} \|\mathbf{v}_1\|^2 - \frac{1}{2} \|\mathbf{v}_2\|^2 \right). \quad (25)$$

Define the shorthand

$$\begin{aligned} \mathbf{x}_1 &:= \mathbf{U}^\top \mathbf{x}_i \sim N(0, \mathbf{I}_m), \\ \mathbf{x}_2 &:= \mathbf{U}_\perp^\top \mathbf{x}_i \sim N(0, \mathbf{I}_{n-m}), \\ \mathbf{g}_1 &:= \mathbf{U}^\top \mathbf{g} \sim N(0, \mathbf{I}_m), \\ \mathbf{g}_2 &:= \mathbf{U}_\perp^\top \mathbf{g} \sim N(0, \mathbf{I}_{n-m}). \end{aligned}$$

Then optimization (25) can be rewritten as

$$\min_{\alpha \in \mathbb{R}^d} \max_{\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^n} \frac{2}{d} \left(\alpha_i \mathbf{v}_1^\top \mathbf{x}_1 + \sqrt{\rho} (\alpha_i - 1) \mathbf{v}_2^\top \mathbf{x}_2 + \|\alpha_{\sim i}\| (\mathbf{v}_1^\top \mathbf{g}_1 + \sqrt{\rho} \mathbf{v}_2^\top \mathbf{g}_2) + \sqrt{\|\mathbf{v}_1\|^2 + \rho \|\mathbf{v}_2\|^2} \mathbf{h}^\top \alpha_{\sim i} - \frac{1}{2} \|\mathbf{v}_1\|^2 - \frac{1}{2} \|\mathbf{v}_2\|^2 \right). \quad (26)$$

We fix $\|\mathbf{v}_1\| = \beta_1$ and $\|\mathbf{v}_2\| = \beta_2$ and first optimize over the directions of $\mathbf{v}_1, \mathbf{v}_2$ and then over the norms β_1 and β_2 . This brings us to

$$\min_{\alpha \in \mathbb{R}^d} \max_{\beta_1, \beta_2 \geq 0} \frac{2}{d} \left(\beta_1 \|\alpha_i \mathbf{x}_1 + \|\alpha_{\sim i}\| \mathbf{g}_1\| + \beta_2 \|\sqrt{\rho} (\alpha_i - 1) \mathbf{x}_2 + \|\alpha_{\sim i}\| \sqrt{\rho} \mathbf{g}_2\| + \sqrt{\beta_1^2 + \rho \beta_2^2} \mathbf{h}^\top \alpha_{\sim i} - \frac{1}{2} \beta_1^2 - \frac{1}{2} \beta_2^2 \right). \quad (27)$$

In order to optimize over $\alpha_{\sim i}$, we first fix its norm to $\eta := \|\alpha_{\sim i}\|$ and optimize over its direction, and then optimize over η , which results in:

$$\min_{\eta \geq 0, \alpha_i} \max_{\beta_1, \beta_2 \geq 0} \frac{2}{d} \left(\beta_1 \|\alpha_i \mathbf{x}_1 + \eta \mathbf{g}_1\| + \beta_2 \|\sqrt{\rho} (\alpha_i - 1) \mathbf{x}_2 + \eta \sqrt{\rho} \mathbf{g}_2\| + \eta \sqrt{\beta_1^2 + \rho \beta_2^2} \|\mathbf{h}\| - \frac{1}{2} \beta_1^2 - \frac{1}{2} \beta_2^2 \right). \quad (28)$$

The next step in the CGMT framework is to compute the pointwise limit of the objective functions. Using the concentration of Lipschitz functions of Gaussian vectors we have

$$\frac{1}{\sqrt{d}} \|\alpha_i \mathbf{x}_1 + \eta \mathbf{g}_1\| \xrightarrow{(p)} \sqrt{(\alpha_i^2 + \eta^2) \frac{\psi}{k}},$$

$$\frac{1}{\sqrt{d}} \|\sqrt{\rho}(\alpha_i - 1) \mathbf{x}_2 + \eta \sqrt{\rho} \mathbf{g}_2\| \xrightarrow{(p)} \sqrt{(\rho(\alpha_i - 1)^2 + \rho \eta^2) \psi \left(1 - \frac{1}{k}\right)},$$

where we used Assumption 2.3 by which $n/d \rightarrow \psi$ and $m = n/k$.

We also have $\frac{1}{\sqrt{d}} \|\mathbf{h}\| \xrightarrow{(p)} 1$.

We therefore arrive at the following deterministic optimization problem

$$\min_{\eta \geq 0, \alpha_i, \beta_1, \beta_2 \geq 0} \max \left(\beta_1 \sqrt{(\alpha_i^2 + \eta^2) \frac{\psi}{k}} + \beta_2 \sqrt{(\rho(\alpha_i - 1)^2 + \rho \eta^2) \psi \left(1 - \frac{1}{k}\right)} \right. \\ \left. + \sqrt{\beta_1^2 + \rho \beta_2^2} - \frac{1}{2} \beta_1^2 - \frac{1}{2} \beta_2^2 \right), \quad (29)$$

where we made the change of variables $2\beta_1/\sqrt{d} \rightarrow \beta_1$ and $2\beta_2/\sqrt{d} \rightarrow \beta_2$.

By writing the stationary conditions for the above optimization, and simplifying the resulting system of equations by solving for β_1, β_2 , and substituting for them in the other two equations, we arrive at the following two equations for α_i and η :

$$\begin{cases} \rho + \frac{\psi}{k(1-\alpha_*)} - 1 = \frac{\psi}{k\alpha_*} \rho(k-1) \\ \eta_*^2 + \frac{k\alpha_*^2}{k-1} + \frac{\eta_*^2 \alpha_*^2}{(1-\alpha_*)^2(k-1)} = \frac{\psi}{k} \frac{\eta_*^2}{(1-\alpha_*)^2}. \end{cases}$$

As the final step, recall that by definition $\eta := \|\alpha_{\sim i}\|$ and therefore, $\|\alpha_i\|^2 \xrightarrow{(p)} \alpha_*^2 + \eta_*^2$. As we see it is independent of the index i and therefore,

$$\frac{1}{d} \|\mathbf{M} \mathbf{\Lambda} \mathbf{X}\|^2 = \frac{1}{d} \sum_{i=1}^d \|\mathbf{M} \mathbf{\Lambda} \mathbf{x}_i\|^2 = \frac{1}{d} \sum_{i=1}^d \alpha_i^2 \xrightarrow{(p)} \alpha_*^2 + \eta_*^2.$$

This completes the proof.

D.2 PROOF OF LEMMA B.2

Recall that $\mathbf{M} = (\mathbf{X}^\top \mathbf{E} \mathbf{X})^{-1} \mathbf{X}^\top$. Consider the following optimization problem:

$$\boldsymbol{\alpha} = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \frac{1}{d} \|\mathbf{E}^{1/2} \mathbf{X} \boldsymbol{\alpha} - \mathbf{E}^{-1/2} \mathbf{U} \mathbf{a}\|_2^2. \quad (30)$$

The solution to the above optimization problem has a closed-form solution given by $\boldsymbol{\alpha} = (\mathbf{X}^\top \mathbf{E} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{U} \mathbf{a} = \mathbf{M} \mathbf{U} \mathbf{a}$. So we are interested in characterizing the norm of the optimal solution to the above optimization problem.

Similar to the proof of Lemma B.1, we use the framework of CGMT to characterize $\|\boldsymbol{\alpha}\|$ in the asymptotic regime described in Assumption 2.3.

Using the identity $\|\mathbf{x}\|/2 = \max_{\mathbf{v}} (\mathbf{v}^\top \mathbf{x} - \|\mathbf{v}\|^2/2)$, we rewrite the above optimization as:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \frac{1}{d} \|\mathbf{E}^{1/2} \mathbf{X} \boldsymbol{\alpha} - \mathbf{E}^{-1/2} \mathbf{U} \mathbf{a}\|_2^2 \\ = \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \max_{\mathbf{v} \in \mathbb{R}^n} \frac{2}{d} \left(\mathbf{v}^\top \mathbf{E}^{1/2} \mathbf{X} \boldsymbol{\alpha} - \mathbf{v}^\top \mathbf{E}^{-1/2} \mathbf{U} \mathbf{a} - \frac{1}{2} \|\mathbf{v}\|_2^2 \right), \quad (31)$$

By using Convex-Gaussian-Minimax-Theorem (Thrapoulidis et al., 2015, Theorem 3), the above optimization is equivalent to the following auxiliary optimization problem:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \max_{\mathbf{v} \in \mathbb{R}^n} \frac{2}{d} \left(\|\boldsymbol{\alpha}\| \mathbf{v}^\top \mathbf{E}^{1/2} \mathbf{g} + \|\mathbf{E}^{1/2} \mathbf{v}\| \mathbf{h}^\top \boldsymbol{\alpha} - \mathbf{v}^\top \mathbf{E}^{-1/2} \mathbf{U} \mathbf{a} - \frac{1}{2} \|\mathbf{v}\|_2^2 \right), \quad (32)$$

with $\mathbf{g} \sim N(0, \mathbf{I}_n)$ and $\mathbf{h} \sim N(0, \mathbf{I}_d)$ independent Gaussian vectors.

We also recall that $\mathbf{S}^\top \mathbf{S} = \mathbf{U}\mathbf{U}^\top$ and so

$$\mathbf{E} := \rho \mathbf{I} + (1 - \rho) \mathbf{S}^\top \mathbf{S} = \mathbf{U}\mathbf{U}^\top + \rho \mathbf{U}_\perp \mathbf{U}_\perp^\top,$$

with $\mathbf{U}_\perp \in \mathbb{R}^{n \times (n-m)}$ denotes the orthogonal matrix, whose column space is orthogonal to the column space of \mathbf{U} . We decompose \mathbf{v} to its component in the column space of \mathbf{U} and \mathbf{U}_\perp as

$$\mathbf{v} = \mathbf{U}\mathbf{v}_1 + \mathbf{U}_\perp \mathbf{v}_2, \quad \|\mathbf{v}\|^2 = \|\mathbf{v}_1\|^2 + \|\mathbf{v}_2\|^2.$$

Therefore, $\mathbf{E}^{1/2} \mathbf{v} = \mathbf{U}\mathbf{v}_1 + \sqrt{\rho} \mathbf{U}_\perp \mathbf{v}_2$ and so the above optimization (32) can be written as

$$\begin{aligned} \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \max_{\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^n} \frac{2}{d} & \left(\|\boldsymbol{\alpha}\| \mathbf{v}_1^\top \mathbf{U}^\top \mathbf{g} + \sqrt{\rho} \|\boldsymbol{\alpha}\| \mathbf{v}_2^\top \mathbf{U}_\perp^\top \mathbf{g} + \sqrt{\|\mathbf{v}_1\|^2 + \rho \|\mathbf{v}_2\|^2} \mathbf{h}^\top \boldsymbol{\alpha} \right. \\ & \left. - \mathbf{v}_1^\top \mathbf{a} - \frac{1}{2} \|\mathbf{v}_1\|^2 - \frac{1}{2} \|\mathbf{v}_2\|^2 \right). \end{aligned} \quad (33)$$

We next introduce the following change of variables:

$$\begin{aligned} \mathbf{g}_1 &:= \mathbf{U}^\top \mathbf{g} \sim N(0, \mathbf{I}_m), \\ \mathbf{g}_2 &:= \mathbf{U}_\perp^\top \mathbf{g} \sim N(0, \mathbf{I}_{n-m}). \end{aligned}$$

Rewriting the optimization in terms of \mathbf{g}_1 and \mathbf{g}_2 we get

$$\begin{aligned} \min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \max_{\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^n} \frac{2}{d} & \left(\|\boldsymbol{\alpha}\| \mathbf{v}_1^\top \mathbf{g}_1 + \sqrt{\rho} \|\boldsymbol{\alpha}\| \mathbf{v}_2^\top \mathbf{g}_2 + \sqrt{\|\mathbf{v}_1\|^2 + \rho \|\mathbf{v}_2\|^2} \mathbf{h}^\top \boldsymbol{\alpha} \right. \\ & \left. - \mathbf{v}_1^\top \mathbf{a} - \frac{1}{2} \|\mathbf{v}_1\|^2 - \frac{1}{2} \|\mathbf{v}_2\|^2 \right). \end{aligned} \quad (34)$$

We next do the maximization on \mathbf{v}_1 and \mathbf{v}_2 by first fixing the norms to $\beta_1 := \|\mathbf{v}_1\|$ and $\beta_2 := \|\mathbf{v}_2\|$ and maximize over the directions and then maximize over β_1, β_2 . This gives us

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^d} \max_{\beta_1, \beta_2 \geq 0} \frac{2}{d} \left(\beta_1 \|\boldsymbol{\alpha}\| \|\mathbf{g}_1 - \mathbf{a}\| + \beta_2 \sqrt{\rho} \|\boldsymbol{\alpha}\| \|\mathbf{g}_2\| + \sqrt{\beta_1^2 + \rho \beta_2^2} \mathbf{h}^\top \boldsymbol{\alpha} - \frac{\beta_1^2 + \beta_2^2}{2} \right). \quad (35)$$

For minimization over $\boldsymbol{\alpha}$, we first fix its norm to $\eta := \|\boldsymbol{\alpha}\|$ and optimize over its direction, and then over η :

$$\min_{\eta \geq 0} \max_{\beta_1, \beta_2 \geq 0} \frac{2}{d} \left(\beta_1 \|\eta \mathbf{g}_1 - \mathbf{a}\| + \beta_2 \sqrt{\rho} \eta \|\mathbf{g}_2\| - \eta \sqrt{\beta_1^2 + \rho \beta_2^2} \|\mathbf{h}\| - \frac{\beta_1^2 + \beta_2^2}{2} \right). \quad (36)$$

The next step in the CGMT framework is to compute the pointwise limit of the objective function. By concentration of Lipschitz functions of Gaussian vectors we have

$$\begin{aligned} \frac{1}{\sqrt{d}} \|\eta \mathbf{g}_1 - \mathbf{a}\| &\stackrel{(p)}{\rightarrow} \sqrt{\frac{\|\mathbf{a}\|^2}{d} + \eta^2 \frac{\psi}{k}}, \\ \frac{1}{\sqrt{d}} \|\mathbf{g}_2\| &\stackrel{(p)}{\rightarrow} \sqrt{\psi \left(1 - \frac{1}{k}\right)}, \\ \frac{1}{\sqrt{d}} \|\mathbf{h}\| &\stackrel{(p)}{\rightarrow} 1, \end{aligned}$$

where we used Assumption 2.3 by which $n/d \rightarrow \psi$, and Assumption 2.4 by which $m = n/k$. Using these limits in (36), we arrive at the following deterministic optimization problem:

$$\min_{\eta \geq 0} \max_{\beta_1, \beta_2 \geq 0} \beta_1 \sqrt{\frac{\|\mathbf{a}\|^2}{d} + \eta^2 \frac{\psi}{k}} + \beta_2 \sqrt{\rho} \eta \sqrt{\psi \left(1 - \frac{1}{k}\right)} - \eta \sqrt{\beta_1^2 + \rho \beta_2^2} - \frac{\beta_1^2 + \beta_2^2}{2}, \quad (37)$$

where we applied the change of variables $2\beta_1/\sqrt{d} \rightarrow \beta_1$ and $2\beta_2/\sqrt{d} \rightarrow \beta_2$.

In order to find the optimal solution we solve the stationary conditions. By setting derivative with respect to η to zero we obtain

$$\frac{\beta_1 \eta \frac{\psi}{k}}{\sqrt{\frac{\|\mathbf{a}\|^2}{d} + \eta^2 \frac{\psi}{k}}} + \eta \sqrt{\rho \psi \left(1 - \frac{1}{k}\right)} - \sqrt{\beta_1^2 + \rho \beta_2^2} = 0. \quad (38)$$

In addition by setting the derivative with respect to β_1 and β_2 to zero, we obtain

$$\begin{aligned}\sqrt{\frac{\|\mathbf{a}\|^2}{d} + \eta^2 \frac{\psi}{k}} &= \left(\frac{\eta}{\sqrt{\beta_1^2 + \rho\beta_2^2}} + 1 \right) \beta_1, \\ \eta \sqrt{\rho\psi \left(1 - \frac{1}{k}\right)} &= \left(\frac{\rho\eta}{\sqrt{\beta_1^2 + \rho\beta_2^2}} + 1 \right) \beta_2.\end{aligned}\tag{39}$$

By substituting for β_1 and β_2 from (39) into (38) we get

$$\frac{\eta \frac{\psi}{k}}{\eta + c} - 1 + \frac{\rho\eta\psi \left(1 - \frac{1}{k}\right)}{\rho\eta + c} = 0,\tag{40}$$

where $c = \sqrt{\beta_1^2 + \rho\beta_2^2}$.

Also by substituting for β_1 and β_2 from (39) into the definition $c = \sqrt{\beta_1^2 + \rho\beta_2^2}$, we have

$$\frac{\frac{\|\mathbf{a}\|^2}{d} + \eta^2 \frac{\psi}{k}}{(\eta + c)^2} + \frac{\rho^2 \eta^2 \psi \left(1 - \frac{1}{k}\right)}{(\rho\eta + c)^2} = 1.\tag{41}$$

We next make the change of variable: $c = \eta u$, and rewriting equations (40) and (41) as follows:

$$\begin{cases} \frac{\psi}{1+u} + \frac{\rho\psi(k-1)}{\rho+u} &= k, \\ \frac{\frac{k\|\mathbf{a}\|^2}{d\eta^2} + \psi}{(1+u)^2} + \frac{\rho^2\psi(k-1)}{(\rho+u)^2} &= k.\end{cases}$$

Defining $v := \frac{k\|\mathbf{a}\|^2}{\psi d \eta^2}$ we get the system of equations given in the lemma statement.

As the final step, recall that as we discussed at the beginning of the proof, $\boldsymbol{\alpha}_* = \mathbf{M}\mathbf{U}\mathbf{a}$. Therefore,

$$\frac{n}{\|\mathbf{a}\|^2} \|\mathbf{M}\mathbf{U}\mathbf{a}\|_2^2 = \frac{n}{\|\mathbf{a}\|^2} \|\boldsymbol{\alpha}_*\|_2^2 \stackrel{(p)}{\rightarrow} \frac{n}{\|\mathbf{a}\|^2} \eta_*^2 = \frac{k}{v_*},$$

which completes the proof.