

## A Algorithmic Details

To summarise the main steps of the submission, we now detail a pseudo-code for distributed multitask reinforcement learning.

---

### Algorithm 1 Distributed Multitask Reinforcement Learning

---

**Input:** A set of  $T$  MDPs, dimensions of the latent space  $k$ , parameter initialisation for  $\Theta_{\text{sh}}$ ,  $\tilde{\theta}_1, \dots, \tilde{\theta}_T, \phi_1, \dots, \phi_T$ , precision parameter  $\epsilon$ .

**Step 1:** Optimise for the variational parameters  $\phi_1, \dots, \phi_T$  by solving:

$$\max_{\phi_1: \phi_T} \sum_{t=1}^T \mathbb{E}_{q_{\phi_t}(\tau_t)} \left[ \log p(\hat{\mathcal{R}}_t | \tau_t) \right] - \sum_{t=1}^T \mathcal{D}_{\text{KL}}(q_{\phi_t}(\tau_t) || p_{\theta_t}(\tau_t)).$$

**Step 2:** Given updated  $\phi_1, \dots, \phi_T$ , solve for  $\Theta_{\text{sh}}$  as follows:

**Step 2.1:** Distribute tasks among a graph  $\mathcal{G}$  of  $n$  computational units

**Step 2.2:** Update  $\Theta_{\text{sh}}$  using our distributed Newton method up-to precision  $\epsilon$  (Section 4.2)

**Step 3:** Given updated  $\phi_1, \dots, \phi_T$ , and  $\Theta_{\text{sh}}$ , determine the task specific coefficients by solving:

$$\max_{\theta_1: \theta_T} \sum_{t=1}^T \mathbb{E}_{q_{\phi_t}(\tau_t)} \left[ \log p(\hat{\mathcal{R}}_t | \tau_t) \right] - \sum_{t=1}^T \mathcal{D}_{\text{KL}}(q_{\phi_t}(\tau_t) || p_{\theta_t}(\tau_t)).$$

**Output:** Variational, shared, and task specific parameters.

---

## B Additional Experiments

This section details additional experimental results reflecting consensus errors on various systems. We ran our experiments on the benchmarks depicted in Figure 1.

### B.1 Additional Results

Clearly, our algorithm converges faster than others to low-consensus error and optimal objective values.

## C Theoretical Guarantees

For the clarity of the presentation we split the Appendix in several sections. In the first section, we provide theoretical analysis for Distributed Chebyshev Solver for solving SDD systems. In the second section we provide the convergence analysis for Distributed Newton network and prove Theorem ??

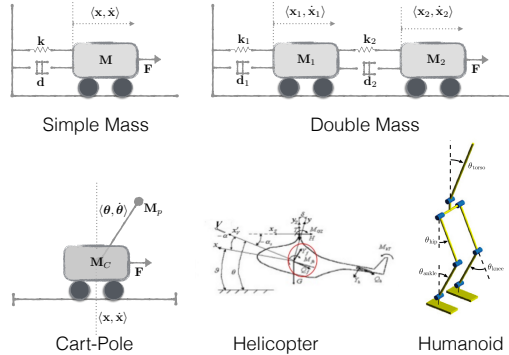


Figure 1: A high-level depiction of the benchmark dynamical systems used in our experiments.

	SDD-Newton	ADD-Newton	Netw-Newton	Dist-ADMM	Dist-Average	Dist-Gradient
SM	$10^1$	$10^2$	$\sim 10^2$	$\sim 10^4$	$\sim 10^4$	$\sim 10^5$
DM	$10^2$	$10^3$	$\sim 10^5$	$\sim 10^4$	$\sim 10^5$	$\sim 10^5$
CP	$10^3$	$\sim 10^4$	$\sim 10^5$	$\sim 10^5$	$\sim 10^5$	$\sim 10^5$

Figure 2: Number of iterations needed for convergence to low consensus showing that our method outperforms state-of-the-art techniques.

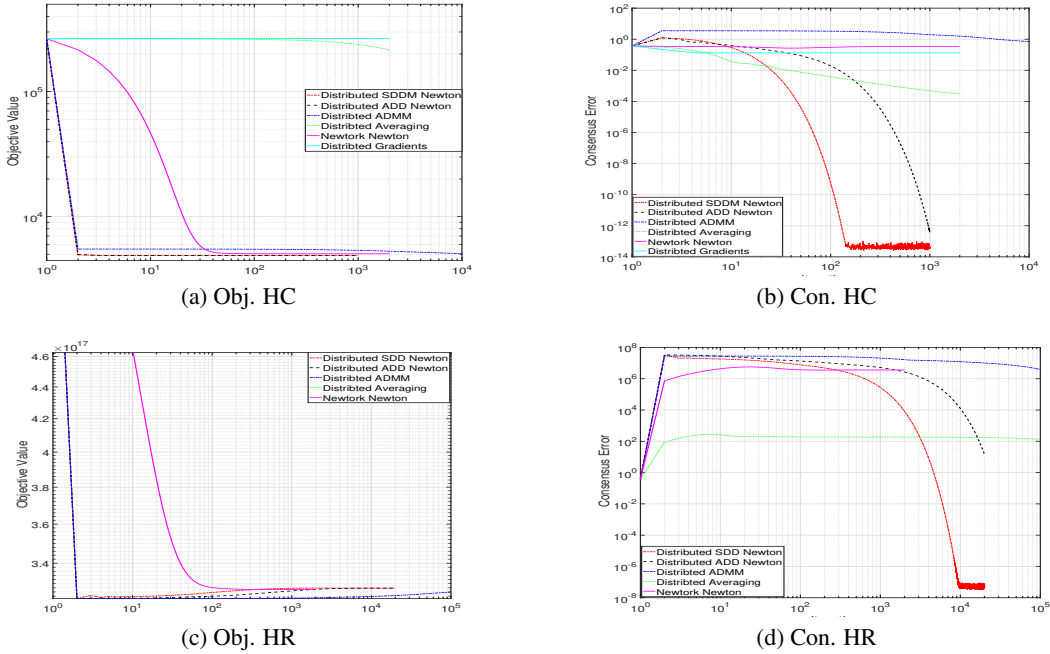


Figure 3: Figures (a) and (b) report the objective value and consensus error versus iterations on HC systems. Figures (c) and (d) demonstrate the same criteria on the humanoid tasks. In all these cases, our method outperforms others in literature.

### C.1 Distributed Chebyshev Solver

The approximated solution vector  $\mathbf{d}_s^{(m)}$  can be represented as a concatenation of  $dk$  vectors:  $\mathbf{d}_s^{(m)} = [\mathbf{d}_s^{(m),1,\top}, \dots, \mathbf{d}_s^{(m),dk,\top}]^\top$ , where each vector  $\mathbf{d}_s^{(m),i}$  is an  $\epsilon$ -approximated solution of the system

$$\mathcal{L}\mathbf{d}_s^{(m),i} = \mathbf{b}_s^i \quad (1)$$

with vector  $\mathbf{b}_s^i$  is the  $i^{th}$  chunk of vector  $\mathbf{b}_s = \sum_{i=1}^n \sum_{t=1}^{T_i} \nabla^2 \mathcal{J}_{MTRL}(\mathbf{y}(\lambda_s)) \mathbf{y}(\lambda_s)$ . Please notice, that each vector  $\mathbf{b}_s^i \in \mathbb{R}^n$  and it is distributed across the nodes of graph  $\mathcal{G}$ . Indeed, the  $r^{th}$  component of this vector can be computed as follows:

$$[\mathbf{b}_s^i]_r = \sum_{j=1}^{dk} \frac{\partial^2 \left[ \sum_{t=1}^{T_i} \mathcal{J}_{MTRL}(\mathbf{y}(\lambda_s)) \right]}{\partial [\mathbf{y}_i]_r \partial [\mathbf{y}_j]_r} [\mathbf{y}(\lambda_s)]_j \quad (2)$$

where we represented primal variable  $\mathbf{y}(\lambda_s)$  as concatenation  $\mathbf{y}(\lambda_s) = [\mathbf{y}(\lambda_s)_1^\top, \dots, \mathbf{y}(\lambda_s)_{dk}^\top]^\top$ . Indeed, Equation (2) can be computed locally by node  $r \in \mathcal{V}$  because it stores the  $r^{th}$  components of vectors  $\mathbf{y}(\lambda_s)_1, \dots, \mathbf{y}(\lambda_s)_{dk}$  as well as it stores the local primal objective  $\sum_{t=1}^{T_i} \mathcal{J}_{MTRL}^{(t)}(\mathbf{y}(\lambda_s))$ .

This observation allows us to distribute the computation of an  $\epsilon$ -approximated solution vector  $\mathbf{d}_s^{(m),i}$  using Chebyshev polynomials:

$$\mathbf{d}_s^{(m),i} = \mathcal{L}^\dagger (\mathbf{I} - \mathcal{Q}_m(\mathcal{L})) \mathbf{b}_k^i, \quad i = 1 \dots, dk. \quad (3)$$

where  $m = \lceil \frac{1}{2} (\sqrt{\kappa(\mathcal{L})} + 1) \ln \frac{2}{\epsilon} \rceil$  and

$$\mathcal{Q}_m(z) = \frac{T_m\left(\frac{(\mu_n + \mu_1) - 2z}{\mu_n - \mu_1}\right)}{T_m\left(\frac{\mu_n + \mu_1}{\mu_n - \mu_1}\right)} \quad (4)$$

where  $\mu_1, \mu_n$  are smallest and largest non-zero eigenvalues of graph Laplacian  $\mathcal{L}$ . Please notice, that for any  $z \in [\mu_1, \mu_p]$  we have

$$|\mathcal{Q}_m(z)|^2 \leq T_m^{-2} \left( \frac{\mu_p + \mu_1}{\mu_p - \mu_1} \right) = T_m^{-2} \left( \frac{\kappa(\mathcal{L}) + 1}{\kappa(\mathcal{L}) - 1} \right) \leq 4e^{-\frac{4m}{\sqrt{\kappa(\mathcal{L})} + 1}}$$

where  $\kappa(\mathcal{L}) = \frac{\mu_p}{\mu_1}$ . Let  $\mathbf{d}_s^{*,i}$  be the exact solution of system (1), then a solution vector in Equation (3) satisfies  $\|\mathbf{d}_s^{(m),i} - \mathbf{d}_s^{*,i}\|_{\mathcal{L}}^2 \leq 4e^{-\frac{4m}{\sqrt{\kappa(\mathcal{L})} + 1}} \|\mathbf{d}_s^{*,i}\|_{\mathcal{L}}^2$ . Hence, by choosing degree  $m = \lceil \frac{1}{2} (\sqrt{\kappa(\mathcal{L})} + 1) \ln \frac{2}{\epsilon} \rceil$  the solution vector

$$\mathbf{d}_s^{(m),i} = \mathcal{L}^\dagger \left[ \frac{T_m\left(\frac{\mu_p + \mu_1}{\mu_p - \mu_1}\right) \mathbf{I} - T_m\left(\frac{((\mu_p + \mu_1)\mathbf{I} - 2\mathcal{L})}{\mu_p - \mu_1}\right)}{T_m\left(\frac{\mu_p + \mu_1}{\mu_p - \mu_1}\right)} \right] \mathbf{b}_k^i, \quad i = 1, \dots, dk \quad (5)$$

satisfies the  $\epsilon$ -accuracy requirement:  $\|\mathbf{d}_s^{(m),i} - \mathbf{d}_s^{*,i}\|_{\mathcal{L}}^2 \leq \epsilon \|\mathbf{d}_s^{*,i}\|_{\mathcal{L}}^2$ .

Having proposed an approximate solution  $\mathbf{d}_s^{(m),i}$ , at this stage we are ready to commence with the distributed implementation of our solver. However, we recognize the following two challenges hindering its direct distributed implementation. First, we note that computing the minimum and maximum non-zero eigenvalues of  $\mathcal{L}$  requires global information. The second relates to the product with  $\mathcal{L}^\dagger$  needed in Equation 5. Here we detail the solutions to above two problems for the case when  $\mathcal{L}$  is graph Laplacian and derive our distributed SDD solver, which is used later to compute the Newton direction

1. **Parameters  $\mu_1$  and  $\mu_p$ .** As clear from the previous section, our method requires the computation of the second-minimum and maximum eigenvalues of  $\mathcal{L}$ . The computation of these, however, requires global information and hence are difficult to determine in a distributed fashion. As a substitute for the exact values of  $\mu_1$  and  $\mu_p$ , one can use the well-known eigenvalue bounds determined as

$$\begin{aligned} \mu_1 &\geq \underline{\mu} = \frac{4}{n^2} \\ \mu_p &\leq \bar{\mu} = 2n \end{aligned}$$

2. **Multiplication on  $\mathcal{L}^\dagger$ .** We start by noting that the second issue faced relates to the computational inefficiency when attempting to compute the coefficients of  $T_m\left(\frac{\mu_p + \mu_1}{\mu_p - \mu_1}\right) \mathbf{I} - T_m\left(\frac{((\mu_p + \mu_1)\mathbf{I} - 2\mathcal{L})}{\mu_p - \mu_1}\right)$  where performing it naively will potentially lead to linear dependency on the condition number of the processing graph. To illustrate, let us, in fact, consider the naive approach by assuming that each node  $i$  has access to the following decomposition of  $T_m(z)$ :

$$T_m(z) = 1 + \alpha_1 z + \alpha_2 z^2 + \dots + \alpha_m z^m$$

where  $\alpha_1, \dots, \alpha_m$  are coefficients one for each power of the polynomial. For ease of exposition, let us further denote

$$c_1 = \frac{\bar{\mu} + \mu}{\bar{\mu} - \mu} \quad c_2 = \frac{2}{\bar{\mu} - \mu}$$

Using the above, the numerator in Equation (5): can be written as

$$\begin{aligned} \mathcal{L}^\dagger [T_m(c_1)\mathbf{I} - T_m(c_1\mathbf{I} - c_2\mathcal{L})] \mathbf{b}_s^i &= \\ \mathcal{L}^\dagger \left[ \sum_{\nu=1}^m \alpha_\nu c_1^\nu \mathbf{I} - \sum_{\nu=1}^m \alpha_\nu (c_1\mathbf{I} - c_2\mathcal{L})^\nu \right] \mathbf{b}_s^i &= \\ \mathcal{L}^\dagger \left[ \sum_{\nu=1}^m \alpha_\nu [(c_1\mathbf{I})^\nu - (c_1\mathbf{I} - c_2\mathcal{L})^\nu] \right] \mathbf{b}_s^i \end{aligned}$$

The term  $c_1^\nu \mathbf{I}$  is easy to compute. The second, on the other hand, can be computed by rewriting the term  $(c_1\mathbf{I} - c_2\mathcal{L})^i$  explicitly in terms of  $\mathcal{L}$  for each node  $i$ . Unfortunately, this procedure is inefficient as it boils-down to a total of  $\mathcal{O}(m^2)$  matrix vector multiplications of the form  $\mathcal{L}\mathbf{u}$ . Taking into account the expression for  $m$ , we end up with an algorithm exhibiting linear dependency on the condition number  $\kappa(\mathcal{L})$ . Instead, our goal is to show that solution in (5) can be computed in fully distributed way in  $\mathcal{O}(m)$  rounds. The crucial property for us here is the recursive relation of Chebyshev polynomials:

$$\begin{aligned} T_0(z) &= 1, \\ T_1(z) &= z, \\ T_\ell(z) &= 2zT_{\ell-1}(z) - T_{\ell-2}(z) \end{aligned} \tag{6}$$

Denote by

$$\begin{aligned} \Delta_\ell &= \mathcal{L}^\dagger [T_\ell(c_1)\mathbf{I} - T_\ell(c_1\mathbf{I} - c_2\mathcal{L})] \mathbf{b}_s^i \\ \Omega_\ell &= T_\ell(c_1\mathbf{I} - c_2\mathcal{L}) \mathbf{b}_s^i \\ \Theta_\ell &= T_\ell(c_1) \end{aligned}$$

Therefore, the solution vector (5) can be written as  $\mathbf{d}_s^{(m),i} = \frac{\Delta_m}{\Theta_m}$  and recursive relation gives:

$$\begin{aligned} \Delta_\ell &= 2c_1\Delta_{\ell-1} - \Delta_{\ell-2} + 2c_2\Omega_{\ell-1} \\ \Omega_\ell &= 2(c_1\mathbf{I} - c_2\mathcal{L})\Omega_{\ell-1} - \Omega_{\ell-2} \\ \Theta_\ell &= 2c_1\Theta_{\ell-1} - \Theta_{\ell-2} \end{aligned} \tag{7}$$

with initials given by:

$$\begin{aligned} \Delta_1 &= c_2 \mathbf{b}_s^i & \Omega_1 &= [c_1\mathbf{I} - c_2\mathcal{L}] \mathbf{b}_s^i & \Theta_1 &= c_1 \\ \Delta_0 &= \mathbf{0} & \Omega_0 &= \mathbf{b}_k^i & \Theta_0 &= 1 \end{aligned}$$

Algorithm 2 summarizes these results and provides a fully distributed computation of vector (5) in  $\mathcal{O}(m)$  rounds. Clearly, lines 8-10 are executing relations (7) in a fully distributed way. Indeed, each matrix vector multiplication  $(c_1\mathbf{I} - c_2\mathcal{L})\mathbf{u}$  can be computed locally by a single message exchange between the neighboring nodes. Moreover, the total number of such multiplications is bounded by  $\mathcal{O}(m)$  and this fact establishes the following

**Theorem 1.** *The distributed SDD solver described in Algorithm 2 uses local communication exchange to compute an  $\epsilon$ -approximate solution of the SDD system (1) in the following number of rounds*

$$\mathcal{O} \left( \sqrt{\kappa(\mathcal{L})} \log \left( \frac{1}{\epsilon} \right) \right) \leq \mathcal{O} \left( \sqrt{nd_{\max} \text{diam}(\mathcal{G})} \log \left( \frac{1}{\epsilon} \right) \right)$$

where  $\kappa(\mathcal{L})$  is condition number of  $\mathcal{G}$ ,  $d_{\max}$ ,  $\text{diam}(\mathcal{G})$  are its maximal degree and diameter.

This result provides us with time complexity  $\mathcal{O} \left( d_{\max} \sqrt{\kappa(\mathcal{L})} \log \left( \frac{1}{\epsilon} \right) \right)$  and message complexity  $\mathcal{O} \left( |\mathcal{E}| \sqrt{\kappa(\mathcal{L})} \log \left( \frac{1}{\epsilon} \right) \right)$ .

---

**Algorithm 2 : Chebyshev SDD Solver**


---

- 1: **Input:** The  $r^{th}$  row of graph Laplacian  $\mathcal{L}$ , the  $r^{th}$  component of vector  $\mathbf{b}_k^i$ , precision parameter  $\epsilon$ .
  - 2: **Output:** The  $r^{th}$  components of  $\epsilon$ - approximate solution  $\mathbf{d}_s^{(m),i}$ .
  - 3: Set  $\bar{\mu} = 2n$ ,  $\underline{\mu} = \frac{4}{n^2}$  and  $c_1 = \frac{\bar{\mu} + \underline{\mu}}{\bar{\mu} - \underline{\mu}}$ ,  $c_2 = \frac{2}{\bar{\mu} - \underline{\mu}}$ ,  $\kappa(\mathcal{L}) = \frac{\bar{\mu}}{\underline{\mu}}$
  - 4:  $m = \lceil \frac{1}{2}(\sqrt{\kappa(\mathcal{L})} + 1) \ln \frac{2}{\epsilon} \rceil$ .
  - 5:  $[\Delta_0]_r = 0$   $[\Omega_0]_r = [\mathbf{b}_s^i]_r$   $\Theta_0 = 1$ .
  - 6:  $[\Delta_1]_r = c_2[\mathbf{b}_s^i]_r$   $[\Omega_1]_r = [(c_1\mathbf{I} - c_2\mathcal{L})\mathbf{b}_s^i]_r$   $\Theta_1 = c_1$ .
  - 7: **for**  $\ell = 2$  **to**  $m$  **do**
  - 8:  $\Theta_\ell = 2c_1\Theta_{\ell-1} - \Theta_{\ell-2}$ .
  - 9:  $[\Omega_\ell]_r = [2(c_1\mathbf{I} - c_2\mathcal{L})\Omega_{\ell-1}]_r - [\Omega_{\ell-2}]_r$ .
  - 10:  $[\Delta_\ell]_r = 2c_1[\Delta_{\ell-1}]_r - [\Delta_{\ell-2}]_r + 2c_2[\Omega_{\ell-1}]_r$
  - 11: **end for**
  - 12: Set  $[\mathbf{d}_s^{(m),i}]_r = \frac{[\Delta_m]_r}{\Theta_m}$
- 

## C.2 Convergence Analysis of Distributed Newton Method

Before to proceed to the prove of the Theorem ?? we will establish several intermediate results which are crucial for our convergence analysis. We pose the following assumptions on the local functions:

**Assumption 1.** The cost functions,  $f_r(\cdot) = \sum_{t=1}^{T_i} \mathcal{J}_{MTRL}^{(t)}(\cdot)$ , in Equation (??) are

1. twice continuously differentiable, i.e.,  $\gamma \mathbf{I}_{dk \times dk} \leq \nabla^2 f_r(\cdot) \leq \Gamma \mathbf{I}_{dk \times dk}$ , with  $\gamma$  and  $\Gamma$  are constants; and
2. Hessian Lipschitz continuous, i.e.,  $\|\nabla^2 f_r(\mathbf{x}) - \nabla^2 f_r(\hat{\mathbf{x}})\|_2 \leq \delta \|\mathbf{x} - \hat{\mathbf{x}}\|_2$  for all  $\mathbf{x}, \hat{\mathbf{x}} \in \mathbb{R}^{dk}$

### C.2.1 Primal Dual Transition

Recall, that transition between primal and dual variables is given by a system of differential equations (??) and let  $\phi_1^{(i)}, \dots, \phi_{dk}^{(i)}$  denote the solution of this system.

**Lemma** Let  $z_1 = [\mathcal{L}\lambda_1]_r, z_2 = [\mathcal{L}\lambda_2]_r, \dots, z_{dk} = [\mathcal{L}\lambda_{dk}]_r$ . Under Assumption 1, the functions  $\phi_1^{(r)}, \dots, \phi_{dk}^{(r)}$  exhibit bounded partial derivatives with respect to  $z_1, \dots, z_{dk}$ . In other words, for any  $r = 1, \dots, dk$ :

$$\left| \frac{\partial \phi_r^{(i)}}{\partial z_1} \right| \leq \frac{\sqrt{dk}}{\gamma}, \quad \dots, \quad \left| \frac{\partial \phi_r^{(i)}}{\partial z_{dk}} \right| \leq \frac{\sqrt{dk}}{\gamma}$$

for any  $[z_1, \dots, z_{dk}] \in \mathbb{R}^{dk}$ .

**Proof.** Using the definition of  $z_1, \dots, z_{dk}$  we can write:

$$\begin{cases} \frac{\partial f_i}{\partial \phi_1^{(i)}} = -z_1 \\ \frac{\partial f_i}{\partial \phi_2^{(i)}} = -z_2 \\ \vdots \\ \frac{\partial f_i}{\partial \phi_{dk}^{(i)}} = -z_{dk} \end{cases} \quad (8)$$

Taking the derivative with respect to  $z_1$  in each equation of system (8) gives:

$$\begin{cases} \frac{\partial^2 f_i}{\partial (\phi_1^{(i)})^2} \frac{\partial \phi_1^{(i)}}{\partial z_1} + \frac{\partial^2 f_i}{\partial \phi_1^{(i)} \partial \phi_2^{(i)}} \frac{\partial \phi_2^{(i)}}{\partial z_1} + \dots + \frac{\partial^2 f_i}{\partial \phi_1^{(i)} \partial \phi_{dk}^{(i)}} \frac{\partial \phi_{dk}^{(i)}}{\partial z_1} = -1 \\ \frac{\partial^2 f_i}{\partial \phi_2^{(i)} \partial \phi_1^{(i)}} \frac{\partial \phi_1^{(i)}}{\partial z_1} + \frac{\partial^2 f_i}{\partial (\phi_2^{(i)})^2} \frac{\partial \phi_2^{(i)}}{\partial z_1} + \dots + \frac{\partial^2 f_i}{\partial \phi_2^{(i)} \partial \phi_{dk}^{(i)}} \frac{\partial \phi_{dk}^{(i)}}{\partial z_1} = 0 \\ \vdots \\ \frac{\partial^2 f_i}{\partial \phi_{dk}^{(i)} \partial \phi_1^{(i)}} \frac{\partial \phi_1^{(i)}}{\partial z_1} + \frac{\partial^2 f_i}{\partial \phi_p^{(i)} \partial \phi_2^{(i)}} \frac{\partial \phi_2^{(i)}}{\partial z_1} + \dots + \frac{\partial^2 f_i}{\partial (\phi_{dk}^{(i)})^2} \frac{\partial \phi_{dk}^{(i)}}{\partial z_1} = 0 \end{cases}$$

Let  $\mathbf{u}_1 = [\frac{\partial \phi_1^{(i)}}{\partial z_1}, \frac{\partial \phi_2^{(i)}}{\partial z_1}, \dots, \frac{\partial \phi_{dk}^{(i)}}{\partial z_1}]^\top$  then the above result can be written in matrix vector form:

$$[\nabla^2 f_i] \mathbf{u}_1 = -\mathbf{e}_1$$

where  $\mathbf{e}_1 = [1, 0 \dots, 0] \in \mathbb{R}^{dk}$ . Similarly we have:

$$[\nabla^2 f_i] \mathbf{u}_2 = -\mathbf{e}_2 \quad [\nabla^2 f_i] \mathbf{u}_3 = -\mathbf{e}_3, \dots \quad [\nabla^2 f_i] \mathbf{u}_{dk} = -\mathbf{e}_{dk}$$

with  $\mathbf{u}_r = [\frac{\partial \phi_1^{(i)}}{\partial z_r}, \frac{\partial \phi_2^{(i)}}{\partial z_r}, \dots, \frac{\partial \phi_{dk}^{(i)}}{\partial z_r}]^\top$ . Combining all these equations gives:

$$[\nabla^2 f_i] \mathbf{U} = -\mathbf{I}_{dk \times dk} \tag{9}$$

where

$$\mathbf{U} = \begin{bmatrix} \frac{\partial \phi_1^{(i)}}{\partial z_1} & \frac{\partial \phi_1^{(i)}}{\partial z_2} & \dots & \frac{\partial \phi_1^{(i)}}{\partial z_{dk}} \\ \frac{\partial \phi_2^{(i)}}{\partial z_1} & \frac{\partial \phi_2^{(i)}}{\partial z_2} & \dots & \frac{\partial \phi_2^{(i)}}{\partial z_{dk}} \\ \vdots & & \ddots & \vdots \\ \frac{\partial \phi_{dk}^{(i)}}{\partial z_1} & \frac{\partial \phi_{dk}^{(i)}}{\partial z_2} & \dots & \frac{\partial \phi_{dk}^{(i)}}{\partial z_{dk}} \end{bmatrix}$$

Notice, Equation (9) implies:

$$\mathbf{U} = -[\nabla^2 f_i]^{-1}$$

Hence, using Assumption 1:  $\|\mathbf{U}\|_2 \leq \frac{1}{\gamma}$ , and:

$$|U_{ij}| \leq \|\mathbf{U}\|_F \leq \sqrt{dk} \|\mathbf{U}\|_2 \leq \frac{\sqrt{dk}}{\gamma}$$

### C.2.2 Dual Function Properties

In this section we establish important properties of the dual to the problem (??).

**Lemma 1.** *The function  $q(\boldsymbol{\lambda}) = q(\lambda_1, \dots, \lambda_{dk})$  abides by the following properties:*

1. *Let  $\mathbf{y}(\boldsymbol{\lambda})$  be the primal variable corresponding to dual vector  $\boldsymbol{\lambda}$ . Then the gradient and the Hessian of  $q(\boldsymbol{\lambda})$  are given by*

$$\nabla q(\boldsymbol{\lambda}) = \mathbf{g}(\boldsymbol{\lambda}) = \mathbf{M} \mathbf{y}(\boldsymbol{\lambda})$$

$$\nabla^2 q(\boldsymbol{\lambda}) = \mathbf{H}(\boldsymbol{\lambda}) = -\mathbf{M} [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1} \mathbf{M}$$

where  $f(\mathbf{y}(\boldsymbol{\lambda})) = \sum_{i=1}^n \sum_{t=1}^{T_i} \mathcal{J}_{MTRL}^{(t)}(\mathbf{y}(\boldsymbol{\lambda}))$ .

2. *Denote  $\mu_n(\mathcal{L})$  as the largest eigenvalue of the unweighted Laplacian of  $\mathbb{G}$  and constants  $\delta, \gamma$  are given in Assumption 1. Then, for constant  $B = dk\delta \left(\frac{\mu_n(\mathcal{L})}{\gamma}\right)^3$  and for any  $\bar{\boldsymbol{\lambda}}, \boldsymbol{\lambda} \in \mathbb{R}^{ndk}$ :*

$$\|\mathbf{H}(\bar{\boldsymbol{\lambda}}) - \mathbf{H}(\boldsymbol{\lambda})\|_2 \leq B \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2$$

*i.e.  $\mathbf{H}(\boldsymbol{\lambda})$  is Lipschitz continuous with constant  $B$ .*

**Proof:** To avoid confusion with indexes let us focus on more general optimization problem:

$$\min_{\mathbf{x}} f(\mathbf{x}) \tag{10}$$

$$s.t. \quad \mathbf{A} \mathbf{x} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{n \times p}, \quad \mathbf{b} \in \mathbb{R}^n$$

where  $f(\mathbf{x})$  twice differentiable strongly convex function, and unknown variable  $\mathbf{x} \in \mathbb{R}^p$ . One can see, that problem (??) is a special case of (10) with  $f(\cdot) = \sum_{i=1}^n f_i(\cdot)$ ,  $\mathbf{A} = \mathbf{M}$ , and  $\mathbf{b} = \mathbf{0}^1$ . Let  $q(\boldsymbol{\lambda})$  be the corresponding dual for (10), with dual variable  $\boldsymbol{\lambda} \in \mathbb{R}^n$ . We will show that:

$$\nabla^2 q(\boldsymbol{\lambda}) = -\mathbf{A} [\nabla^2 f(\mathbf{x}(\boldsymbol{\lambda}))]^{-1} \mathbf{A}^\top \tag{11}$$

$$\nabla q(\boldsymbol{\lambda}) = \mathbf{A} \mathbf{x}(\boldsymbol{\lambda}) - \mathbf{b}$$

---

<sup>1</sup>In this case  $n = p = dk$

where  $\mathbf{x}(\boldsymbol{\lambda}) = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) + \boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{x} - \mathbf{b})$  minimizes the Lagrangian of problem (10). Let us denote  $\mathbf{x}(\boldsymbol{\lambda}) = \mathbf{x}^+$

$$\mathbf{A} = \begin{bmatrix} a_{11} & \cdots & a_{1p} \\ a_{21} & \cdots & a_{2p} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{np} \end{bmatrix}, \quad \mathbf{x}^+ = \begin{bmatrix} x_1^+(\boldsymbol{\lambda}) \\ x_2^+(\boldsymbol{\lambda}) \\ \vdots \\ x_p^+(\boldsymbol{\lambda}) \end{bmatrix}, \quad \nabla f(\mathbf{x}^+) = \begin{bmatrix} z_1(\mathbf{x}^+) \\ z_2(\mathbf{x}^+) \\ \vdots \\ z_p(\mathbf{x}^+) \end{bmatrix} \quad (12)$$

The optimal primal variable  $\mathbf{x}(\boldsymbol{\lambda})$  satisfies:

$$\nabla f(\mathbf{x}(\boldsymbol{\lambda})) + \mathbf{A}^\top \boldsymbol{\lambda} = \mathbf{0}. \quad (13)$$

Using Fenchel's conjugate, the dual function can be written as:

$$q(\boldsymbol{\lambda}) = -\mathbf{b}^\top \boldsymbol{\lambda} - f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) \quad (14)$$

Therefore,

$$\nabla q(\boldsymbol{\lambda}) = -\mathbf{b} - \nabla f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) \quad (15)$$

Denoting  $\mathbf{u} = -\mathbf{A}^\top \boldsymbol{\lambda}$ , then the  $k^{th}$  component of vector  $\nabla f^*(-\mathbf{A}^\top \boldsymbol{\lambda})$  can be written as:

$$[\nabla f^*(-\mathbf{A}^\top \boldsymbol{\lambda})]_s = \sum_{j=1}^p \frac{\partial f^*}{\partial u_j} \frac{\partial u_j}{\partial \lambda_k} = - [a_{k1} \quad a_{k2} \quad \cdots \quad a_{kp}] \times \begin{bmatrix} \frac{\partial f^*}{\partial u_1} \\ \frac{\partial f^*}{\partial u_2} \\ \vdots \\ \frac{\partial f^*}{\partial u_p} \end{bmatrix} \Big|_{-\mathbf{A}^\top \boldsymbol{\lambda}}$$

Applying result (13) and the relation between the gradients of function and its Fenchel's conjugate:

$$\begin{aligned} \nabla f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) &= -\mathbf{A} \nabla_{\mathbf{u}} f^*(\mathbf{u})|_{-\mathbf{A}^\top \boldsymbol{\lambda}} = -\mathbf{A} \nabla_{\mathbf{u}} f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) = \\ &= -\mathbf{A} \nabla_{\mathbf{u}} f^*(\nabla f(\mathbf{x}^+)) = -\mathbf{A} \mathbf{x}(\boldsymbol{\lambda}) \end{aligned} \quad (16)$$

Therefore, the result (15) gives:

$$\nabla q(\boldsymbol{\lambda}) = -\mathbf{b} + \mathbf{A} \mathbf{x}(\boldsymbol{\lambda}) \quad (17)$$

which establishes the claim for the dual gradient.

Taking the gradient in (17) gives:

$$\nabla^2 q(\boldsymbol{\lambda}) = \mathbf{A} \underbrace{\begin{bmatrix} \frac{\partial x_1^+(\boldsymbol{\lambda})}{\partial \lambda_1} & \frac{\partial x_1^+(\boldsymbol{\lambda})}{\partial \lambda_2} & \cdots & \frac{\partial x_1^+(\boldsymbol{\lambda})}{\partial \lambda_n} \\ \frac{\partial x_2^+(\boldsymbol{\lambda})}{\partial \lambda_1} & \frac{\partial x_2^+(\boldsymbol{\lambda})}{\partial \lambda_2} & \cdots & \frac{\partial x_2^+(\boldsymbol{\lambda})}{\partial \lambda_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_p^+(\boldsymbol{\lambda})}{\partial \lambda_1} & \frac{\partial x_p^+(\boldsymbol{\lambda})}{\partial \lambda_2} & \cdots & \frac{\partial x_p^+(\boldsymbol{\lambda})}{\partial \lambda_n} \end{bmatrix}}_{\mathbf{F}(\mathbf{x}^+)} \quad (18)$$

Hence, we target matrix  $\mathbf{F}(\mathbf{x}^+)$  to obtain the form of dual Hessian. Equation (13) gives:

$$\nabla f(\mathbf{x}^+) = -\mathbf{A}^\top \boldsymbol{\lambda}$$

On the next step, we take partial derivative  $\frac{\partial}{\partial \lambda_j}$  for both sides of the above equation for  $j = 1, \dots, n$ . For simplicity, consider  $\frac{\partial}{\partial \lambda_1}$ :

$$\begin{aligned} \frac{\partial}{\partial \lambda_1} \nabla f(\mathbf{x}^+) &= \begin{bmatrix} \frac{\partial}{\partial \lambda_1} z_1(\mathbf{x}^+) \\ \frac{\partial}{\partial \lambda_1} z_2(\mathbf{x}^+) \\ \vdots \\ \frac{\partial}{\partial \lambda_1} z_p(\mathbf{x}^+) \end{bmatrix} = \\ &= \begin{bmatrix} \frac{\partial z_1(\mathbf{x}^+)}{\partial x_1^+(\lambda)} \frac{\partial x_1^+(\lambda)}{\partial \lambda_1} + \frac{\partial z_1(\mathbf{x}^+)}{\partial x_2^+(\lambda)} \frac{\partial x_2^+(\lambda)}{\partial \lambda_1} + \dots + \frac{\partial z_1(\mathbf{x}^+)}{\partial x_p^+(\lambda)} \frac{\partial x_p^+(\lambda)}{\partial \lambda_1} \\ \frac{\partial z_2(\mathbf{x}^+)}{\partial x_1^+(\lambda)} \frac{\partial x_1^+(\lambda)}{\partial \lambda_1} + \frac{\partial z_2(\mathbf{x}^+)}{\partial x_2^+(\lambda)} \frac{\partial x_2^+(\lambda)}{\partial \lambda_1} + \dots + \frac{\partial z_2(\mathbf{x}^+)}{\partial x_p^+(\lambda)} \frac{\partial x_p^+(\lambda)}{\partial \lambda_1} \\ \vdots \\ \frac{\partial z_p(\mathbf{x}^+)}{\partial x_1^+(\lambda)} \frac{\partial x_1^+(\lambda)}{\partial \lambda_1} + \frac{\partial z_p(\mathbf{x}^+)}{\partial x_2^+(\lambda)} \frac{\partial x_2^+(\lambda)}{\partial \lambda_1} + \dots + \frac{\partial z_p(\mathbf{x}^+)}{\partial x_p^+(\lambda)} \frac{\partial x_p^+(\lambda)}{\partial \lambda_1} \end{bmatrix} = \\ &= \underbrace{\begin{bmatrix} \frac{\partial z_1(\mathbf{x}^+)}{\partial x_1^+(\lambda)} & \frac{\partial z_1(\mathbf{x}^+)}{\partial x_2^+(\lambda)} & \dots & \frac{\partial z_1(\mathbf{x}^+)}{\partial x_p^+(\lambda)} \\ \frac{\partial z_2(\mathbf{x}^+)}{\partial x_1^+(\lambda)} & \frac{\partial z_2(\mathbf{x}^+)}{\partial x_2^+(\lambda)} & \dots & \frac{\partial z_2(\mathbf{x}^+)}{\partial x_p^+(\lambda)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_p(\mathbf{x}^+)}{\partial x_1^+(\lambda)} & \frac{\partial z_p(\mathbf{x}^+)}{\partial x_2^+(\lambda)} & \dots & \frac{\partial z_p(\mathbf{x}^+)}{\partial x_p^+(\lambda)} \end{bmatrix}}_{\nabla^2 f(\mathbf{x}^+)} \begin{bmatrix} \frac{\partial x_1^+(\lambda)}{\partial \lambda_1} \\ \frac{\partial x_2^+(\lambda)}{\partial \lambda_1} \\ \vdots \\ \frac{\partial x_p^+(\lambda)}{\partial \lambda_1} \end{bmatrix} = \nabla^2 f(\mathbf{x}^+) \begin{bmatrix} \frac{\partial x_1^+(\lambda)}{\partial \lambda_1} \\ \frac{\partial x_2^+(\lambda)}{\partial \lambda_1} \\ \vdots \\ \frac{\partial x_p^+(\lambda)}{\partial \lambda_1} \end{bmatrix} = \\ &= \frac{\partial}{\partial \lambda_1} (-\mathbf{A}^\top \boldsymbol{\lambda}) = - \begin{bmatrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1p} \end{bmatrix} \end{aligned}$$

Repeating this for  $\frac{\partial}{\partial \lambda_2}, \dots, \frac{\partial}{\partial \lambda_p}$  gives:

$$\nabla^2 f(\mathbf{x}^+) \underbrace{\begin{bmatrix} \frac{\partial x_1^+(\lambda)}{\partial \lambda_1} & \frac{\partial x_1^+(\lambda)}{\partial \lambda_2} & \dots & \frac{\partial x_1^+(\lambda)}{\partial \lambda_n} \\ \frac{\partial x_2^+(\lambda)}{\partial \lambda_1} & \frac{\partial x_2^+(\lambda)}{\partial \lambda_2} & \dots & \frac{\partial x_2^+(\lambda)}{\partial \lambda_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_p^+(\lambda)}{\partial \lambda_1} & \frac{\partial x_p^+(\lambda)}{\partial \lambda_2} & \dots & \frac{\partial x_p^+(\lambda)}{\partial \lambda_n} \end{bmatrix}}_{\mathbf{F}(\mathbf{x}^+)} = - \underbrace{\begin{bmatrix} a_{11} & a_{21} & \dots & a_{n1} \\ a_{12} & a_{22} & \dots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1p} & a_{2p} & \dots & a_{np} \end{bmatrix}}_{\mathbf{A}^\top}$$

Therefore, for matrix  $\mathbf{F}(\mathbf{x}^+) = -[\nabla^2 f(\mathbf{x}^+)]^{-1} \mathbf{A}^\top$  and combining this result with (18) gives:

$$\nabla^2 q(\boldsymbol{\lambda}) = -\mathbf{A}[\nabla^2 f(\mathbf{x}^+)]^{-1} \mathbf{A}^\top$$

Now we are ready to prove the second statement of the lemma. Using  $\mathbf{M} \preceq \mu_n(\mathcal{L})\mathbf{I}$ :

$$\begin{aligned} \|[\mathbf{H}(\bar{\boldsymbol{\lambda}}) - \mathbf{H}(\boldsymbol{\lambda})]\mathbf{v}\|_2^2 &= \|[\mathbf{M}([\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1})\mathbf{M}\mathbf{v}]\|_2^2 = \\ &= \mathbf{v}^\top \mathbf{M}([\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}) \mathbf{M}^2([\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}) \mathbf{M}\mathbf{v} \leq \\ &= \mu_n^2(\mathcal{L}) \mathbf{v}^\top \mathbf{M}([\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1})^2 \mathbf{M}\mathbf{v} \leq \\ &= \mu_n^2(\mathcal{L}) \mu_{\max}^2(|[\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}|) \mathbf{v}^\top \mathbf{M}^2 \mathbf{v} \leq \\ &= \mu_n^4(\mathcal{L}) \mu_{\max}^2(|[\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}|) \|\mathbf{v}\|_2^2 \end{aligned}$$

Therefore,

$$\|\mathbf{H}(\bar{\boldsymbol{\lambda}}) - \mathbf{H}(\boldsymbol{\lambda})\|_2 \leq \mu_n^2(\mathcal{L}) \mu_{\max}(|[\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}|) \quad (19)$$

To bound the term  $\mu_{\max}(|[\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}|)$  we study the properties of primal Hessian more carefully:



**Claim 2.** For primal Hessian  $\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))$  the following properties are true

$$\gamma \preceq \nabla^2 f(\mathbf{y}(\boldsymbol{\lambda})) \preceq \Gamma \quad (20)$$

$$\mu_{\max}(|[\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}|) \leq \quad (21)$$

$$\delta \max_{i \in \mathbb{V}} \sqrt{\sum_{k=1}^p ([\mathbf{y}_s]_i(\bar{\boldsymbol{\lambda}}) - [\mathbf{y}_s]_i(\boldsymbol{\lambda}))^2}$$

for any  $\bar{\boldsymbol{\lambda}}, \boldsymbol{\lambda} \in \mathbb{R}^p$ .

**Proof.** Firstly, notice that for any  $j \neq i$  and any  $r = 1 \dots, p$ :

$$\frac{\partial^2 f}{\partial [\mathbf{y}_1]_i \partial [\mathbf{y}_r]_j} = \frac{\partial^2 f}{\partial [\mathbf{y}_2]_i \partial [\mathbf{y}_r]_j} = \dots = \frac{\partial^2 f}{\partial [\mathbf{y}_p]_i \partial [\mathbf{y}_r]_j} = 0$$

Hence, the sparsity pattern of primal Hessian allows the symmetric reordering of rows and columns such that  $\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))$  is transformed into the block diagonal matrix:

$$\mathbf{W}(\boldsymbol{\lambda}) = \begin{bmatrix} \nabla^2 f_1(\boldsymbol{\lambda}) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \nabla^2 f_2(\boldsymbol{\lambda}) & \dots & \mathbf{0} \\ \vdots & & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \nabla^2 f_p(\boldsymbol{\lambda}) \end{bmatrix}$$

The matrix  $\mathbf{W}(\boldsymbol{\lambda})$  preserves the important properties of  $\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))$ . Particularly, the spectrum of these two matrices are the same. Indeed, let  $\mathbf{T}_{ij}$  is the operator that swaps  $i^{th}$  and  $j^{th}$  rows of some arbitrary matrix  $\mathbf{A}$  and let  $\bar{\mathbf{A}}$  be the result of such transformation. Then,  $\bar{\mathbf{A}} = \mathbf{T}_{ij} \mathbf{A} \mathbf{T}_{ij}$ , and using  $\mathbf{T}_{ij}^2 = \mathbf{I}$ :

$$\begin{aligned} \det(\bar{\mathbf{A}} - \mu \mathbf{I}) &= \det(\mathbf{T}_{ij} \mathbf{A} \mathbf{T}_{ij} - \mu \mathbf{I}) = \det(\mathbf{T}_{ij} (\mathbf{A} - \mu \mathbf{I}) \mathbf{T}_{ij}) = \\ &= \det(\mathbf{A} - \mu \mathbf{I}) \det(\mathbf{T}_{ij}^2) = \det(\mathbf{A} - \mu \mathbf{I}) \end{aligned}$$

Since  $\mathbf{W}(\boldsymbol{\lambda})$  is constructed from  $\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))$  by symmetric reordering rows and columns, then  $\text{Spectrum}(\mathbf{W}(\boldsymbol{\lambda})) = \text{Spectrum}(\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda})))$ . Using the Assumption 1 it implies:

$$\gamma \preceq \mathbf{W}(\boldsymbol{\lambda}) \preceq \Gamma$$

To prove (21), notice that if  $\bar{\mathbf{A}} = \mathbf{T}_{ij} \mathbf{A} \mathbf{T}_{ij}$  and  $\mathbf{A}$  is invertible, then so  $\bar{\mathbf{A}}$  and using  $\mathbf{T}_{ij}^{-1} = \mathbf{T}_{ij}$ :

$$\begin{aligned} \det(\bar{\mathbf{A}}^{-1} - \mu \mathbf{I}) &= \det(\mathbf{T}_{ij}^{-1} \mathbf{A}^{-1} \mathbf{T}_{ij}^{-1} - \mu \mathbf{I}) = \det(\mathbf{T}_{ij} (\mathbf{A}^{-1} - \mu \mathbf{I}) \mathbf{T}_{ij}) = \\ &= \det(\mathbf{A}^{-1} - \mu \mathbf{I}) \end{aligned}$$

Denote  $\{\mathbf{T}_1, \dots, \mathbf{T}_l\}$  is a collection of operators that swaps the rows of matrix  $\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))$  to transform it to  $\mathbf{W}(\boldsymbol{\lambda})$ , i.e.

$$\mathbf{W}(\boldsymbol{\lambda}) = \mathbf{T}_1 \dots \mathbf{T}_l \nabla^2 f(\mathbf{y}(\boldsymbol{\lambda})) \mathbf{T}_l \dots \mathbf{T}_1$$

Then  $[\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1} = \mathbf{T}_l \dots \mathbf{T}_1 \mathbf{W}^{-1}(\boldsymbol{\lambda}) \mathbf{T}_1 \dots \mathbf{T}_l$ , and using the Assumption 1:

$$\begin{aligned} \mu_{\max}(|[\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}|) &= \\ \mu_{\max}(\mathbf{T}_l \dots \mathbf{T}_1 |\mathbf{W}^{-1}(\bar{\boldsymbol{\lambda}}) - \mathbf{W}^{-1}(\boldsymbol{\lambda})| \mathbf{T}_1 \dots \mathbf{T}_l) &\leq \mu_{\max}(|\mathbf{W}^{-1}(\bar{\boldsymbol{\lambda}}) - \mathbf{W}^{-1}(\boldsymbol{\lambda})|) \leq \\ \max_{i \in \mathbb{V}} \mu_{\max}(|[\nabla^2 f_i([\mathbf{y}_1]_i(\bar{\boldsymbol{\lambda}}), \dots, [\mathbf{y}_p]_i(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f_i([\mathbf{y}_1]_i(\boldsymbol{\lambda}), \dots, [\mathbf{y}_p]_i(\boldsymbol{\lambda}))]^{-1}|) &= \\ \max_{i \in \mathbb{V}} |||[\nabla^2 f_i([\mathbf{y}_1]_i(\bar{\boldsymbol{\lambda}}), \dots, [\mathbf{y}_p]_i(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f_i([\mathbf{y}_1]_i(\boldsymbol{\lambda}), \dots, [\mathbf{y}_p]_i(\boldsymbol{\lambda}))]^{-1}|||_2 &\leq \\ \frac{\delta}{\gamma^2} \max_{i \in \mathbb{V}} ||([\mathbf{y}_1]_i(\bar{\boldsymbol{\lambda}}), \dots, [\mathbf{y}_p]_i(\bar{\boldsymbol{\lambda}})) - ([\mathbf{y}_1]_i(\boldsymbol{\lambda}), \dots, [\mathbf{y}_p]_i(\boldsymbol{\lambda}))||_2 &= \\ \frac{\delta}{\gamma^2} \max_{i \in \mathbb{V}} \sqrt{\sum_{k=1}^p ([\mathbf{y}_s]_i(\bar{\boldsymbol{\lambda}}) - [\mathbf{y}_s]_i(\boldsymbol{\lambda}))^2} \end{aligned}$$

which establishes (21).

Consider the term  $([\mathbf{y}_s]_i(\bar{\boldsymbol{\lambda}}) - [\mathbf{y}_s]_i(\boldsymbol{\lambda}))$ . Using the above results we can write:

$$\begin{aligned}
|y_k(i)(\bar{\boldsymbol{\lambda}}) - y_k(i)(\boldsymbol{\lambda})| &= |\phi_s^{(i)}([\mathcal{L}\bar{\boldsymbol{\lambda}}_1]_i, \dots, [\mathcal{L}\bar{\boldsymbol{\lambda}}_p]_i) - \phi_s^{(i)}([\mathcal{L}\boldsymbol{\lambda}_1]_i, \dots, [\mathcal{L}\boldsymbol{\lambda}_p]_i)| \leq \\
&\frac{\sqrt{p}}{\gamma} \sqrt{\sum_{r=1}^p ([\mathcal{L}\bar{\boldsymbol{\lambda}}_r]_i - [\mathcal{L}\boldsymbol{\lambda}_r]_i)^2} = \frac{\sqrt{p}}{\gamma} \sqrt{\sum_{r=1}^p [\mathbf{L}_{\mathcal{G}}(\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r)]_i^2} \leq \frac{\sqrt{p}}{\gamma} \sqrt{\sum_{r=1}^p \|\mathcal{L}(\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r)\|_2^2} = \\
&\frac{\sqrt{p}}{\gamma} \sqrt{\sum_{r=1}^p (\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r)^\top \mathcal{L}^2 (\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r)} \leq \frac{\sqrt{p}}{\gamma} \sqrt{\mu_n^2(\mathcal{L}) \sum_{r=1}^p (\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r)^\top (\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r)} = \\
&= \mu_n(\mathcal{L}) \frac{\sqrt{p}}{\gamma} \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2
\end{aligned}$$

where

$$\sum_{r=1}^p (\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r)^\top (\bar{\boldsymbol{\lambda}}_r - \boldsymbol{\lambda}_r) = \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2^2$$

is used. Hence,

$$([\mathbf{y}_s]_i(\bar{\boldsymbol{\lambda}}) - [\mathbf{y}_s]_i(\boldsymbol{\lambda}))^2 \leq \mu_n^2(\mathcal{L}) \frac{p}{\gamma^2} \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2^2$$

Combining this result with (21) gives:

$$\mu_{\max}(\|[\nabla^2 f(\mathbf{y}(\bar{\boldsymbol{\lambda}}))]^{-1} - [\nabla^2 f(\mathbf{y}(\boldsymbol{\lambda}))]^{-1}\|) \leq \frac{\delta}{\gamma^2} \mu_n(\mathcal{L}) \frac{p}{\gamma} \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2$$

and applying it to (19) gives:

$$\|\mathbf{H}(\bar{\boldsymbol{\lambda}}) - \mathbf{H}(\boldsymbol{\lambda})\|_2 \leq \delta p \left( \frac{\mu_n(\mathcal{L})}{\gamma} \right)^3 \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2 = \delta dk \left( \frac{\mu_n(\mathcal{L})}{\gamma} \right)^3 \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2 = B \|\bar{\boldsymbol{\lambda}} - \boldsymbol{\lambda}\|_2$$

In other words, dual Hessian is Lipschitz continuous with constant  $B = \delta dk \left( \frac{\mu_n(\mathcal{L})}{\gamma} \right)^3$

### C.2.3 Dual Gradient Bounds

The following Lemma studies the change of the norm of dual gradient for Distributed Newton iteration scheme and plays a crucial role for the convergence analysis:

**Lemma 2.** *Let us consider iteration scheme given by  $\boldsymbol{\lambda}_{s+1} = \boldsymbol{\lambda}_s + \alpha_s \mathbf{d}_s^{(m)}$  and denote*

$$\boldsymbol{\epsilon}_s = \mathbf{H}_s \mathbf{d}_s^{(m)} + \mathbf{g}_s$$

*be the approximation error vector corresponding to  $\epsilon$ -approximated Newton direction vector  $\mathbf{d}_s^{(m)}$  and  $\mathbf{g}_s = \mathbf{g}(\boldsymbol{\lambda}_s) = \nabla q(\boldsymbol{\lambda}_s)$ . Then for any  $\alpha_k \in (0, 1]$*

$$\|\mathbf{g}_{s+1}\|_2 \leq \tag{22}$$

$$(1 - \alpha_k) \|\mathbf{g}_s\|_2 + \alpha_k^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_s\|_2^2 + \alpha_k \|\boldsymbol{\epsilon}_s\|_2 + \alpha_s^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\boldsymbol{\epsilon}_s\|_2^2$$

where  $B$  is defined in Lemma 1 and  $\mu_2(\mathcal{L})$  is the smallest nonzero eigenvalue of unweighted Laplacian of  $\mathcal{G}$ .

**Proof.** Using definition of  $\boldsymbol{\epsilon}_s$  for the dual gradient we have:

$$\begin{aligned}
\mathbf{g}(\boldsymbol{\lambda}_s + \alpha_s \mathbf{d}_s^{(m)}) &= \mathbf{g}(\boldsymbol{\lambda}_s) + \int_0^1 \mathbf{H}(\boldsymbol{\lambda}_s + t\alpha_s \mathbf{d}_s^{(m)}) \alpha_s \mathbf{d}_s^{(m)} dt = \\
&\mathbf{g}(\boldsymbol{\lambda}_s) + \int_0^1 [\mathbf{H}(\boldsymbol{\lambda}_s + t\alpha_s \mathbf{d}_s^{(m)}) - \mathbf{H}(\boldsymbol{\lambda}_s)] \alpha_s \mathbf{d}_s^{(m)} dt + \alpha_s \int_0^1 \mathbf{H}(\boldsymbol{\lambda}_s) \mathbf{d}_s^{(m)} dt = \\
&\mathbf{g}(\boldsymbol{\lambda}_s) + \int_0^1 [\mathbf{H}(\boldsymbol{\lambda}_s + t\alpha_s \mathbf{d}_s^{(m)}) - \mathbf{H}(\boldsymbol{\lambda}_s)] \alpha_s \mathbf{d}_s^{(m)} dt + \alpha_s (\boldsymbol{\epsilon}_s - \mathbf{g}(\boldsymbol{\lambda}_s))
\end{aligned}$$

Applying  $\mathbf{g}_{s+1} = \mathbf{g}(\boldsymbol{\lambda}_s + \alpha_s \mathbf{d}_s^{(m)})$ ,  $\mathbf{g}_s = \mathbf{g}(\boldsymbol{\lambda}_s)$  and Lemma 1:

$$\begin{aligned} \|\mathbf{g}_{s+1}\|_2 &\leq (1 - \alpha_s)\|\mathbf{g}_s\|_2 + \alpha_s\|\boldsymbol{\epsilon}_s\|_2 + \frac{1}{2}\alpha_s^2 B \|\mathbf{d}_s^{(m)}\|_2^2 = \\ &(1 - \alpha_s)\|\mathbf{g}_s\|_2 + \alpha_s\|\boldsymbol{\epsilon}_s\|_2 + \frac{1}{2}\alpha_s^2 B \|\mathbf{H}^\dagger(\boldsymbol{\lambda}_s)(\mathbf{g}_s - \boldsymbol{\epsilon}_s)\|_2^2 \leq \\ &(1 - \alpha_s)\|\mathbf{g}_s\|_2 + \alpha_s\|\boldsymbol{\epsilon}_s\|_2 + \alpha_s^2 B \|\mathbf{H}^\dagger(\boldsymbol{\lambda}_s)\|_2^2 (\|\mathbf{g}_s\|_2^2 + \|\boldsymbol{\epsilon}_s\|_2^2) \end{aligned}$$

Investigating the explicit form of dual Hessian gives  $\|\mathbf{H}^\dagger(\boldsymbol{\lambda}_s)\|_2 \leq \frac{\Gamma}{\mu_2^2(\mathcal{L})}$ . Hence,

$$\|\mathbf{g}_{s+1}\|_2 \leq (1 - \alpha_s)\|\mathbf{g}_s\|_2 + \alpha_s^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_s\|_2^2 + \alpha_s\|\boldsymbol{\epsilon}_s\|_2 + \alpha_s^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\boldsymbol{\epsilon}_s\|_2^2.$$

### C.2.4 Newton Method Proofs

Similar to centralized Newton method, the step size  $\alpha_s$  in iteration scheme  $\boldsymbol{\lambda}_{s+1} = \boldsymbol{\lambda}_s + \alpha_s \mathbf{d}_s^{(m)}$  should be chosen carefully in order to attain quadratic convergence. In this part of the Appendix, we consider the distributed version of Armijo rule is given in Algorithm 3. We use  $\mathbf{g} = \mathbf{M}\mathbf{y} =$

---

#### Algorithm 3 : Distributed Line Search

---

**Input:** The constants  $\sigma \in (0, \frac{1}{2}]$  and  $\beta \in (0, 1)$ , parameters  $\epsilon, \Gamma, \gamma, \delta$ . The  $i^{th}$  component of dual gradient chunks:  $\{[\mathcal{L}\mathbf{y}_r]_i\}_{r=1}^{dk}$   
**Output:** step size  $\alpha_s$ .  
Set  $m_i = 0$ .  
Compute  $\eta_i = \max_r \{[\mathcal{L}\mathbf{y}_r]_i\}$ .  
Compute  $\max_i \{\eta_i\}$  using maximal consensus protocol.  
**while**  $\max_r \{[\mathcal{L}\mathbf{y}_r]_i\} > (1 - \sigma\beta^{m_i})\sqrt{n} \max_i \{\eta_i\} + 2\epsilon \frac{n\gamma^2}{dk\delta\Gamma}$  **do**  
     $m_i = m_i + 1$ .  
**end while**  
Compute  $\hat{m} = \max_i \{m_i\}$  using maximal consensus protocol.  
Set  $\alpha_s = \beta^{\hat{m}}$ .

---

$((\mathcal{L}\mathbf{y}_1)^\top, \dots, (\mathcal{L}\mathbf{y}_{dk})^\top)^\top$ . Algorithm 3 requires only  $\mathcal{O}(\text{diam}(\mathcal{G}))$  time steps and conducts only exact computations. The following Lemma studies the change of step size given by the proposed backtracking line search procedure:

**Lemma 3.** *Let step size  $\alpha_s$  is chosen according to Algorithm 3 and let  $\mathbf{g}_s$  be the dual gradient evaluated at  $\boldsymbol{\lambda}_s$ . Then*

1. *If  $\|\mathbf{g}_s\|_2 \leq \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$  then  $\alpha_s = 1$*
2. *If  $\|\mathbf{g}_s\|_2 > \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$  then  $\alpha_s \geq \beta \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2 \max_i \{\eta_i\}}$*

where  $B$  is a constant defined in Lemma 1 and  $\mu_2(\mathcal{L}), \mu_n(\mathcal{L})$  are the smallest and largest nonzero eigenvalues of the unweighted Laplacian of  $\mathcal{G}$ .

**Proof.** Combining  $\|\mathbf{g}_s\|_2 \leq \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$  with Lemma 2 implies:

$$\|\mathbf{g}_{s+1}\|_2 \leq \left(\frac{3}{2} - \alpha_s\right) \|\mathbf{g}_s\|_2 + \alpha_s\|\boldsymbol{\epsilon}_s\|_2 + \alpha_s^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\boldsymbol{\epsilon}_s\|_2^2$$

Since  $\|\epsilon_s\|_2 \leq \epsilon \frac{\mu_n(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\Gamma}{\gamma}} \|\mathbf{g}_s\|_2$ ,  $\|\mathbf{g}_s\|_2 \leq \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$  and  $\alpha_s \leq 1$ , then

$$\begin{aligned} \|\mathbf{g}_{s+1}\|_2 &\leq \left(\frac{3}{2} - \alpha_s\right) \|\mathbf{g}_s\|_2 + \alpha_s \epsilon \frac{\mu_n(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\Gamma}{\gamma}} \|\mathbf{g}_s\|_2 + \alpha_s^2 \epsilon^2 B \frac{\Gamma^3}{\mu_2^6(\mathcal{L})} \frac{\mu_n^2(\mathcal{L})}{\gamma} \|\mathbf{g}_s\|_2^2 \leq \\ &\left(\frac{3}{2} - \alpha_s\right) \|\mathbf{g}_s\|_2 + \epsilon \frac{\mu_n(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\Gamma}{\gamma}} \|\mathbf{g}_s\|_2 + \epsilon^2 B \frac{\Gamma^3}{\mu_2^6(\mathcal{L})} \frac{\mu_n^2(\mathcal{L})}{\gamma} \|\mathbf{g}_s\|_2^2 \leq \\ &\left(\frac{3}{2} - \alpha_s\right) \|\mathbf{g}_s\|_2 + \frac{1}{2} \epsilon \frac{\mu_2^2(\mathcal{L}) \gamma^2}{\mu_n(\mathcal{L}) dk \Gamma \delta} \left[ \frac{\mu_2(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\gamma}{\Gamma}} + \frac{\epsilon}{2} \right] = \left(\frac{3}{2} - \alpha_s\right) \|\mathbf{g}_s\|_2 + \hat{B} \end{aligned}$$

where we denote  $\hat{B} = \frac{1}{2} \epsilon \frac{\mu_2^2(\mathcal{L}) \gamma^2}{\mu_n(\mathcal{L}) dk \Gamma \delta} \left[ \frac{\mu_2(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\gamma}{\Gamma}} + \frac{\epsilon}{2} \right] \leq 2\epsilon \frac{n\gamma^2}{dk\Gamma\delta}$  for  $\epsilon \leq \frac{4}{n^3} \sqrt{\frac{\gamma}{\Gamma}}$ . Since  $\|\mathbf{g}_{s+1}\|_2 \geq \max_r \{[\mathcal{L}\mathbf{y}_r]_i\}$  and  $\|\mathbf{g}_s\|_2 \leq \sqrt{n} \max_i \{\eta_i\}$ , then

$$\max_r \{[\mathcal{L}\mathbf{y}_r]_i\} \leq \left(\frac{3}{2} - \alpha_s\right) \sqrt{n} \max_i \{\eta_i\} + 2\epsilon \frac{n\gamma^2}{dk\Gamma\delta}$$

Notice that if  $m_i = 0$  the  $\frac{3}{2} - \beta^{m_i} \leq 1 - \sigma\beta^{m_i}$ . Therefore, for  $m_i = 0$  we have

$$\max_r \{[\mathcal{L}\mathbf{y}_r]_i\} \leq (1 - \sigma\beta^{m_i}) \sqrt{n} \max_i \{\eta_i\} + 2\epsilon \frac{n\gamma^2}{dk\Gamma\delta}$$

In other words, Algorithm 3 returns  $\alpha_s = \beta^0 = 1$ .

For the case  $\|\mathbf{g}_s\|_2 > \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$  consider  $\bar{\alpha}_s = \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2 \sqrt{n} \max_i \{\eta_i\}}$ . Because  $\|\mathbf{g}_s\|_2 \leq \sqrt{n} \max_i \{\eta_i\}$  and  $\|\mathbf{g}_s\|_2 > \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$  then  $\bar{\alpha}_s < 1$ . Hence, applying  $\bar{\alpha}_s$  with  $\epsilon \leq \frac{4}{n^3} \sqrt{\frac{\gamma}{\Gamma}}$  for (22) gives:

$$\begin{aligned} \|\mathbf{g}_{s+1}\|_2 &\leq (1 - \bar{\alpha}_s) \|\mathbf{g}_s\|_2 + \bar{\alpha}_s^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_s\|_2^2 + \bar{\alpha}_s \|\epsilon_s\|_2 + \bar{\alpha}_s^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\epsilon_s\|_2^2 = \\ &\|\mathbf{g}_s\|_2 + \bar{\alpha}_s \|\epsilon_s\|_2 + \bar{\alpha}_s^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\epsilon_s\|_2^2 - \bar{\alpha}_s \|\mathbf{g}_s\|_2 \left[ 1 - \bar{\alpha}_s B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_s\|_2 \right] \leq \\ &\|\mathbf{g}_s\|_2 + \bar{\alpha}_s \epsilon \frac{\mu_n(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\Gamma}{\gamma}} \|\mathbf{g}_s\|_2 + \bar{\alpha}_s^2 \epsilon^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \frac{\mu_n^2(\mathcal{L})}{\mu_2^2(\mathcal{L})} \frac{\Gamma}{\gamma} \|\mathbf{g}_s\|_2^2 - \\ &\bar{\alpha}_s \|\mathbf{g}_s\|_2 \left[ 1 - \frac{\|\mathbf{g}_s\|_2}{2\sqrt{n} \max_i \{\eta_i\}} \right] \leq \|\mathbf{g}_s\|_2 + \bar{\alpha}_s \epsilon \frac{\mu_n(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\Gamma}{\gamma}} \|\mathbf{g}_s\|_2 + \\ &\bar{\alpha}_s^2 \epsilon^2 B \frac{\Gamma^2}{\mu_2^4(\mathcal{L})} \frac{\mu_n^2(\mathcal{L})}{\mu_2^2(\mathcal{L})} \frac{\Gamma}{\gamma} \|\mathbf{g}_s\|_2^2 - \frac{1}{2} \bar{\alpha}_s \|\mathbf{g}_s\|_2 = \left(1 - \frac{\bar{\alpha}_s}{2}\right) \|\mathbf{g}_s\|_2 + \\ &\epsilon \frac{\mu_n(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\Gamma}{\gamma}} \|\mathbf{g}_s\|_2 \frac{2B\Gamma^2}{\mu_2^2(\mathcal{L})} \sqrt{n} \max_i \{\eta_i\} + \epsilon^2 \frac{\mu_n^2(\mathcal{L})}{\mu_2^2(\mathcal{L})} \frac{\Gamma}{\gamma} \frac{1}{4 \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})}} \frac{\|\mathbf{g}_s\|_2^2}{n \max_i \{\eta_i^2\}} \leq \\ &\left(1 - \frac{\bar{\alpha}_s}{2}\right) \|\mathbf{g}_s\|_2 + \hat{B} \leq \left(1 - \frac{\bar{\alpha}_s}{2}\right) \|\mathbf{g}_s\|_2 + 2\epsilon \frac{n\gamma^2}{dk\delta\Gamma} \end{aligned}$$

In other words, we establishes:

$$\|\mathbf{g}_{s+1}\|_2 \leq (1 - \sigma\bar{\alpha}_s) \|\mathbf{g}_s\|_2 + 2\epsilon \frac{n\gamma^2}{dk\delta\Gamma}$$

Applying again  $\|\mathbf{g}_{s+1}\|_2 \geq \max_r \{[\mathcal{L}\mathbf{y}_r]_i\}$  and  $\|\mathbf{g}_s\|_2 \leq \sqrt{n} \max_i \{\eta_i\}$  gives:

$$\max_r \{[\mathcal{L}\mathbf{y}_r]_i\} \leq (1 - \sigma\bar{\alpha}_s) \sqrt{n} \max_i \{\eta_i\} + 2\epsilon \frac{n\gamma^2}{dk\Gamma\delta}$$

Therefore, Algorithm 3 returns  $\alpha_s \geq \beta\bar{\alpha}_s = \beta \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2 \max_i \{\eta_i\}}$ .

### C.2.5 Proof of the Main Theorem

In this appendix we prove Theorem ??:

**Theorem** Let  $\gamma, \Gamma, \delta, B$  be the constants defined in Assumption 1 and Lemma 1,  $\mu_2(\mathcal{L})$  and  $\mu_n(\mathcal{L})$  representing the smallest and largest nonzero eigenvalues of the unweighted Laplacian of  $\mathcal{G}$ ,  $\epsilon \leq \frac{\beta}{8} \frac{\gamma^3}{\Gamma^2 p \delta} \frac{\mu_2^4(\mathcal{L})}{\mu_n^3(\mathcal{L})}$  be the precision parameter for the SDD solver. Consider our iteration scheme with the step size  $\alpha_s$  is calculated by Algorithm 3. Then, this iteration scheme exhibits two convergence phases:

1. **Strict Decreases Phase** If  $\|\mathbf{g}_s\|_2 > \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$ , then

$$\|\mathbf{g}_{s+1}\|_2 - \|\mathbf{g}_k\|_2 \leq -\frac{\beta}{8\sqrt{np}\delta} \frac{\gamma^3}{\Gamma^2} \frac{\mu_2^4(\mathcal{L})}{\mu_n^3(\mathcal{L})}$$

where parameter  $\beta \in (0, 1)$ .

2. **Quadratic Decreases Phase** If  $\|\mathbf{g}_s\|_2 \leq \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$ , then for any  $o \geq 1$ :

$$\|\mathbf{g}_{s+o}\|_2 \leq \frac{1}{2^{2^o} \frac{B\Gamma^2}{\mu_2^2(\mathcal{L})}} + \hat{B} + \frac{\tilde{\Lambda}}{\frac{B\Gamma^2}{\mu_2^2(\mathcal{L})}} \left[ \frac{2^{2^l-1} - 1}{2^{2^l}} \right]$$

where

$$\begin{aligned} \hat{B} &= \frac{1}{2} \epsilon \frac{\mu_2^2(\mathcal{L}) \gamma^2}{\mu_n(\mathcal{L}) p \Gamma \delta} \left[ \frac{\mu_2(\mathcal{L})}{\mu_n(\mathcal{L})} \sqrt{\frac{\gamma}{\Gamma}} + \frac{\epsilon}{2} \right] \sim \mathcal{O}(\epsilon) \\ \tilde{\Lambda} &= \hat{B} \frac{4B\Gamma^2}{\mu_2^4(\mathcal{L})} \left[ 1 + \hat{B} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \right] \sim \mathcal{O}(\epsilon) \end{aligned}$$

**Proof.** We will proof the above theorem by handling each of the cases separately. We start by considering the case when  $\|\mathbf{g}_s\|_2 > \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2}$ . Then, according to Lemma 3:  $\alpha_s \geq \beta \frac{\mu_2^4(\mathcal{L})}{2B\Gamma^2 \max_i \{\eta_i\}}$  and Equation (22) we have:

$$\|\mathbf{g}_{s+1}\|_2 \leq (1 - \frac{1}{2} \beta \bar{\alpha}_s) \|\mathbf{g}_s\|_2 + 2\epsilon \frac{n\gamma^2}{p\delta\Gamma}$$

Choosing  $\epsilon \leq \frac{\beta}{8} \frac{\gamma^3}{\Gamma^2 p \delta} \frac{\mu_2^4(\mathcal{L})}{\mu_n^3(\mathcal{L})}$  implies  $2\epsilon \frac{n\gamma^2}{p\delta\Gamma} \leq \frac{1}{4} \beta \bar{\alpha}_s \|\mathbf{g}_s\|_2$  and

$$\begin{aligned} \|\mathbf{g}_{s+1}\|_2 - \|\mathbf{g}_s\|_2 &\leq -\frac{1}{4} \beta \bar{\alpha}_s \|\mathbf{g}_s\|_2 \leq -\frac{1}{4} \beta \frac{B\Gamma^2}{2 \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \sqrt{n} \max_i \{\eta_i\}} \|\mathbf{g}_s\|_2 \\ &= -\frac{1}{8} \beta \frac{1}{\frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \sqrt{n}} = -\frac{\beta}{8\sqrt{np}\delta} \frac{\gamma^3}{\Gamma^2} \frac{\mu_2^4(\mathcal{L})}{\mu_n^3(\mathcal{L})} \end{aligned}$$

The the quadratic decrease phase we use the result of Lemma 3 and induction:

1. For  $m = 1$  applying  $\alpha_s = 1$  in Equation (22):

$$\|\mathbf{g}_{s+1}\|_2 \leq \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_s\|_2^2 + \hat{B} \leq \frac{1}{4 \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})}} + \hat{B}$$

This result validates the claim for  $m = 1$ .

2. Let us assume it is correct for some  $m > 0$ .

3. Using  $\alpha_{s+m+1} = 1$  in Equation (22) and denoting  $u = 2^{2^m}$  gives :

$$\begin{aligned} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_{s+m+1}\|_2 &\leq \left[ \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_{s+m}\|_2 \right]^2 + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} \leq \\ &\left[ \frac{1}{u} + \hat{B} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} + \tilde{\Lambda} \frac{\frac{1}{2}u - 1}{u} \right]^2 + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} = \\ &\frac{1}{u^2} + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} + \tilde{\Lambda} \frac{\frac{1}{2}u^2 - 1}{u^2} - \tilde{\Lambda} \frac{\frac{1}{2}u^2 - 1}{u^2} + \tilde{\Lambda} \frac{u - 2}{u^2} + \hat{B} \frac{2B\Gamma^2}{\mu_2^4(\mathcal{L})} \frac{1}{u} + \\ &\left( \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \right)^2 \left[ \hat{B}^2 + 2\hat{B}\tilde{\Lambda} \frac{1}{\frac{B\Gamma^2}{\mu_2^4(\mathcal{L})}} \frac{(u-2)}{u} + \tilde{\Lambda}^2 \frac{1}{\left( \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \right)^2} \frac{(u-2)^2}{4u^2} \right] \end{aligned}$$

Since  $\hat{B} + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B}^2 = \frac{\tilde{\Lambda}}{4 \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})}}$ , then

$$\begin{aligned} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \|\mathbf{g}_{s+m+1}\|_2 &\leq \\ \frac{1}{u^2} + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} + \tilde{\Lambda} \frac{\frac{1}{2}u^2 - 1}{u^2} - \tilde{\Lambda} \frac{\frac{1}{2}u^2 - 1}{u^2} + \tilde{\Lambda} \frac{u - 2}{u^2} + \hat{B} \frac{2B\Gamma^2}{\mu_2^4(\mathcal{L})} \frac{1}{u} + \\ &\left( \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \right)^2 \left[ \tilde{\Lambda} \frac{1}{4 \left( \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \right)^2} - \frac{\hat{B}}{\frac{B\Gamma^2}{\mu_2^4(\mathcal{L})}} + \frac{2\hat{B}\tilde{\Lambda}}{\mu_2^4(\mathcal{L})} \left( \frac{1}{2} - \frac{1}{u} \right) + \frac{\tilde{\Lambda}^2}{\left( \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \right)^2} \frac{(u-2)^2}{u^2} \right] = \\ &\frac{1}{u^2} + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} + \tilde{\Lambda} \frac{(u^2 - 2)}{2u^2} + \frac{\tilde{\Lambda}}{u^2} \left[ -\frac{1}{2}u^2 + u - 1 \right] + \frac{\tilde{\Lambda}}{4} + \hat{B} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \frac{2}{u} + \\ &\hat{B}\tilde{\Lambda} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \left[ 1 - \frac{2}{u} \right] + \tilde{\Lambda}^2 \left( \frac{u-2}{2u} \right)^2 = \frac{1}{u^2} + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} + \tilde{\Lambda} \frac{(u^2 - 2)}{2u^2} - \\ &\frac{\tilde{\Lambda}}{u^2} \left( \frac{u}{2} - 1 \right)^2 + \tilde{\Lambda}^2 \left( \frac{1}{2} - \frac{1}{u} \right)^2 + \hat{B} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \left[ -1 + \frac{2}{u} + \tilde{\Lambda} - \frac{2}{u} \tilde{\Lambda} \right] = \\ &\frac{1}{u^2} + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} + \tilde{\Lambda} \frac{(u^2 - 2)}{2u^2} - \left( \frac{1}{2} - \frac{1}{u} \right)^2 (\tilde{\Lambda} - \tilde{\Lambda}^2) - \hat{B} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} (1 - \tilde{\Lambda}) \left( 1 - \frac{2}{u} \right) \leq \\ &\frac{1}{u^2} + \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} \hat{B} + \tilde{\Lambda} \frac{(u^2 - 2)}{2u^2} = \frac{1}{2^{2^{m+1}}} + \hat{B} \frac{B\Gamma^2}{\mu_2^4(\mathcal{L})} + \tilde{\Lambda} \left[ \frac{2^{2^{m+1}} - 1}{2^{2^{m+1}}} \right] \end{aligned}$$

The last step follows due to  $u > 2$  and  $\hat{\Lambda} < 1$  (choosing  $\epsilon$  small enough).

Hence, our claim is correct.