

Deriving the Block Coordinate Descent Rules for Tree-Structured NMF with sparsity regularization

Objective

Given $t \in \{1, \dots, T\}$ tasks, each with input matrix $X^{(t)} \in \mathbb{R}^{n_t \times m}$, related to each other in a task hierarchy/tree with a set of nodes $c \in \{r\} \cup \mathcal{B} \cup \mathcal{T}$ where r is the root node, \mathcal{B} a set of internal (or branch) nodes $b \in \mathcal{B}$, and \mathcal{T} a set of the task-specific leaf nodes, the objective is:

$$O = \sum_{t=1}^T \left[\left\| X^{(t)} - U^{(t)} V^{(t)\top} \right\|_F^2 + \lambda \sum_{i=1}^m \left\| V^{(t)}[i, :] \right\|_1 \right] + \alpha \sum_c \left\| V^{(c)} - V^{Pa(c)} \right\|_F^2 \quad (1)$$

where $U^{(t)} \in \mathbb{R}_{\geq 0}^{n_t \times k}$, $V^{(\cdot)} \in \mathbb{R}_{\geq 0}^{m \times k}$, $k \ll n, m$. $V^{(t)}[i, :]$ is the i row of task-specific factor matrix $V^{(t)}$.

The regularization term involving λ tries to enforce sparsity in each row of task-specific $V^{(t)}$, ultimately so that only one latent dimension “lights up” for each row of $V^{(t)}$. Higher λ will enforce stricter sparsity.

The regularization term involving α will:

- a. constrain a task-specific latent feature factor $V^{(t)}$ in a leaf node of the task hierarchy to be similar to $V^{Pa(t)}$ in its parent node;
- b. constrain an internal node’s latent feature factor $V^{(b)}$ to be similar to its direct child nodes’ $V^{(c)}$ and its parent node’s $V^{Pa(b)}$; and
- c. constrain the root node’s latent feature factor $V^{(r)}$ to be similar to all of its direct child nodes’ $V^{(c)}$ s.

Breaking down to task-level and column-level subproblems

The objective can be written as:

$$O = \sum_{t=1}^T \left[\left\| X^{(t)} - \sum_k u_k^{(t)} v_k^{(t)\top} \right\|_F^2 + \lambda \sum_{i=1}^m \sum_k |V^{(t)}[i, k]| \right] + \alpha \sum_c \sum_k \left\| v_k^{(c)} - v_k^{Pa(c)} \right\|_2^2 \quad (2)$$

$$= \sum_{t=1}^T \left[\left\| X^{(t)} - \sum_k u_k^{(t)} v_k^{(t)\top} \right\|_F^2 + \lambda \sum_k \sum_{i=1}^m |V^{(t)}[i, k]| \right] + \alpha \sum_c \sum_k \left\| v_k^{(c)} - v_k^{Pa(c)} \right\|_2^2 \quad (3)$$

$$= \sum_{t=1}^T \left[\left\| X^{(t)} - \sum_k u_k^{(t)} v_k^{(t)\top} \right\|_F^2 + \lambda \sum_k \left\| v_k^{(t)} \right\|_1 \right] + \alpha \sum_c \sum_k \left\| v_k^{(c)} - v_k^{Pa(c)} \right\|_2^2 \quad (4)$$

Where $u_k^{(t)} \in \mathbb{R}^{n_t}$ is the k th column vector of $U^{(t)}$ and $v_k^{(t)} \in \mathbb{R}^m$ is the k th column vector of $V^{(t)}$. Now we ‘pull out’ terms involving the k th column in all factors:

$$O = \sum_{t=1}^T \left[\left\| X^{(t)} - u_k^{(t)} v_k^{(t)\top} - \sum_{j \neq k} u_j^{(t)} v_j^{(t)\top} \right\|_F^2 + \lambda \left\| v_k^{(t)} \right\|_1 + \lambda \sum_{j \neq k} \left\| v_j^{(t)} \right\|_1 \right] \quad (5)$$

$$+ \alpha \sum_c \left(\left\| v_k^{(c)} - u_k^{Pa(c)} \right\|_2^2 + \sum_{j \neq k} \left\| v_j^{(c)} - v_j^{Pa(c)} \right\|_2^2 \right) \quad (6)$$

Now we’ll substitute with $R_k^{(t)} = X^{(t)} - \sum_{j \neq k} u_j^{(t)} v_j^{(t)\top}$:

$$O = \sum_{t=1}^T \left[\left\| R_k^{(t)} - u_k^{(t)} v_k^{(t)\top} \right\|_F^2 + \lambda \left\| v_k^{(t)} \right\|_1 + \lambda \sum_{j \neq k} \left\| v_j^{(t)} \right\|_1 \right] \quad (7)$$

$$+ \alpha \sum_c \left\| v_k^{(c)} - u_k^{Pa(c)} \right\|_2^2 + \alpha \sum_c \sum_{j \neq k} \left\| v_j^{(c)} - v_j^{Pa(c)} \right\|_2^2 \quad (8)$$

We can now attempt to optimize $u_k^{(t)}$ and $v_k^{(\cdot)}$, fixing all other parameters to be constant.

Optimize $v_k^{(t)}$

To find $v_k^{(t)}$ for each leaf node task t that minimizes the objective, we find the derivative of the objective with respect to $v_k^{(t)}$ and set it to 0, then solve. First we expand the objective into matrix multiplications:

$$O = \left\| R_k^{(t)} - u_k^{(t)} v_k^{(t)\top} \right\|_F^2 + \lambda \left\| v_k^{(t)} \right\|_1 + \alpha \left\| v_k^{(t)} - v_k^{Pa(t)} \right\|_2^2 + C \quad (9)$$

$$= \text{Tr} \left[\left(R_k^{(t)} - u_k^{(t)} v_k^{(t)\top} \right)^\top \left(R_k^{(t)} - u_k^{(t)} v_k^{(t)\top} \right) \right] + \lambda \left\| v_k^{(t)} \right\|_1 \quad (10)$$

$$+ \alpha \left(v_k^{(t)} - v_k^{Pa(t)} \right)^\top \left(v_k^{(t)} - v_k^{Pa(t)} \right) + C \quad (11)$$

Here C subsumes all elements of the objective that does not involve $v_k^{(t)}$ (including terms involving tasks other than t), since they will be zeroed out when the derivative is taken with respect to $v_k^{(t)}$. Now we keep expanding:

$$O = \text{Tr} \left[R_k^{(t)\top} R_k^{(t)} - 2R_k^{(t)\top} u_k^{(t)} v_k^{(t)\top} + \left(u_k^{(t)} v_k^{(t)\top} \right)^\top \left(u_k^{(t)} v_k^{(t)\top} \right) \right] \quad (12)$$

$$+ \lambda \left\| v_k^{(t)} \right\|_1 + \alpha \left(v_k^{(t)\top} v_k^{(t)} - 2v_k^{(t)\top} v_k^{Pa(t)} + v_k^{Pa(t)\top} v_k^{Pa(t)} \right) + C \quad (13)$$

$$= \text{Tr} \left(R_k^{(t)\top} R_k^{(t)} \right) - 2 \text{Tr} \left(R_k^{(t)\top} u_k^{(t)} v_k^{(t)\top} \right) + \text{Tr} \left(v_k^{(t)} u_k^{(t)\top} u_k^{(t)} v_k^{(t)\top} \right) \quad (14)$$

$$+ \lambda \left\| v_k^{(t)} \right\|_1 + \alpha v_k^{(t)\top} v_k^{(t)} - 2\alpha v_k^{(t)\top} v_k^{Pa(t)} + \alpha v_k^{Pa(t)\top} v_k^{Pa(t)} + C \quad (15)$$

$$= \text{Tr} \left(R_k^{(t)\top} R_k^{(t)} \right) - 2 \left(R_k^{(t)\top} u_k^{(t)} \right)^\top v_k^{(t)} + \left(u_k^{(t)\top} u_k^{(t)} \right) \left(v_k^{(t)\top} v_k^{(t)} \right) \quad (16)$$

$$+ \lambda \mathbf{1}_m^\top v_k^{(t)} + \alpha v_k^{(t)\top} v_k^{(t)} - 2\alpha v_k^{(t)\top} v_k^{Pa(t)} + \alpha v_k^{Pa(t)\top} v_k^{Pa(t)} + C \quad (17)$$

where $\mathbf{1}_m$ is a vector of size m , filled with 1's. What allows us to expand $\left\| v_k^{(t)} \right\|_1$ from (16) to $\mathbf{1}_m^\top v_k^{(t)}$ (17) is the fact that we're enforcing $v_k^{(t)}$ to be non-negative at initialization and at each iteration.

Now the fun part:

$$\frac{\partial O}{\partial v_k^{(t)}} = 0 - 2R_k^{(t)\top} u_k^{(t)} + 2v_k^{(t)} u_k^{(t)\top} u_k^{(t)} + \lambda \mathbf{1}_m + 2\alpha v_k^{(t)} - 2\alpha v_k^{Pa(t)} + 0 + 0 \quad (18)$$

$$0 = -R_k^{(t)\top} u_k^{(t)} + \left(u_k^{(t)\top} u_k^{(t)} + \alpha \right) v_k^{(t)} + \frac{\lambda}{2} \mathbf{1}_m - \alpha v_k^{Pa(t)} \quad (19)$$

$$v_k^{(t)} = \frac{R_k^{(t)\top} u_k^{(t)} + \alpha v_k^{Pa(t)} - \frac{\lambda}{2} \mathbf{1}_m}{\left\| u_k^{(t)} \right\|_2^2 + \alpha} \quad (20)$$

With the non-negativity constraint $v_k^{(t)} \geq 0$, we want $R_k^{(t)\top} u_k^{(t)} + \alpha v_k^{Pa(t)} - \frac{\lambda}{2} \mathbf{1}_m \geq 0$, because if $R_k^{(t)\top} u_k^{(t)} + \alpha v_k^{Pa(t)} - \frac{\lambda}{2} \mathbf{1}_m < 0$, O will increase in (16) and (17). So the finalized update rule is:

$$v_k^{(t)} = \frac{\left[R_k^{(t)\top} u_k^{(t)} + \alpha v_k^{Pa(t)} - \frac{\lambda}{2} \mathbf{1}_m \right]_+}{\left\| u_k^{(t)} \right\|_2^2 + \alpha} \quad (21)$$

Optimize $u_k^{(t)}$

We can derive the update rule for $u_k^{(t)}$ in leaf node task t similarly but much more simply. From (17), we take the derivative of O_t with respect to $u_k^{(t)}$; all regularization terms will zero out since they do not involve $u_k^{(t)}$. Hence the final update rule for $u_k^{(t)}$ is:

$$u_k^{(t)} = \frac{\left[R_k^{(t)} v_k^{(t)} \right]_+}{\left\| v_k^{(t)} \right\|_2^2} \quad (22)$$

Optimize $v_k^{(r)}$

For the overall consensus factor in the root of the task hierarchy, $v_k^{(r)}$, we can again ignore terms that do not involve $v_k^{(r)}$ in the objective (4). Note that we're going to collect the terms involving nodes c whose parent is the root node, i.e. $\text{Pa}(c) = r$:

$$O = \alpha \sum_{c \in \text{Child}(r)} \left\| v_k^{(c)} - v_k^{(r)} \right\|_2^2 + C \quad (23)$$

$$= \alpha \sum_{c \in \text{Child}(r)} \left(v_k^{(c)} - v_k^{(r)} \right)^\top \left(v_k^{(c)} - v_k^{(r)} \right) + C \quad (24)$$

$$= \alpha \sum_{c \in \text{Child}(r)} \left[v_k^{(c)\top} v_k^{(c)} - 2v_k^{(c)\top} v_k^{(r)} + v_k^{(r)\top} v_k^{(r)} \right] + C \quad (25)$$

$$= C - \sum_{c \in \text{Child}(r)} 2\alpha v_k^{(c)\top} v_k^{(r)} + \sum_{c \in \text{Child}(r)} \alpha v_k^{(r)\top} v_k^{(r)} \quad (26)$$

Now we take the derivative, set to 0, and solve:

$$\frac{\partial O}{\partial v_k^{(r)}} = 0 - \sum_{c \in \text{Child}(r)} 2\alpha v_k^{(c)} + \sum_{c \in \text{Child}(r)} 2\alpha v_k^{(r)} \quad (27)$$

$$0 = - \sum_{c \in \text{Child}(r)} v_k^{(c)} + |\text{Child}(r)| \cdot v_k^{(r)} \quad (28)$$

$$v_k^{(r)} = \frac{\sum_{c \in \text{Child}(r)} v_k^{(c)}}{|\text{Child}(r)|} \quad (29)$$

where $|\text{Child}(r)|$ is the number of direct child nodes of the root node r .

Optimize $v_k^{(b)}$

For the latent feature factor in an internal/branch node of the task hierarchy, $v_k^{(b)}$, same drill as before: we ignore terms that do not involve $v_k^{(b)}$ for the particular node b of interest in the objective (4). This time we collect terms involving the parent node of b , i.e. $\text{Pa}(b)$, and nodes c whose parent is b , i.e. $\text{Pa}(c) = b$:

$$O = \alpha \left(\left\| v_k^{(b)} - v_k^{\text{Pa}(b)} \right\|_2^2 + \sum_{c \in \text{Child}(b)} \left\| v_k^{(c)} - v_k^{(b)} \right\|_2^2 \right) + C \quad (30)$$

$$= \alpha \left(v_k^{(b)} - v_k^{\text{Pa}(b)} \right)^\top \left(v_k^{(b)} - v_k^{\text{Pa}(b)} \right) + \alpha \sum_{c \in \text{Child}(b)} \left(v_k^{(c)} - v_k^{(b)} \right)^\top \left(v_k^{(c)} - v_k^{(b)} \right) + C \quad (31)$$

$$= \alpha \left[v_k^{(b)\top} v_k^{(b)} - 2v_k^{(b)\top} v_k^{\text{Pa}(b)} + v_k^{\text{Pa}(b)\top} v_k^{\text{Pa}(b)} \right] + \alpha \sum_{c \in \text{Child}(b)} \left[v_k^{(c)\top} v_k^{(c)} - 2v_k^{(c)\top} v_k^{(b)} + v_k^{(b)\top} v_k^{(b)} \right] + C \quad (32)$$

$$= \alpha v_k^{(b)\top} v_k^{(b)} - 2\alpha v_k^{(b)\top} v_k^{\text{Pa}(b)} - \sum_{c \in \text{Child}(b)} 2\alpha v_k^{(c)\top} v_k^{(b)} + \sum_{c \in \text{Child}(b)} \alpha v_k^{(b)\top} v_k^{(b)} + C \quad (33)$$

Now we take the derivative, set to 0, and solve:

$$\frac{\partial O}{\partial v_k^{(b)}} = 2\alpha v_k^{(b)} - 2\alpha v_k^{\text{Pa}(b)} - \sum_{c \in \text{Child}(b)} 2\alpha v_k^{(c)} + \sum_{c \in \text{Child}(b)} 2\alpha v_k^{(b)} \quad (34)$$

$$0 = v_k^{(b)} - v_k^{\text{Pa}(b)} - \sum_{c \in \text{Child}(b)} v_k^{(c)} - |\text{Child}(b)| \cdot v_k^{(b)} \quad (35)$$

$$= (1 + |\text{Child}(b)|)v_k^{(b)} - v_k^{\text{Pa}(b)} - \sum_{c \in \text{Child}(b)} v_k^{(c)} \quad (36)$$

$$v_k^{(b)} = \frac{v_k^{\text{Pa}(b)} + \sum_{c \in \text{Child}(b)} v_k^{(c)}}{1 + |\text{Child}(b)|} \quad (37)$$

where $|\text{Child}(b)|$ is the number of direct child nodes of b .