

8. Additional Demonstration Plots

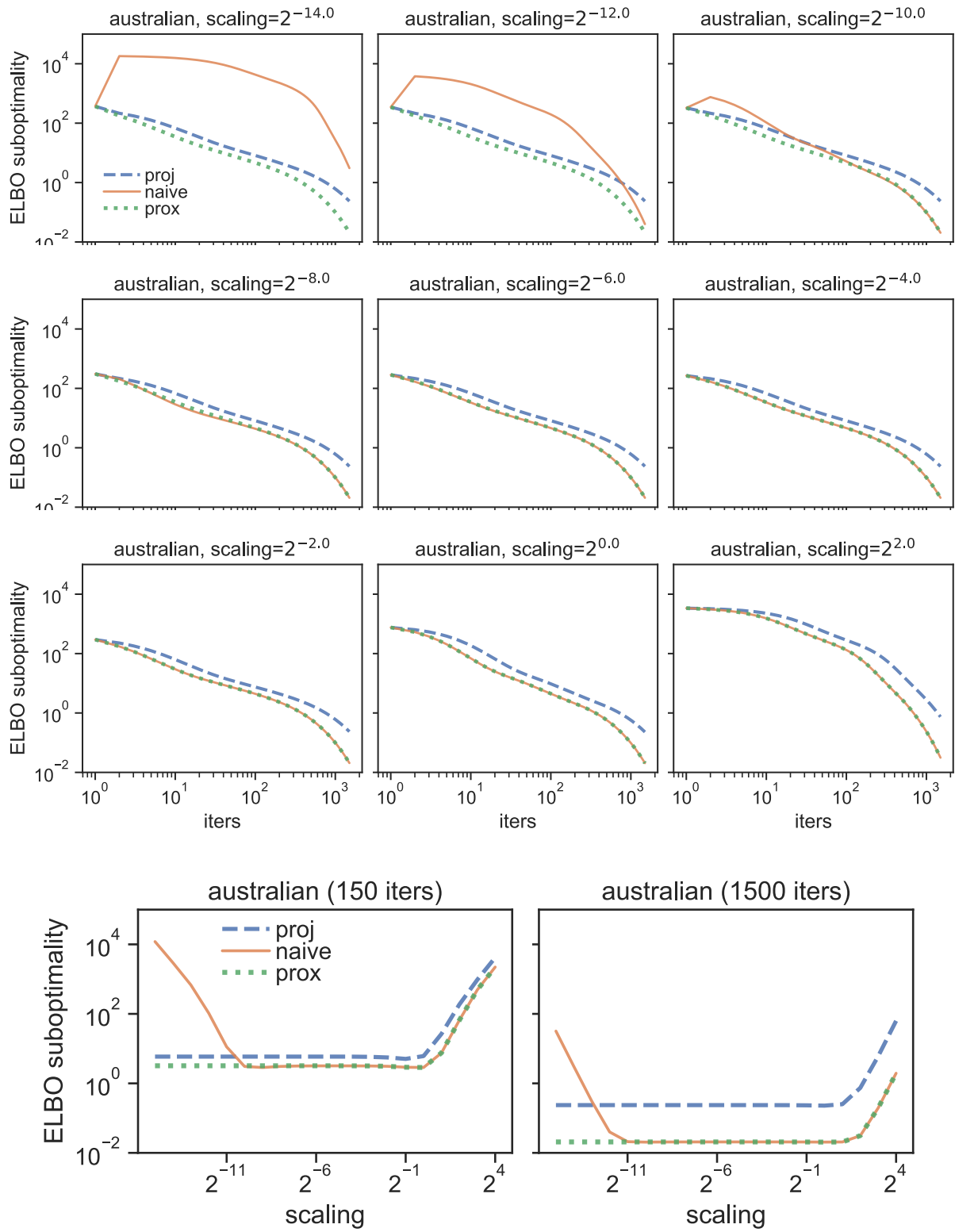


Figure 4. Looseness of the objective obtained by naive gradient descent ($\gamma = 1/M$), projected gradient descent ($\gamma = 1/(2M)$) and proximal gradient descent ($\gamma = 1/M$). Optimization starts with $\mathbf{m} = 0$ and $C = \rho I$ where ρ is a scaling factor.

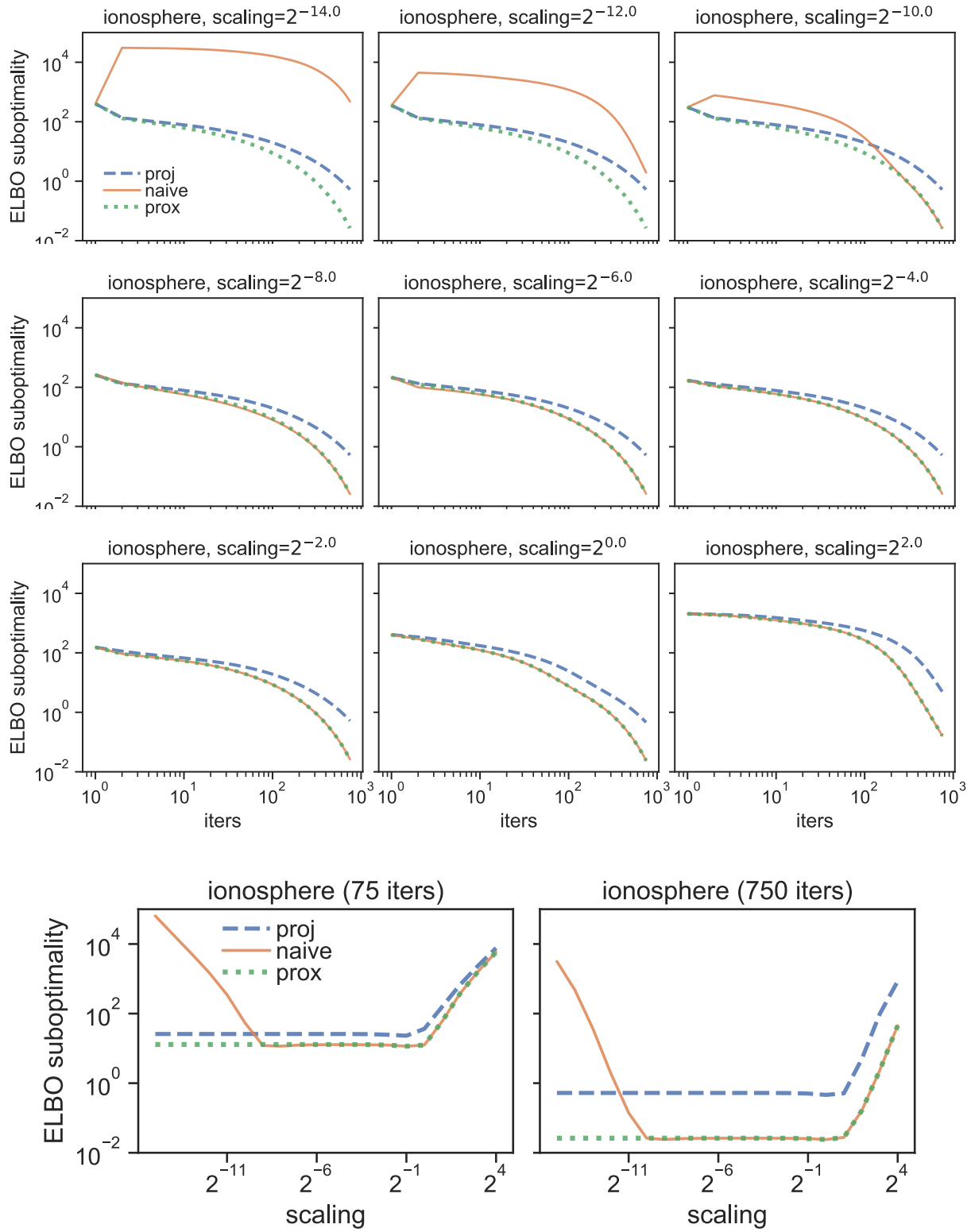


Figure 5. Looseness of the objective obtained by naive gradient descent ($\gamma = 1/M$), projected gradient descent ($\gamma = 1/(2M)$) and proximal gradient descent ($\gamma = 1/M$). Optimization starts with $\mathbf{m} = 0$ and $C = \rho I$ where ρ is a scaling factor.

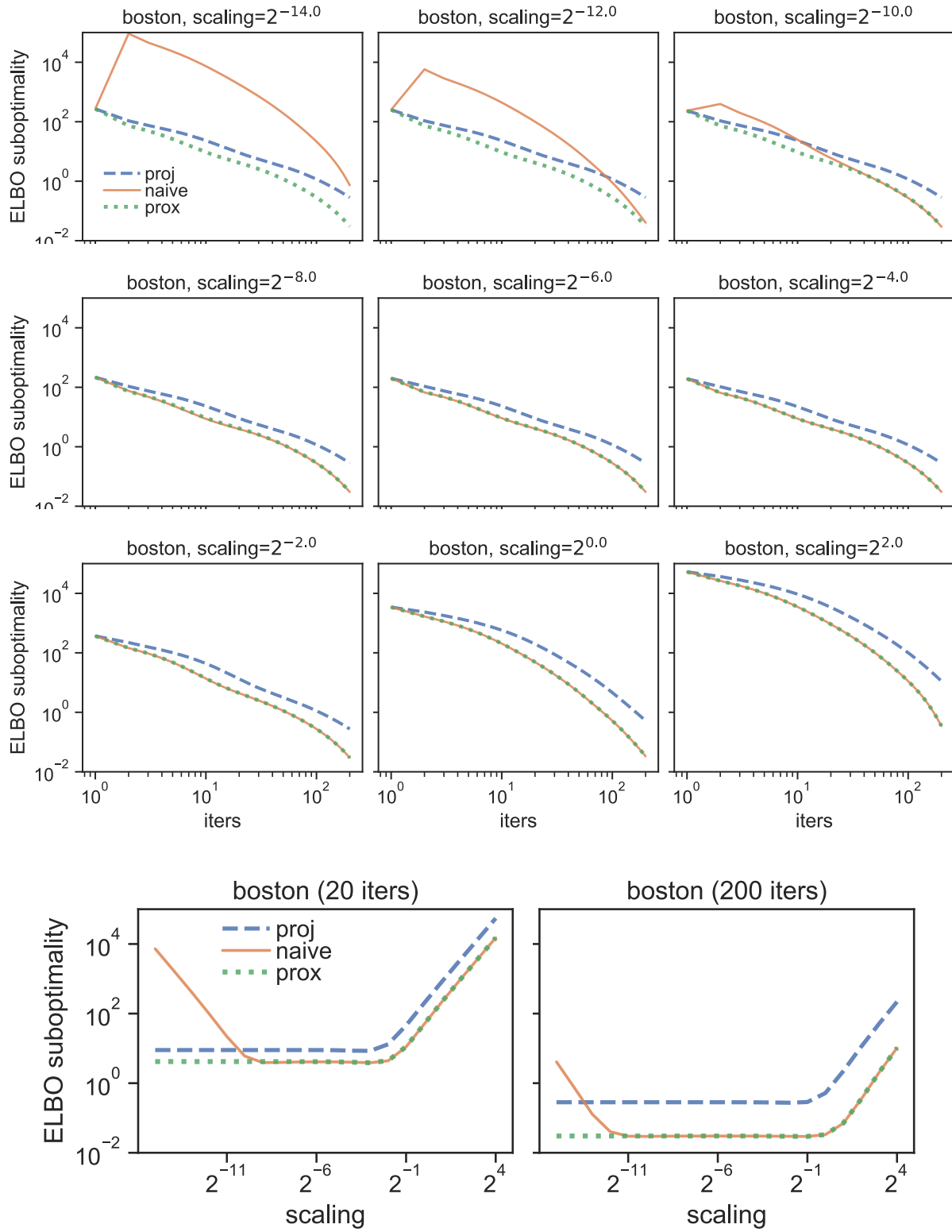


Figure 6. Looseness of the objective obtained by naive gradient descent ($\gamma = 1/M$), projected gradient descent ($\gamma = 1/(2M)$) and proximal gradient descent ($\gamma = 1/M$). Optimization starts with $\mathbf{m} = 0$ and $C = \rho I$ where ρ is a scaling factor.

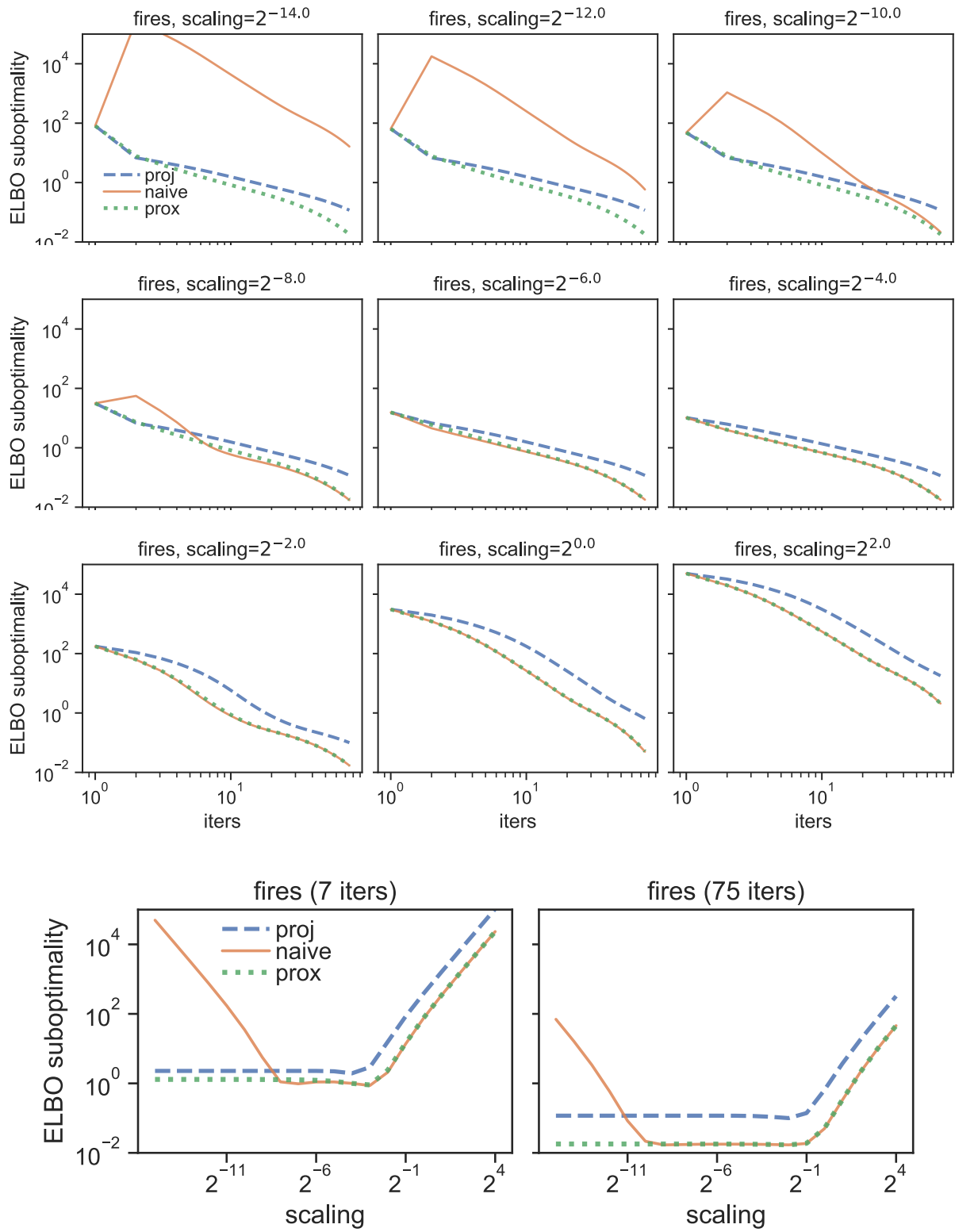


Figure 7. Looseness of the objective obtained by naive gradient descent ($\gamma = 1/M$), projected gradient descent ($\gamma = 1/(2M)$) and proximal gradient descent ($\gamma = 1/M$). Optimization starts with $\mathbf{m} = 0$ and $C = \rho I$ where ρ is a scaling factor.

9. Proofs for Technical Lemmas

This section gives proofs for the technical lemmas used in the main result. Firstly, we show that $\langle \cdot, \cdot \rangle_s$ is a valid inner-product.

Lemma 2. $\langle \mathbf{a}, \mathbf{b} \rangle_s = \mathbb{E}_{\mathbf{u} \sim_s} \mathbf{a}(\mathbf{u})^\top \mathbf{b}(\mathbf{u})$ is a valid inner-product on squared-integrable $\mathbf{a} : \mathbb{R}^d \rightarrow \mathbb{R}^k$.

Proof. The space of square integrable functions is $\{\mathbf{a} : \mathbb{R}^d \rightarrow \mathbb{R}^k \mid \mathbb{E}_{\mathbf{u} \sim_s} a_i(\mathbf{u})^2 \leq \infty \forall i \in \{1, \dots, k\}\}$.

Since each component $a_i(\mathbf{u})$ and $b_i(\mathbf{u})$ is square-integrable with respect to $s(\mathbf{u})$ we know (by Cauchy-Schwarz) that $\mathbb{E}_{\mathbf{u} \sim_s} a_i(\mathbf{u})b_i(\mathbf{u}) \leq \sqrt{\mathbb{E}_{\mathbf{u} \sim_s} a_i(\mathbf{u})^2} \sqrt{\mathbb{E}_{\mathbf{u} \sim_s} b_i(\mathbf{u})^2}$ is finite and real. Therefore, we have by linearity of expectation that

$$\begin{aligned} \sum_{i=1}^k \mathbb{E}_{\mathbf{u} \sim_s} a_i(\mathbf{u})b_i(\mathbf{u}) &= \mathbb{E}_{\mathbf{u} \sim_s} \sum_{i=1}^k a_i(\mathbf{u})b_i(\mathbf{u}) \\ &= \mathbb{E}_{\mathbf{u} \sim_s} \mathbf{a}(\mathbf{u})^\top \mathbf{b}(\mathbf{u}) \\ &= \langle \mathbf{a}, \mathbf{b} \rangle_s \end{aligned}$$

is finite and real for all $\mathbf{a}, \mathbf{b} \in V_s$. To show that $(V_s, \langle \cdot, \cdot \rangle_s)$ is a valid inner-product space, it is easy to establish all the necessary properties of the inner-product, namely for all $\mathbf{a}, \mathbf{b}, \mathbf{c} \in V_s$,

$$\langle \mathbf{a}, \mathbf{b} \rangle = \langle \mathbf{b}, \mathbf{a} \rangle$$

$$\langle \theta \mathbf{a}, \mathbf{b} \rangle = \theta \langle \mathbf{a}, \mathbf{b} \rangle \text{ for } \theta \in \mathbb{R}$$

$$\langle \mathbf{a} + \mathbf{b}, \mathbf{c} \rangle = \langle \mathbf{a}, \mathbf{c} \rangle + \langle \mathbf{b}, \mathbf{c} \rangle$$

$$\langle \mathbf{a}, \mathbf{a} \rangle \geq 0$$

$\langle \mathbf{a}, \mathbf{a} \rangle = 0 \Leftrightarrow \mathbf{a} = \mathbf{0}$. (Where $\mathbf{0}(\varepsilon)$ is a function that always returns a vector of k zeros.) \square

Next, we give three technical Lemmas, which do most of the work of the proof.

Lemma 3. Let $\mathbf{a}_i(\mathbf{u}) = \frac{d}{dw_i} \mathbf{t}_w(\mathbf{u})$. This is independent of w and $\frac{dl(\mathbf{w})}{dw_i} = \langle \mathbf{a}_i, \nabla f \circ \mathbf{t}_w \rangle_s$.

Proof. Now, we can write $l(\mathbf{w})$ as

$$l(\mathbf{w}) = \mathbb{E}_{\mathbf{z} \sim_{q_w}} f(\mathbf{z}) = \mathbb{E}_{\mathbf{u} \sim_s} f(\mathbf{t}_w(\mathbf{u})).$$

Since $\mathbf{t}_w(\mathbf{u}) = C\mathbf{u} + \mathbf{m}$ is an affine function, it's easy to see that both $\frac{d}{dC_{ij}} \mathbf{t}_w(\mathbf{u})$ and $\frac{d}{dm_i} \mathbf{t}_w(\mathbf{u})$ are independent of w . Therefore, the gradient of $l(\mathbf{w})$ can be written as

$$\begin{aligned} \nabla_{w_i} l(\mathbf{w}) &= \nabla_{w_i} \mathbb{E}_{\mathbf{u} \sim_s} f(\mathbf{t}_w(\mathbf{u})) \\ &= \mathbb{E}_{\mathbf{u} \sim_s} \nabla_{w_i} \mathbf{t}_w(\mathbf{u})^\top \nabla f(\mathbf{t}_w(\mathbf{u})). \\ &= \langle \mathbf{a}_i, \nabla f \circ \mathbf{t}_w \rangle_s. \end{aligned}$$

\square

Lemma 4. If s is standardized, then the functions $\{\mathbf{a}_i\}$ are orthonormal in $\langle \cdot, \cdot \rangle_s$.

Proof. It is easy to calculate that

$$\begin{aligned} \frac{d}{dm_i} \mathbf{t}_w(\mathbf{u}) &= \mathbf{e}_i \\ \frac{d}{dC_{ij}} \mathbf{t}_w(\mathbf{u}) &= \mathbf{e}_i u_j, \end{aligned}$$

where \mathbf{e}_i is the indicator vector in the i -th component. Therefore, we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{u} \sim_s} \left(\frac{d}{dm_i} \mathbf{t}_w(\mathbf{u}) \right)^\top \left(\frac{d}{dm_j} \mathbf{t}_w(\mathbf{u}) \right) &= \mathbb{E}_{\mathbf{u} \sim_s} \mathbf{e}_i^\top \mathbf{e}_j \\ &= I[i = j] \\ \mathbb{E}_{\mathbf{u} \sim_s} \left(\frac{d}{dC_{ij}} \mathbf{t}_w(\mathbf{u}) \right)^\top \left(\frac{d}{dC_{kl}} \mathbf{t}_w(\mathbf{u}) \right) &= \mathbb{E}_{\mathbf{u} \sim_s} u_j \mathbf{e}_i^\top \mathbf{e}_k \\ &= I[i = k] \mathbb{E}_{\mathbf{u} \sim_s} u_j \\ &= 0 \\ &\text{(since zero mean)} \\ \mathbb{E}_{\mathbf{u} \sim_s} \left(\frac{d}{dC_{ij}} \mathbf{t}_w(\mathbf{u}) \right)^\top \left(\frac{d}{dC_{kl}} \mathbf{t}_w(\mathbf{u}) \right) &= \mathbb{E}_{\mathbf{u} \sim_s} u_j u_l \mathbf{e}_i^\top \mathbf{e}_k \\ &= I[i = k] \mathbb{E}_{\mathbf{u} \sim_s} u_j u_l \\ &= I[i = k] I[j = l] \\ &\text{(since unit variance and zero mean)} \end{aligned}$$

These three identities are equivalent to stating that $\{\mathbf{a}_i\}$ are orthonormal in $\langle \cdot, \cdot \rangle_s$. \square

Lemma 5. If s is standardized, then $\mathbb{E}_{\mathbf{u} \sim_s} \|\mathbf{t}_w(\mathbf{u}) - \mathbf{t}_v(\mathbf{u})\|_2^2 = \|\mathbf{w} - \mathbf{v}\|_2^2$.

Proof. Let $\Delta \mathbf{m}$ and ΔS denote the difference of the \mathbf{m} and S parts of \mathbf{w} , respectively. We want to calculate

$$\begin{aligned} \mathbb{E}_{\mathbf{u} \sim_s} \|\mathbf{t}_w(\mathbf{u}) - \mathbf{t}_v(\mathbf{u})\|_2^2 &= \mathbb{E}_{\mathbf{u} \sim_s} \|\Delta C \varepsilon + \Delta \mathbf{m}\|_2^2 \\ &= \mathbb{E}_{\mathbf{u} \sim_s} \left(\|\Delta C \mathbf{u}\|_2^2 + 2\Delta \mathbf{m}^\top \Delta C \mathbf{u} + \|\Delta \mathbf{m}\|_2^2 \right). \end{aligned}$$

It is easy to see that the expectation of the middle term is zero, and the last is a constant. The expectation of the first

term is

$$\begin{aligned}
 \mathbb{E}_{\mathbf{u} \sim s} \|(\Delta C)\mathbf{u}\|_2^2 &= \mathbb{E}_{\mathbf{u} \sim s} \mathbf{u}^\top (\Delta C)^\top (\Delta C) \mathbf{u} \\
 &= \mathbb{E}_{\mathbf{u} \sim s} \text{tr}(\mathbf{u}^\top (\Delta C)^\top (\Delta C) \mathbf{u}) \\
 &= \mathbb{E}_{\mathbf{u} \sim s} \text{tr}((\Delta C)^\top (\Delta C) \mathbf{u} \mathbf{u}^\top) \\
 &= \text{tr}((\Delta C)^\top (\Delta C)) = \|\nabla C\|_F^2. \\
 &\quad (\text{since zero mean and unit variance})
 \end{aligned}$$

Putting this together gives that

$$\begin{aligned}
 \mathbb{E}_{\mathbf{u} \sim s} \|\mathbf{t}_{\mathbf{w}}(\mathbf{u}) - \mathbf{t}_{\mathbf{v}}(\mathbf{u})\|_2^2 &= \|\Delta C\|_F^2 + \|\Delta \mathbf{m}\|_2^2 \\
 &= \|\mathbf{w} - \mathbf{v}\|_2^2.
 \end{aligned}$$

□

10. Proof for Example Function

Theorem 6. Let $q_w = \text{LocScale}(\mathbf{m}, C, s)$ with parameters $w = (\mathbf{m}, C)$ and a standardized base distribution s and let $f(\mathbf{z}) = \frac{a}{2} \|\mathbf{z} - \mathbf{z}^*\|_2^2$. Then $l(w) = \mathbb{E}_{\mathbf{z} \sim q_w} f(\mathbf{z}) = \frac{a}{2} (\|\mathbf{m} - \mathbf{z}^*\|_2^2 + \|C\|_F^2)$.

Proof. For a general distribution, we have that

$$\begin{aligned} \mathbb{E} f(\mathbf{z}) &= \frac{a}{2} \mathbb{E} \|\mathbf{z} - \mathbb{E}[\mathbf{z}] + \mathbb{E}[\mathbf{z}] - \mathbf{z}^*\|_2^2 \\ &= \frac{a}{2} \mathbb{E} \left(\|\mathbf{z} - \mathbb{E}[\mathbf{z}]\|_2^2 \right. \\ &\quad \left. + 2(\mathbf{z} - \mathbb{E}[\mathbf{z}])^\top (\mathbb{E}[\mathbf{z}] - \mathbf{z}^*) + \|\mathbb{E}[\mathbf{z}] - \mathbf{z}^*\|_2^2 \right) \\ &= \frac{a}{2} \left(\text{tr } \mathbb{V}[\mathbf{z}] + \|\mathbb{E}[\mathbf{z}] - \mathbf{z}^*\|_2^2 \right). \end{aligned}$$

Now, if q_w is a location-scale family, we have that $\mathbf{z} = C\mathbf{u} + \mathbf{m}$. Thus,

$$\begin{aligned} \text{tr } \mathbb{V}[\mathbf{z}] &= \text{tr } \mathbb{V}[C\mathbf{u} + \mathbf{m}] \\ &= \text{tr } \mathbb{V}[C\mathbf{u}] \\ &= \text{tr } C \mathbb{V}[\mathbf{u}] C^\top \\ &= \text{tr } C C^\top \mathbb{V}[\mathbf{u}]. \end{aligned}$$

Meanwhile, we have that

$$\begin{aligned} \|\mathbb{E}[\mathbf{z}] - \mathbf{z}^*\|_2^2 &= \|\mathbb{E}[C\mathbf{u} + \mathbf{m}] - \mathbf{z}^*\|_2^2 \\ &= \|C \mathbb{E}[\mathbf{u}] + \mathbf{m} - \mathbf{z}^*\|_2^2 \end{aligned}$$

Thus,

$$\mathbb{E} f(\mathbf{z}) = \frac{a}{2} \left(\text{tr } C \mathbb{V}[\mathbf{u}] C^\top + \|C \mathbb{E}[\mathbf{u}] + \mathbf{m} - \mathbf{z}^*\|_2^2 \right).$$

The case where s is standardized follows from substituting $\mathbb{E}[\mathbf{u}] = 0$ and $\mathbb{V}[\mathbf{u}] = I$ and applying the fact that $\text{tr } C C^\top = \|C\|_F^2$. \square

11. Proofs for Solution Guarantees

Lemma 8. Let $q_{\mathbf{w}} = \text{LocScale}(\mathbf{m}, C, s)$ with parameters $\mathbf{w} = (\mathbf{m}, C)$ and a standardized and spherically symmetric base distribution s . Let $l(\mathbf{w}) = \mathbb{E}_{\mathbf{z} \sim q_{\mathbf{w}}} f(\mathbf{z})$. Suppose C is diagonal and f is M -smooth. Then, $|\frac{dl(\mathbf{w})}{dC_{ii}}| \leq M|C_{ii}|$.

Proof. Define \mathbf{w}' to be \mathbf{w} but with C_{ii} set to zero. We will first show that $\frac{dl(\mathbf{w}')}{dC_{ii}} = 0$. Using the definition of $\mathbf{t}_{\mathbf{w}}$ and the fact that $\frac{d}{dC_{ij}} \mathbf{t}_{\mathbf{w}}(u) = \mathbf{e}_i u_j$ gives that

$$\frac{d}{dC_{ii}} l(\mathbf{w}') = \mathbb{E}_{\mathbf{u} \sim s} \frac{d}{dC_{ii}} f(\mathbf{t}_{\mathbf{w}'}(\mathbf{u})) \quad (11)$$

$$= \mathbb{E}_{\mathbf{u} \sim s} \mathbf{u}_i \mathbf{e}_i^\top \nabla f(\mathbf{t}_{\mathbf{w}'}(\mathbf{u})) \quad (12)$$

$$= 0. \quad (13)$$

The final equality above follows from the facts that $\mathbb{E} \mathbf{u}_i = 0$ and $\mathbf{u}_i \perp \mathbf{e}_i^\top \nabla f(\mathbf{t}_{\mathbf{w}'}(\mathbf{u}))$ (Since $\mathbf{t}_{\mathbf{w}'}(\mathbf{u})$ ignores \mathbf{u}_i) so the expectation in Eq. (11) is over two independent random variables, one with mean zero. Now, by Thm. 1, l is also M -smooth, thus

$$\begin{aligned} \left| \frac{dl(\mathbf{w})}{dC_{ii}} \right| &= \left| \frac{dl(\mathbf{w}')}{dC_{ii}} - \frac{dl(\mathbf{w})}{dC_{ii}} \right| \\ &\leq \|\nabla l(\mathbf{w}') - \nabla l(\mathbf{w})\|_2 \\ &\leq M \|\mathbf{w}' - \mathbf{w}\|_2 \\ &= M |C_{ii}|. \end{aligned}$$

□

Theorem 7. Let $q_{\mathbf{w}} = \text{LocScale}(\mathbf{m}, C, s)$ with parameters $\mathbf{w} = (\mathbf{m}, C)$ and a standardized and spherically symmetric base distribution s . Suppose \mathbf{w} minimizes $l(\mathbf{w}) + h(\mathbf{w})$ from Eq. (1) and $\log p(\mathbf{z}, \mathbf{x})$ is M -smooth over \mathbf{z} . Then, $\mathbf{w} \in \mathcal{W}_M$.

Proof. First, suppose that C is diagonal. Since \mathbf{w} minimizes $l + h$, $\nabla l(\mathbf{w}) = -\nabla h(\mathbf{w})$. The gradient of h with respect to C is $-C^{-\top}$. Thus, $|\frac{dl(\mathbf{w})}{dC_{ii}}| = |\frac{dh(\mathbf{w})}{dC_{ii}}| = \frac{1}{|C_{ii}|}$. But by Lem. 8, $|\frac{dl(\mathbf{w})}{dC_{ii}}| \leq M|C_{ii}|$. This establishes the claim for diagonal C .

Now, consider some non-diagonal C . Let the singular value decomposition be $C = USV^\top$. Define $f_U(\mathbf{z}) = f(U\mathbf{z})$ and define l_U with respect to f_U . Let $\mathbf{w}' = (S, U^\top \mathbf{m})$. Then, the following statements are equivalent to $\mathbf{w} \in \text{argmin}_{\mathbf{w}} l(\mathbf{w}) + h(\mathbf{w})$:

$$\begin{aligned} (C, \mathbf{m}) &\in \text{argmin}_{(C, \mathbf{m})} \mathbb{E}_{\mathbf{u} \sim s} f(C\mathbf{u} + \mathbf{m}) - \log |C| \\ &\Leftrightarrow (S, \mathbf{m}) \in \text{argmin}_{(S, \mathbf{m})} \mathbb{E}_{\mathbf{u} \sim s} f(USV^\top \mathbf{u} + \mathbf{m}) - \log |USV^\top| \\ &\Leftrightarrow (S, \mathbf{m}) \in \text{argmin}_{(S, \mathbf{m})} \mathbb{E}_{\mathbf{u} \sim s} f(US\mathbf{u} + \mathbf{m}) - \log |S| \\ &\Leftrightarrow (S, \mathbf{m}) \in \text{argmin}_{(S, \mathbf{m})} \mathbb{E}_{\mathbf{u} \sim s} f_U(S\mathbf{u} + U^\top \mathbf{m}) - \log |S| \\ &\Leftrightarrow \mathbf{w}' \in \text{argmin}_{\mathbf{w}'} l_U(\mathbf{w}') + h(\mathbf{w}'). \end{aligned}$$

Thus, \mathbf{w} minimizing $l + h$ is equivalent to \mathbf{w}' minimizing $l_U + h$. Since f_U is M -smooth and S is diagonal, we know that $S_{ii} \geq \frac{1}{\sqrt{M}}$ for all. □

12. Proofs with Convexity

Theorem 10. *Let $q_{\mathbf{w}} = \text{LocScale}(\mathbf{m}, C, s)$ with parameters $\mathbf{w} = (\mathbf{m}, C)$ and a standardized and spherically symmetric base distribution s . Suppose \mathbf{w} minimizes $l(\mathbf{w}) + h(\mathbf{w})$ from Eq. (1) and $-\log p(\mathbf{z}, \mathbf{x})$ is c -strongly convex over \mathbf{z} . Then, $\|C\|_F^2 + \|\mathbf{m} - \mathbf{z}^*\|_2^2 \leq \frac{d}{c}$, where $\mathbf{z}^* = \text{argmax}_{\mathbf{z}} \log(\mathbf{z}, \mathbf{x})$.*

It's easy to see that l is minimized by $\bar{\mathbf{w}} = (\mathbf{z}^*, \mathbf{0}_{d \times d})$. By Thm. 9, $l(\mathbf{w})$ is c -strongly convex. Thus applying a standard inner-product result on strong convexity (Nesterov, 2014, Thm. 2.1.9),

$$\begin{aligned}
 c \|\mathbf{w} - \bar{\mathbf{w}}\|_2^2 &\leq \langle \nabla l(\mathbf{w}) - \nabla l(\bar{\mathbf{w}}), \mathbf{w} - \bar{\mathbf{w}} \rangle \\
 &\quad (\text{since } l \text{ is strongly convex}) \\
 &= \langle \nabla l(\mathbf{w}), \mathbf{w} - \bar{\mathbf{w}} \rangle \\
 &\quad (\text{since } \nabla l(\bar{\mathbf{w}}) = 0) \\
 &= -\langle \nabla h(\mathbf{w}), \mathbf{w} - \bar{\mathbf{w}} \rangle \\
 &\quad (\text{since } \nabla l(\mathbf{w}) + \nabla h(\mathbf{w}) = 0) \\
 &= \text{tr}(C^{-\top} C) \\
 &\quad (\text{since } \nabla_C h(\mathbf{w}) = -C^{-\top}, \nabla_{\mathbf{m}} h(\mathbf{w}) = 0). \\
 &= \text{tr } I = d.
 \end{aligned}$$

The result follows from observing that $\|\mathbf{w} - \bar{\mathbf{w}}\|_2^2 = \|C\|_F^2 + \|\mathbf{m} - \mathbf{z}^*\|_2^2$.

13. Convergence Considerations

Lemma 12. Let $q_w = \text{LocScale}(\mathbf{m}, C, s)$ with parameters $\mathbf{w} = (\mathbf{m}, C)$. Then, $h(\mathbf{w}) = \mathbb{E}_{\mathbf{z} \sim q_w} [\log q_w(\mathbf{z})]$ is M -smooth over \mathcal{W}_M .

Proof. Take $\mathbf{w} = (C, \mathbf{m}) \in \mathcal{W}_M$ and $\mathbf{v} = (B, \mathbf{n}) \in \mathcal{W}_M$. We write $h(C)$ since $h(\mathbf{w})$ is independent of \mathbf{m} . The gradient is $\nabla h(C) = C^{-T}$. Now, use that $\|AX\|_F \leq \|A\|_2 \|X\|_F$ to get that

$$\begin{aligned} \|\nabla h(B) - \nabla h(C)\|_F &= \|B^{-1} - C^{-1}\|_F \\ &= \|B^{-1}(B - C)C^{-1}\|_F \\ &\leq \|B^{-1}\|_2 \|C^{-1}\|_2 \|B - C\|_F. \end{aligned}$$

But, since $\mathbf{w} \in \mathcal{W}_M$, $\|C^{-1}\|_2 = \frac{1}{\sigma_{\min}(C)} \leq \sqrt{M}$ and similarly for B . This establishes that $\|\nabla h(B) - \nabla h(C)\|_F \leq M \|B - C\|_F$, equivalent to the result. \square

Theorem 13. Suppose $h(\mathbf{w})$ corresponds to a location-scale family with a standardized s , and $\mathbf{w} = (\mathbf{m}, C)$.

- If C has singular value decomposition $C = USV^\top$, then $\text{proj}_{\mathcal{W}_M}(\mathbf{w}) = (\mathbf{m}, UTV^\top)$, where T is a diagonal matrix with $T_{ii} = \max\left(S_{ii}, \frac{1}{\sqrt{M}}\right)$.
- If C is triangular with a positive diagonal, then $\text{prox}_\gamma(\mathbf{w}) = (\mathbf{m}, C + \Delta C)$, where ΔC is a diagonal matrix with $\Delta C_{ii} = \frac{1}{2} \left(\sqrt{C_{ii}^2 + 4\gamma} - C_{ii} \right)$.

Proof. (Proximal Operator) We know that $h(\mathbf{w}) = \text{Const.} - \log |C|$. Write $\mathbf{w} = (\mathbf{m}, C)$ and $\mathbf{v} = (\mathbf{n}, B)$. Then, we can write the proximal operator as

$$\text{prox}_\lambda(\mathbf{w}) = \underset{\mathbf{v}}{\text{argmin}} -\log |B| + \frac{1}{2\lambda} \|\mathbf{v} - \mathbf{w}\|_2^2$$

Now, assuming that C is triangular, the solution will leave all entries of \mathbf{w} other than the diagonal entries of C unchanged. Then, we will have that $\log |B| = \sum_{i=1}^d \log B_{ii}$. Since

$$\underset{x>0}{\text{argmin}} -\log x + \frac{1}{2\lambda} (x - y)^2 = \frac{y + \sqrt{y^2 + 4\lambda}}{2}$$

The solution is to set

$$\begin{aligned} B_{ii} &= \frac{1}{2} \left(C_{ii} + \sqrt{C_{ii}^2 + 4\lambda} \right) \\ &= C_{ii} + \frac{1}{2} \left(\sqrt{C_{ii}^2 + 4\lambda} - C_{ii} \right). \end{aligned}$$

(Projection Operator) Von-Neumann's trace inequality states that $|\text{tr} A^\top B| \leq \sum_i \sigma_i(A) \sigma_i(B)$. Consider any candidate solution B with SVD QTP^\top . Then, we can write that

$$\begin{aligned} \|B - C\|_F^2 &= \text{tr} (B - C)^\top (B - C) \\ &= \|B\|_F^2 - 2 \text{tr}(B^\top C) + \|C\|_F^2 \\ &\geq \|T\|_F^2 - 2 \sum_i T_{ii} S_{ii} + \sum_i S_{ii}^2 \\ &= \sum_i (T_{ii} - S_{ii})^2. \end{aligned}$$

We can minimize this lower bound by choosing $T_{ii} = \max(1/\sqrt{M}, S_{ii})$, with a corresponding value of $\sum_i \max(0, 1/\sqrt{M} - S_{ii})^2$. Thus any valid solution will have $\|B - C\|_F^2$ at least this large.

However, suppose we choose $B = UT_i V^\top$ with T_{ii} as above. Then,

$$\|B - C\|_F^2 = \|UTV^\top - USV^\top\|_F^2 = \sum_i (T_{ii} - S_{ii})^2,$$

so this value B is optimal. \square