

## Learning (from) Deep Hierarchical Structure among Features: Supplementary Material

### Proof for Theorem 1

**Proof.** We prove the convexity via the first-order condition. For any  $\bar{w} \in \mathbb{R}$  and  $\bar{\mathbf{z}} \geq \mathbf{0}$ , We can easily compute  $\frac{\partial f(\bar{w}, \bar{\mathbf{z}})}{\partial w} = \frac{2\bar{w}}{\prod_{i=1}^m \bar{z}_i^{\theta_i}}$  and  $\frac{\partial f(\bar{w}, \bar{\mathbf{z}})}{\partial z_j} = -\frac{\theta_j \bar{w}^2}{\bar{z}_j \prod_{i=1}^m \bar{z}_i^{\theta_i}}$ . Then for any  $\hat{w}, \hat{\mathbf{z}} \in \mathbb{R}$  and  $\hat{\mathbf{z}}, \bar{\mathbf{z}} \geq \mathbf{0}$ , we have

$$\begin{aligned} f(\hat{w}, \hat{\mathbf{z}}) - f(\bar{w}, \bar{\mathbf{z}}) - (\hat{w} - \bar{w}) \frac{\partial f(\bar{w}, \bar{\mathbf{z}})}{\partial w} - \sum_{i=1}^m (\hat{z}_i - \bar{z}_i) \frac{\partial f(\bar{w}, \bar{\mathbf{z}})}{\partial z_i} \\ = \frac{\hat{w}^2}{\prod_{i=1}^m \hat{z}_i^{\theta_i}} - \frac{\bar{w}^2}{\prod_{i=1}^m \bar{z}_i^{\theta_i}} - \frac{2\bar{w}(\hat{w} - \bar{w})}{\prod_{i=1}^m \bar{z}_i^{\theta_i}} + \frac{\bar{w}^2 \sum_{j=1}^m \theta_j \left(\frac{\hat{z}_j}{\bar{z}_j} - 1\right)}{\prod_{i=1}^m \bar{z}_i^{\theta_i}} \\ = \frac{1}{\prod_{i=1}^m \hat{z}_i^{\theta_i}} \hat{w}^2 - \frac{2}{\prod_{i=1}^m \bar{z}_i^{\theta_i}} \bar{w} \hat{w} + \frac{\sum_{j=1}^m \theta_j \frac{\hat{z}_j}{\bar{z}_j}}{\prod_{i=1}^m \bar{z}_i^{\theta_i}} \bar{w}^2, \end{aligned}$$

where the second equality holds due to the property of  $\theta$  that  $\sum_{i=1}^m \theta_i = 1$ . Since  $\ln(\cdot)$  is concave and  $\theta$  satisfies that  $\theta_i \geq 0$  for  $i = 1, \dots, m$  and  $\sum_{i=1}^m \theta_i = 1$ , we can have the following inequalities:

$$\begin{aligned} \sum_{i=1}^m \theta_i \ln \frac{\hat{z}_i}{\bar{z}_i} &\leq \ln \left( \sum_{j=1}^m \frac{\hat{z}_j}{\bar{z}_j} \theta_j \right) \\ \Rightarrow \sum_{i=1}^m \theta_i \ln \hat{z}_i &\leq \sum_{i=1}^m \theta_i \ln \bar{z}_i + \ln \left( \sum_{j=1}^m \frac{\hat{z}_j}{\bar{z}_j} \theta_j \right) \\ \Rightarrow \prod_{i=1}^m \hat{z}_i^{\theta_i} &\leq \left( \prod_{i=1}^m \bar{z}_i^{\theta_i} \right) \left( \sum_{j=1}^m \frac{\hat{z}_j}{\bar{z}_j} \theta_j \right) \\ \Rightarrow \frac{\prod_{i=1}^m \hat{z}_i^{\theta_i}}{\left( \prod_{i=1}^m \bar{z}_i^{\theta_i} \right)^2} &\leq \frac{\sum_{j=1}^m \frac{\hat{z}_j}{\bar{z}_j} \theta_j}{\prod_{i=1}^m \bar{z}_i^{\theta_i}}. \end{aligned}$$

Then we can have

$$\begin{aligned} f(\hat{w}, \hat{\mathbf{z}}) - f(\bar{w}, \bar{\mathbf{z}}) - (\hat{w} - \bar{w}) \frac{\partial f(\bar{w}, \bar{\mathbf{z}})}{\partial w} - \sum_{i=1}^m (\hat{z}_i - \bar{z}_i) \frac{\partial f(\bar{w}, \bar{\mathbf{z}})}{\partial z_i} \\ = \frac{\left( \hat{w} - \frac{\prod_{i=1}^m \hat{z}_i^{\theta_i}}{\prod_{i=1}^m \bar{z}_i^{\theta_i}} \bar{w} \right)^2}{\prod_{i=1}^m \hat{z}_i^{\theta_i}} + \left( \frac{\sum_{j=1}^m \theta_j \frac{\hat{z}_j}{\bar{z}_j}}{\prod_{i=1}^m \bar{z}_i^{\theta_i}} - \frac{\prod_{i=1}^m \hat{z}_i^{\theta_i}}{\left( \prod_{i=1}^m \bar{z}_i^{\theta_i} \right)^2} \right) \bar{w}^2 \end{aligned}$$

$\geq 0$ ,

in which we reach the conclusion.  $\square$

### Proof for Theorem 2

**Proof.** It is obvious that the first and third terms of the objective function in problem (1) as well as the linear constraints are convex. Moreover, based on Theorem 1, each summand in the second term of the objective function in problem (1) is jointly convex with respect to  $\mathbf{w}$  and  $\sigma$ , making the whole problem convex.  $\square$

### Proof for Theorem 3

**Proof.** We first rewrite problem (1) where  $m = 3$  as

$$\begin{aligned} \min_{\mathbf{w}, \sigma} \frac{1}{n} \sum_{i=1}^n l(\mathbf{x}_i, y_i, \mathbf{w}) + \frac{\lambda_1}{2} \sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \frac{w_k^2}{(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} \\ + \frac{\lambda_2}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t. } \sum_{i=1}^{s_1} d_i^{(1)} \sigma_{1,i}^{(1)} = 1, \sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} d_j^{(2)} \sigma_{i,j}^{(2)} = 1, \sum_{j=1}^{s_2} \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = 1 \\ \sigma_{1,i}^{(1)} \geq 0 \forall i, \sigma_{i,j}^{(2)} \geq 0 \forall i, j, \sigma_{j,k}^{(3)} \geq 0 \forall j, k. \end{aligned} \quad (13)$$

The Lagrangian of problem (13) is

$$\begin{aligned} \mathcal{L} = \frac{1}{n} \sum_{i=1}^n l(\mathbf{x}_i, y_i, \mathbf{w}) + \frac{\lambda_1}{2} \sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \frac{w_k^2}{(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} \\ + \frac{\lambda_2}{2} \|\mathbf{w}\|_2^2 + \nu_1 \left( \sum_{i=1}^{s_1} d_i^{(1)} \sigma_{1,i}^{(1)} - 1 \right) + \nu_2 \left( \sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} d_j^{(2)} \sigma_{i,j}^{(2)} - 1 \right) \\ + \nu_3 \left( \sum_{j=1}^{s_2} \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} - 1 \right) - \sum_{i=1}^{s_1} \varepsilon_{1,i}^{(1)} \sigma_{1,i}^{(1)} - \sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} \varepsilon_{i,j}^{(2)} \sigma_{i,j}^{(2)} \\ - \sum_{j=1}^{s_2} \sum_{k \in \mathcal{C}_j^2} \varepsilon_{j,k}^{(3)} \sigma_{j,k}^{(3)}. \end{aligned}$$

The optimality conditions for  $\sigma_{1,i}^{(1)}$ ,  $\sigma_{i,j}^{(2)}$  and  $\sigma_{j,k}^{(3)}$  are

$$\frac{\partial \mathcal{L}}{\partial \sigma_{1,i}^{(1)}} = 0 \Rightarrow d_i^{(1)} \nu_1 - \varepsilon_{1,i}^{(1)} - \frac{\lambda_1}{6} \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \frac{w_k^2}{(\sigma_{1,i}^{(1)})^{\frac{4}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} = 0 \quad (14)$$

$$\frac{\partial \mathcal{L}}{\partial \sigma_{i,j}^{(2)}} = 0 \Rightarrow d_j^{(2)} \nu_2 - \varepsilon_{i,j}^{(2)} - \frac{\lambda_1}{6} \sum_{k \in \mathcal{C}_j^2} \frac{w_k^2}{(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{4}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} = 0 \quad (15)$$

$$\frac{\partial \mathcal{L}}{\partial \sigma_{j,k}^{(3)}} = 0 \Rightarrow \nu_3 - \varepsilon_{j,k}^{(3)} - \frac{\lambda_1 w_k^2}{6(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{4}{3}}} = 0. \quad (16)$$

Multiplying Eqs. (14), (15) and (16) by  $\sigma_{1,i}^{(1)}$ ,  $\sigma_{i,j}^{(2)}$  and  $\sigma_{j,k}^{(3)}$  respectively gives

$$d_i^{(1)} \nu_1 \sigma_{1,i}^{(1)} - \frac{\lambda_1}{6} \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \frac{w_k^2}{(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} = 0 \quad (17)$$

$$d_j^{(2)} \nu_2 \sigma_{i,j}^{(2)} - \frac{\lambda_1}{6} \sum_{k \in \mathcal{C}_j^2} \frac{w_k^2}{(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} = 0 \quad (18)$$

$$\nu_3 \sigma_{j,k}^{(3)} - \frac{\lambda_1 w_k^2}{6(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} = 0, \quad (19)$$

where we use the KKT conditions that  $\varepsilon_{1,i}^{(1)} \sigma_{1,i}^{(1)} = 0$ ,  $\varepsilon_{i,j}^{(2)} \sigma_{i,j}^{(2)} = 0$  and  $\varepsilon_{j,k}^{(3)} \sigma_{j,k}^{(3)} = 0$ . Based on Eq. (19), we have

$$\lambda_1 w_k^2 = 6\nu_3 (\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{4}{3}}.$$

By plugging this equation into Eqs. (17) and (18), we can get

$$d_i^{(1)} \nu_1 \sigma_{1,i}^{(1)} - \nu_3 \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = 0 \quad (20)$$

$$d_j^{(2)} \nu_2 \sigma_{i,j}^{(2)} - \nu_3 \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = 0. \quad (21)$$

By taking sum over  $i$  on Eq. (20), we can get  $\nu_1 = \nu_3$  since  $\sum_{i=1}^{s_1} d_i^{(1)} \sigma_{1,i}^{(1)} = 1$  and  $\sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = \sum_{j=1}^{s_2} \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = 1$  due to linear equality constraints in problem (13). By taking sum over  $i$  and  $j$  on Eq. (21), we have  $\nu_2 = \nu_3$  since  $\sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} d_j^{(2)} \sigma_{i,j}^{(2)} = 1$  and  $\sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = \sum_{j=1}^{s_2} \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = 1$  due to linear equality constraints in problem (13). By combining those two results, we have  $\nu_1 = \nu_2 = \nu_3$  and hence we use  $\nu$  to denote them. Based on Eqs. (17-19), we have  $\nu > 0$ . Based on Eq. (19), we have

$$\sigma_{j,k}^{(3)} = \frac{\lambda_1^{\frac{3}{4}} |w_k|^{\frac{3}{2}}}{6^{\frac{3}{4}} \nu^{\frac{3}{4}} (\sigma_{1,i}^{(1)})^{\frac{1}{4}} (\sigma_{i,j}^{(2)})^{\frac{1}{4}}}. \quad (22)$$

Then based on Eqs. (18), (19) and (22), since  $\nu_2 = \nu_3 = \nu > 0$ , we have

$$d_j^{(2)} \sigma_{i,j}^{(2)} = \sum_{k \in \mathcal{C}_j^2} \sigma_{j,k}^{(3)} = \frac{\lambda_1^{\frac{3}{4}}}{6^{\frac{3}{4}} \nu^{\frac{3}{4}} (\sigma_{1,i}^{(1)})^{\frac{1}{4}} (\sigma_{i,j}^{(2)})^{\frac{1}{4}}} \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}},$$

leading to

$$\sigma_{i,j}^{(2)} = \frac{\lambda_1^{\frac{3}{5}}}{6^{\frac{3}{5}} \nu^{\frac{3}{5}} (d_j^{(2)})^{\frac{4}{5}} (\sigma_{1,i}^{(1)})^{\frac{1}{5}}} \left( \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \right)^{\frac{4}{5}}. \quad (23)$$

Based on Eqs. (17), (18) and (23), since  $\nu_1 = \nu_2 = \nu > 0$ , we can get

$$\begin{aligned} d_i^{(1)} \sigma_{1,i}^{(1)} &= \sum_{j \in \mathcal{C}_i^1} d_j^{(2)} \sigma_{i,j}^{(2)} \\ &= \frac{\lambda_1^{\frac{3}{5}}}{6^{\frac{3}{5}} \nu^{\frac{3}{5}} (\sigma_{1,i}^{(1)})^{\frac{1}{5}}} \sum_{j \in \mathcal{C}_i^1} (d_j^{(2)})^{\frac{1}{5}} \left( \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \right)^{\frac{4}{5}}, \end{aligned}$$

leading to

$$\sigma_{1,i}^{(1)} = \frac{\lambda_1^{\frac{1}{2}}}{6^{\frac{1}{2}} \nu^{\frac{1}{2}} (d_i^{(1)})^{\frac{5}{6}}} \left( \sum_{j \in \mathcal{C}_i^1} (d_j^{(2)})^{\frac{1}{5}} \left( \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \right)^{\frac{4}{5}} \right)^{\frac{5}{6}}. \quad (24)$$

By plugging Eq. (24) into the first linear equality constraint in problem (13), we can get

$$\nu = \frac{\lambda_1}{6} \left( \sum_{i=1}^{s_1} (d_i^{(1)})^{\frac{1}{6}} \left( \sum_{j \in \mathcal{C}_i^1} (d_j^{(2)})^{\frac{1}{5}} \left( \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \right)^{\frac{4}{5}} \right)^{\frac{5}{6}} \right)^2. \quad (25)$$

Then based on Eqs. (22-25), we can get

$$\begin{aligned} & (\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}} \\ &= \frac{\lambda_1^{\frac{1}{4}} |w_k|^{\frac{1}{2}} (\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}}}{6^{\frac{1}{4}} \nu^{\frac{1}{4}} (\sigma_{1,i}^{(1)})^{\frac{1}{12}} (\sigma_{i,j}^{(2)})^{\frac{1}{12}}} \\ &= \frac{\lambda_1^{\frac{1}{4}} |w_k|^{\frac{1}{2}}}{6^{\frac{1}{4}} \nu^{\frac{1}{4}}} (\sigma_{1,i}^{(1)})^{\frac{1}{4}} (\sigma_{i,j}^{(2)})^{\frac{1}{4}} \\ &= \frac{\lambda_1^{\frac{2}{5}} (\sigma_{1,i}^{(1)})^{\frac{1}{5}} |w_k|^{\frac{1}{2}}}{6^{\frac{2}{5}} \nu^{\frac{2}{5}} (d_j^{(2)})^{\frac{1}{5}}} \left( \sum_{k' \in \mathcal{C}_j^2} |w_{k'}|^{\frac{3}{2}} \right)^{\frac{1}{5}} \\ &= \frac{\lambda_1^{\frac{1}{2}} |w_k|^{\frac{1}{2}}}{6^{\frac{1}{2}} \nu^{\frac{1}{2}} (d_i^{(1)})^{\frac{1}{6}} (d_j^{(2)})^{\frac{1}{5}}} \left( \sum_{k' \in \mathcal{C}_j^2} |w_{k'}|^{\frac{3}{2}} \right)^{\frac{1}{5}} \times \\ & \left( \sum_{j' \in \mathcal{C}_i^1} (d_{j'}^{(2)})^{\frac{1}{5}} \left( \sum_{k' \in \mathcal{C}_{j'}^2} |w_{k'}|^{\frac{3}{2}} \right)^{\frac{4}{5}} \right)^{\frac{1}{6}} \end{aligned}$$

and

$$\begin{aligned} & \sum_{i=1}^{s_1} \sum_{j \in \mathcal{C}_i^1} \sum_{k \in \mathcal{C}_j^2} \frac{w_k^2}{(\sigma_{1,i}^{(1)})^{\frac{1}{3}} (\sigma_{i,j}^{(2)})^{\frac{1}{3}} (\sigma_{j,k}^{(3)})^{\frac{1}{3}}} \\ &= \frac{6^{\frac{1}{2}} \nu^{\frac{1}{2}}}{\lambda_1^{\frac{1}{2}}} \sum_{i=1}^{s_1} \frac{(d_i^{(1)})^{\frac{1}{6}}}{\left( \sum_{j' \in \mathcal{C}_i^1} (d_{j'}^{(2)})^{\frac{1}{5}} \left( \sum_{k' \in \mathcal{C}_{j'}^2} |w_{k'}|^{\frac{3}{2}} \right)^{\frac{4}{5}} \right)^{\frac{1}{6}}} \times \\ & \sum_{j \in \mathcal{C}_i^1} \frac{(d_j^{(2)})^{\frac{1}{5}}}{\left( \sum_{k' \in \mathcal{C}_j^2} |w_{k'}|^{\frac{3}{2}} \right)^{\frac{1}{5}}} \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \\ &= \frac{6^{\frac{1}{2}} \nu^{\frac{1}{2}}}{\lambda_1^{\frac{1}{2}}} \sum_{i=1}^{s_1} \frac{(d_i^{(1)})^{\frac{1}{6}}}{\left( \sum_{j' \in \mathcal{C}_i^1} (d_{j'}^{(2)})^{\frac{1}{5}} \left( \sum_{k' \in \mathcal{C}_{j'}^2} |w_{k'}|^{\frac{3}{2}} \right)^{\frac{4}{5}} \right)^{\frac{1}{6}}} \times \\ & \sum_{j \in \mathcal{C}_i^1} (d_j^{(2)})^{\frac{1}{5}} \left( \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \right)^{\frac{4}{5}} \\ &= \frac{6^{\frac{1}{2}} \nu^{\frac{1}{2}}}{\lambda_1^{\frac{1}{2}}} \sum_{i=1}^{s_1} (d_i^{(1)})^{\frac{1}{6}} \left( \sum_{j \in \mathcal{C}_i^1} (d_j^{(2)})^{\frac{1}{5}} \left( \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \right)^{\frac{4}{5}} \right)^{\frac{5}{6}} \\ &= \left( \sum_{i=1}^{s_1} (d_i^{(1)})^{\frac{1}{6}} \left( \sum_{j \in \mathcal{C}_i^1} (d_j^{(2)})^{\frac{1}{5}} \left( \sum_{k \in \mathcal{C}_j^2} |w_k|^{\frac{3}{2}} \right)^{\frac{4}{5}} \right)^{\frac{5}{6}} \right)^2. \end{aligned}$$

By plugging this equation into the objective function of problem (13), we reach the conclusion.  $\square$

## The FISTA and GIST Algorithms

The detailed procedures for the FISTA and GIST algorithms are shown in Algorithms 1 and 2.

---

**Algorithm 1** The FISTA Algorithm

---

- 1: Initialize  $r_0, \alpha > 1$ ;
- 2: Set  $\phi^{(0)}$  as the initial value for variable  $\phi$ ;
- 3:  $\psi^{(1)} := \phi^{(0)}$ ;
- 4:  $k := 1$ ;
- 5:  $t_1 := 1$ ;
- 6: **while** not converged **do**
- 7: Find the smallest nonnegative integers  $i_k$  such that with  $\hat{r} = \alpha^{i_k} r_{k-1}$ ,  $F(q_{\hat{r}}(\psi^{(k)})) \leq Q_{\hat{r}}(q_{\hat{r}}(\psi^{(k)}), \psi^{(k)})$ ;
- 8:  $r_k := \alpha^{i_k} r_{k-1}$ ;
- 9:  $\phi^{(k)} := q_{r_k}(\psi^{(k)})$ ;
- 10:  $t_{k+1} := \frac{1 + \sqrt{1 + 4t_k^2}}{2}$ ;
- 11:  $\psi^{(k+1)} := \phi^{(k)} + \left(\frac{t_k - 1}{t_{k+1}}\right) (\phi^{(k)} - \phi^{(k-1)})$
- 12:  $k := k + 1$ ;
- 13: **end while**

---

---

**Algorithm 2** The GIST Algorithm

---

- 1: Choose  $\eta > 1, r_0, \sigma \in (0, 1)$ ;
- 2: Set  $\phi^{(0)}$  as the initial value for variable  $\phi$ ;
- 3:  $k := 0$ ;
- 4: **repeat**
- 5:  $r_{k+1} := r_k$ ;
- 6: **repeat**
- 7:  $\phi^{(k+1)} := \arg \min_{\phi} H_{r_{k+1}}(\phi, \phi^{(k)})$ ;
- 8:  $r_{k+1} := \eta r_{k+1}$ ;
- 9: **until**  $F(\phi^{(k+1)}) \leq F(\phi^{(k)}) - \frac{r_{k+1}\sigma}{2} \|\phi^{(k+1)} - \phi^{(k)}\|_F^2$
- 10:  $k := k + 1$ ;
- 11: **until** Some convergence criterion is satisfied

---

### Optimization Procedure for Problem (4)

Note that the only coupling in the variables of problem (4) comes from the equality constraint. The Lagrangian of problem (4) with respect to the equality constraint is given by

$$\mathcal{L}(\boldsymbol{\rho}, \gamma) = \|\boldsymbol{\rho} - \hat{\boldsymbol{\rho}}\|_2^2 - 2\gamma(\mathbf{a}^T \boldsymbol{\rho} - 1).$$

By setting the derivative of  $\mathcal{L}$  with respect to each  $\rho_i$  to 0, we can see that the minimum is reached when  $\rho_i = \hat{\rho}_i + a_i \gamma$  where  $a_i$  is the  $i$ th entry of  $\mathbf{a}$ . Since each  $\rho_i$  is required to be nonnegative and  $\mathcal{L}(\boldsymbol{\rho}, \gamma)$  is a quadratic function of  $\boldsymbol{\rho}$ , the optimal solution for  $\rho_i$  is given by

$$\rho_i = \max(0, \hat{\rho}_i + a_i \gamma). \quad (26)$$

Plugging the optimal solution of  $\rho_i$  into  $\mathcal{L}(\boldsymbol{\rho}, \gamma)$ , we can obtain the dual problem as

$$\min_{\gamma} \sum_{\gamma \geq -\hat{\rho}_i/a_i} (a_i^2 \gamma^2 + 2a_i \hat{\rho}_i \gamma) - \sum_{\gamma < -\hat{\rho}_i/a_i} \hat{\rho}_i^2 - 2\gamma. \quad (27)$$

Obviously, the objective function of problem (27) is a piecewise linear or quadratic function over regions determined by the sequence  $\{-\frac{\hat{\rho}_i}{a_i}\}$ . The main idea of our method is to determine the functional form of problem (27) over each

region, then compute the local optimum over each region which has an analytical solution, and finally obtain the global optimum by comparing all the local optima. So the main problem is to determine the coefficients of problem (27) over each region efficiently. Without loss of the generality, we assume  $\hat{\rho}_1/a_1 \geq \hat{\rho}_2/a_2 \geq \dots \geq \hat{\rho}_s/a_s$  where  $s$  is the length of  $\mathbf{a}$ . If this is not the case, we can sort the sequence  $\{\hat{\rho}_i/a_i\}$ . When  $\gamma \in [-\hat{\rho}_s/a_s, +\infty)$ , the objective function of problem (27) is  $c_2 \lambda^2 + c_1 \lambda + c_0$ , where  $c_2 = \sum_{i=1}^s a_i^2$ ,  $c_1 = 2(\sum_{i=1}^s a_i \hat{\rho}_i - 1)$ , and  $c_0 = 0$ , and it has an analytical solution as  $\gamma = \max(-\frac{\hat{\rho}_s}{a_s}, -\frac{c_1}{2c_2})$ . When  $\gamma \in (-\infty, -\hat{\rho}_1/a_1)$ , problem (27) has no well-defined solution since the objective function becomes  $-2\gamma - \sum_{i=1}^s \hat{\rho}_i^2$ . So we only need to consider the situation where  $\gamma \in [-\hat{\rho}_1/a_1, -\hat{\rho}_s/a_s]$ . We summarize the algorithm for solving problem (27) in Algorithm 3. This algorithm needs to scan the sequence  $\{-\hat{\rho}_i/a_i\}$  only once which costs  $O(s)$ . So the complexity of the whole algorithm is at most  $O(s \ln s)$ , which is more efficient than existing QP solvers.

After determining the optimal  $\gamma$ , we can obtain the solution for  $\rho_i$  via Eq. (26).

---

**Algorithm 3** Algorithm for problem (27)

---

- 1: Sort  $\{\hat{\rho}_i/a_i\}$  if needed;
- 2:  $c_0 := 0$ ; % coefficient for constant term
- 3:  $c_1 := 2(\sum_{i=1}^s a_i \hat{\rho}_i - 1)$ ; % coefficient for linear term
- 4:  $c_2 := \sum_{i=1}^s a_i^2$ ; % coefficient for quadratic term
- 5:  $\gamma := \max(-\frac{\hat{\rho}_s}{a_s}, -\frac{c_1}{2c_2})$ ;
- 6:  $f := c_0 + c_1 \gamma + c_2 \gamma^2$ ; % value of current minimum
- 7: **for**  $i = s$  to 2 **do**
- 8: % Determine the coefficients over  $[-\hat{\rho}_{i-1}/a_{i-1}, -\hat{\rho}_i/a_i]$ ;
- 9:  $c_0 := c_0 - \hat{\rho}_i^2$ ;
- 10:  $c_1 := c_1 - 2a_i \hat{\rho}_i$ ;
- 11:  $c_2 := c_2 - a_i^2$ ;
- 12:  $\gamma_0 := \min(-\frac{\hat{\rho}_i}{a_i}, \max(-\frac{\hat{\rho}_{i-1}}{a_{i-1}}, -\frac{c_1}{2c_2}))$ ;
- 13:  $f_0 := c_0 + c_1 \gamma_0 + c_2 \gamma_0^2$ ;
- 14: **if**  $f_0 < f$  **then**
- 15:  $\gamma := \gamma_0$ ;
- 16:  $f := f_0$ ;
- 17: **end if**
- 18: **end for**

---

### Optimization Procedure for Problem (5)

Unfortunately, problem (5) is non-convex with respect to all variables. To see this, we present a case that  $f(w, \mathbf{z}, \boldsymbol{\theta}) = \frac{w^2}{\prod_{i=1}^m z_i^{\theta_i}}$  is not jointly convex with respect to  $w \in \mathbb{R}$ ,  $\mathbf{z} \in \mathbb{R}^m$  and  $\boldsymbol{\theta}$ , where  $w$  and  $\{z_i\}$  are required to be positive, and  $\boldsymbol{\theta}$  is required to satisfy  $\theta_i \geq 0$  for  $i = 1, \dots, m$  and  $\sum_{i=1}^m \theta_i = 1$ . By considering a case that  $d = 1$  and  $m = 2$  and defining  $\hat{w} = -48$ ,  $\hat{w} = -74$ ,  $\hat{\mathbf{z}} = (60 \ 86)^T$ ,  $\hat{\mathbf{z}} = (99 \ 93)^T$ ,  $\hat{\boldsymbol{\theta}} = (0.0006 \ 0.9994)^T$ ,  $\hat{\boldsymbol{\theta}} = (0.8286 \ 0.1714)^T$ ,

and  $\alpha = 0.4095$ , we have

$$\begin{aligned} & \alpha f(\hat{\mathbf{w}}, \hat{\mathbf{z}}, \hat{\boldsymbol{\theta}}) + (1 - \alpha)f(\bar{\mathbf{w}}, \bar{\mathbf{z}}, \bar{\boldsymbol{\theta}}) \\ &= 0.4095 \times 26.7968 + (1 - 0.4095) \times 55.9091 \\ &< 46.3551 \\ &= f(\alpha\hat{\mathbf{w}} + (1 - \alpha)\bar{\mathbf{w}}, \alpha\hat{\mathbf{z}} + (1 - \alpha)\bar{\mathbf{z}}, \alpha\hat{\boldsymbol{\theta}} + (1 - \alpha)\bar{\boldsymbol{\theta}}), \end{aligned}$$

which implies that problem (5) is non-convex.

In order to solve problem (5), we use the GIST algorithm. With the abuse of notations, we use a vector  $\phi$  to denote the concatenation of  $\mathbf{w}$ ,  $\boldsymbol{\sigma}$  and  $\boldsymbol{\theta}$ . We define the set of constraints on  $\phi$  as  $\mathcal{S}_\phi = \{\phi \mid \sum_{j=1}^{s_i} d_j^{(i)} \sigma_{\mathcal{F}_j^i, j}^{(i)} = 1 \forall i \in [m], \sigma_{\mathcal{F}_j^i, j}^{(i)} \geq 0 \forall i, j, \sum_{j=1}^m \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} = 1 \forall i \in [d], \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} \geq 0 \forall i, j\}$ . We define  $f(\phi) = \frac{1}{n} \sum_{i=1}^n l(\mathbf{x}_i, y_i, \mathbf{w}) + \frac{\lambda_1}{2} \sum_{i=1}^d \frac{w_i^2}{\prod_{j=1}^m \left( \sigma_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} \right)^{\theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)}}}$  and  $g(\phi) = \begin{cases} \frac{\lambda_2}{2} \|\mathbf{w}\|_2^2 & \text{if } \phi \in \mathcal{S}_\phi \\ +\infty, & \text{otherwise} \end{cases}$  as an extended real-value function. Instead of directly minimizing the original composite objective function  $F(\phi) = f(\phi) + g(\phi)$ , the GIST algorithm shown in Algorithm 2 of the supplementary material minimizes a surrogate function:  $\min_{\phi \in \mathcal{S}_\phi} H_r(\phi, \hat{\phi}) = g(\phi) + f(\hat{\phi}) + (\phi - \hat{\phi})^T \nabla_\phi f(\hat{\phi}) + \frac{r}{2} \|\phi - \hat{\phi}\|_2^2$ , which can be simplified as

$$\begin{aligned} & \min_{\mathbf{w}, \boldsymbol{\sigma}, \boldsymbol{\theta}} \frac{\lambda_2}{2} \|\mathbf{w}\|_2^2 + \frac{r}{2} \|\mathbf{w} - \bar{\mathbf{w}}\|_2^2 + \frac{r}{2} \|\boldsymbol{\sigma} - \bar{\boldsymbol{\sigma}}\|_2^2 + \frac{r}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2 \\ & \text{s.t. } \sum_{j=1}^{s_i} d_j^{(i)} \sigma_{\mathcal{F}_j^i, j}^{(i)} = 1 \forall i \in [m], \sigma_{\mathcal{F}_j^i, j}^{(i)} \geq 0 \forall i, j \\ & \sum_{j=1}^m \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} = 1 \forall i \in [d], \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} \geq 0 \forall i, j, \quad (28) \end{aligned}$$

where  $r$  is a step size determined by the GIST algorithm,  $\bar{\mathbf{w}} = \hat{\mathbf{w}} - \frac{1}{r} \nabla_{\mathbf{w}} f(\hat{\mathbf{w}})$ ,  $\bar{\boldsymbol{\sigma}} = \hat{\boldsymbol{\sigma}} - \frac{1}{r} \nabla_{\boldsymbol{\sigma}} f(\hat{\boldsymbol{\sigma}})$ , and  $\bar{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}} - \frac{1}{r} \nabla_{\boldsymbol{\theta}} f(\hat{\boldsymbol{\theta}})$ . The solution of  $\mathbf{w}$  in problem (28) is  $\mathbf{w} = \frac{r}{\lambda_2 + r} \bar{\mathbf{w}}$ . It is easy to see that  $\boldsymbol{\sigma}$  and  $\boldsymbol{\theta}$  in problem (28) are independent and the subproblem with respect to  $\boldsymbol{\sigma}$  can be decomposed into multiple problems each of which has the same formulation as problem (4), leading to an efficient solution. The subproblem with respect to  $\boldsymbol{\theta}$  based on problem (28) is formulated as

$$\begin{aligned} & \min_{\boldsymbol{\theta}} \frac{1}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2 \\ & \text{s.t. } \sum_{j=1}^m \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} = 1 \forall i \in [d], \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} \geq 0 \forall i, j. \quad (29) \end{aligned}$$

Problem (29) is a QP problem and here we devise an optimization method for problem (29) to achieve the speedup. Since different features share the same  $\theta_{k,r}^{(j)}$  when they have the same ancestor in the tree, problem (29) cannot directly reduce to problem (4). By introducing  $\vartheta_i^{(j)}$  as a copy of

$\theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)}$ , we can reformulate problem (29) as

$$\begin{aligned} & \min_{\boldsymbol{\theta}, \boldsymbol{\vartheta}} \frac{1}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2 \\ & \text{s.t. } \sum_{j=1}^m \vartheta_i^{(j)} = 1 \forall i \in [d], \vartheta_i^{(j)} \geq 0 \forall i, j, \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} = \vartheta_i^{(j)} \forall i, j, \quad (30) \end{aligned}$$

where  $\boldsymbol{\vartheta}$  is a vector containing all  $\vartheta_i^{(j)}$ 's. Because of the linear equality constraints in problem (30), we use the ADMM algorithm (Boyd et al. 2011) to solve it. The augmented Lagrangian function is defined as  $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\vartheta}) = \frac{1}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|_2^2 + \sum_{i=1}^d \sum_{j=1}^m \left( \rho_i^{(j)} (\theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} - \vartheta_i^{(j)}) + \frac{\rho}{2} (\theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} - \vartheta_i^{(j)})^2 \right)$ , where  $\{\rho_i^{(j)}\}_{i=1}^m$  act as Lagrangian multipliers and  $\rho$  is a penalty parameter. Then we need to solve

$$\min_{\boldsymbol{\theta}, \boldsymbol{\vartheta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\vartheta}) \quad \text{s.t. } \sum_{j=1}^m \vartheta_i^{(j)} = 1 \forall i \in [d], \vartheta_i^{(j)} \geq 0 \forall i, j. \quad (31)$$

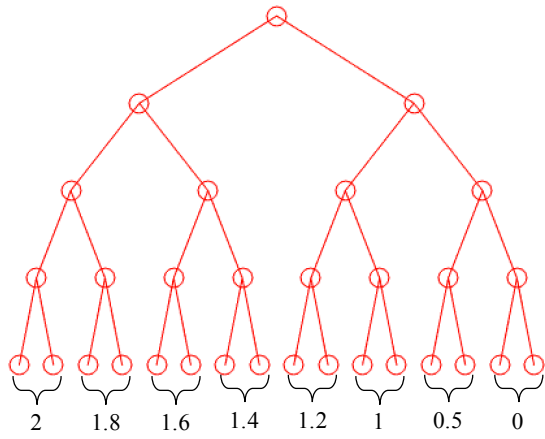
In the ADMM algorithm, problem (31) can be solved alternatively with respect to  $\boldsymbol{\theta}$  and  $\boldsymbol{\vartheta}$ . When  $\boldsymbol{\theta}$  is fixed, the problem with respect to  $\boldsymbol{\vartheta}_i = (\vartheta_i^{(1)}, \dots, \vartheta_i^{(m)})^T$  is formulated as

$$\min_{\boldsymbol{\vartheta}_i} \sum_{j=1}^m \left( \vartheta_i^{(j)} - \tilde{\vartheta}_i^{(j)} \right)^2 \quad \text{s.t. } \sum_{j=1}^m \vartheta_i^{(j)} = 1 \forall i \in [d], \vartheta_i^{(j)} \geq 0 \forall i, j,$$

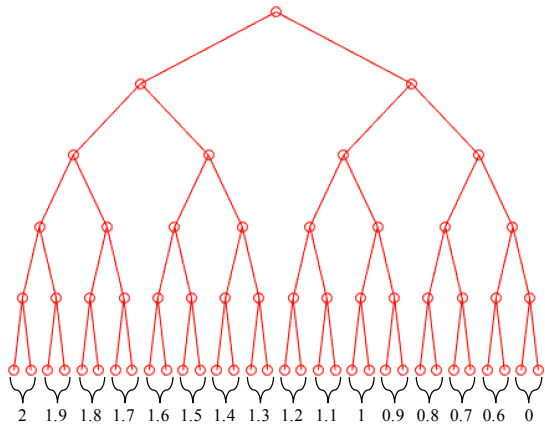
where  $\tilde{\vartheta}_i^{(j)} = \theta_{\mathcal{P}_i^{j-1}, \mathcal{P}_i^j}^{(j)} + \frac{\rho_i^{(j)}}{\rho}$ . This problem takes the same form as problem (4), which makes it have an efficient solution. The subproblem with respect to  $\boldsymbol{\theta}$  is a QP problem without any constraint and by setting the derivative to zero, we can obtain the analytical solution as  $\theta_{k,r}^{(j)} = \frac{1}{1 + \rho d_r^{(j)}} \left( \tilde{\theta}_{k,r}^{(j)} + \sum_{i \in \mathcal{S}_{k,r}^{(j)}} \left( \rho \vartheta_i^{(j)} - \rho_i^{(j)} \right) \right)$ , where  $\mathcal{S}_{k,r}^{(j)} = \{i \mid \mathcal{P}_i^{j-1} = k, \mathcal{P}_i^j = r\}$  with its cardinality equal to  $d_r^{(j)}$ .

## Details in Experiments

The true feature weights for  $\mathbf{w}^*$  when  $m$  equals 4 and 5 are shown in Fig. 3.



(a)  $m = 4$



(b)  $m = 5$

Figure 3: Binary trees generated in synthetic data when  $m$  equals 4 and 5.