

# Convergence Property for Tensor Regression

Jiixin Hu

05/24/2020

## CONVERGENCE PROPERTY

The convergence of Algorithm 1 is guaranteed because the alternating algorithm monotonically increase the objective function after each iteration. Now, we study the convergence property for the actual estimation sequence  $(\mathcal{C}^{(t)}, \{M_k^{(t)}\})$  and  $\mathcal{B}^{(t)} = \mathcal{C}^{(t)} \times_1 M_1^{(t)} \times_2 \cdots \times_K M_K^{(t)}$ . To simplify the analysis, we set the hyper-parameter  $\alpha$  to infinity and define the notation  $\mathcal{A} = (\mathcal{C}, \{M_k\})$ . Define the Forbenius norm of  $\|\mathcal{A}\|_F = \|\mathcal{C}\|_F + \sum_{k=1}^K \|M_k\|_F$ . We propose following assumptions.

Is it true that the Frobenius norm of  $M_k$  is always  $\sqrt{r_k}$ ?

A1. (Regularity Condition) The log-likelihood function  $\mathcal{L}(\mathcal{A})$  is continuous and the set  $\{\mathcal{A} : \mathcal{L}(\mathcal{A}) \geq \mathcal{L}(\mathcal{A}^{(0)})\}$  is compact.

The Hessian of what function w.r.t. what decision variables?

A2. (Strictly local maximum condition) Each block update in Algorithm is well-defined and the corresponding Hessian is non-singular at the solution.

A3. (Local Uniqueness condition) The set of stationary points of  $\mathcal{L}(\mathcal{A})$  is isolated up to orthogonalization.

A4. (Tensor Lipschitz condition) The tensor representation  $\mathcal{B}(\mathcal{A})$  is tensor Lipschitz at  $\mathcal{A}^*$ . That means there exists two constant  $c_1, c_2 > 0$ , s.t.

Does this assumption hold in our context?

$$c_1 \|\mathcal{A}' - \mathcal{A}''\|_F \leq \|\mathcal{B}(\mathcal{A}') - \mathcal{B}(\mathcal{A}'')\|_F \leq c_2 \|\mathcal{A}' - \mathcal{A}''\|_F$$

What if  $\mathcal{A}', \mathcal{A}''$  differ by orthogonalization?  $\mathcal{B}(\mathcal{A}') - \mathcal{B}(\mathcal{A}'') = 0$ , but  $\mathcal{A}' \neq \mathcal{A}''$ . Consider the following counter-example  
 $\mathcal{A}' = [\mathcal{C}; M_1, \dots, M_K]$   
 $\mathcal{A}'' = [\mathcal{C} \times_1 P; M_1 P^{\wedge \{-1\}}, M_2, \dots, M_K]$

for  $\mathcal{A}', \mathcal{A}''$  are sufficiently close to  $\mathcal{A}^*$ .

These conditions are mild for the tensor of order 3 or higher.

**Proposition 1** (Algorithm 1 Convergence). Suppose A1-A3 holds.

1. (Global Convergence) Any sequence  $\mathcal{A}^{(t)}$  generated by Algorithm 1 converges to a stationary point of  $\mathcal{L}(\mathcal{A})$ .

2. (Local Linear Convergence) Let  $\mathcal{A}^*$  be a local maximizer of  $\mathcal{L}(\mathcal{A})$ . There exists an  $\epsilon$ -neighborhood of  $\mathcal{A}^*$ , such that, for any  $\mathcal{A}^{(0)}$  in the neighborhood, the iterates  $\mathcal{A}^{(t)}$  generated by Algorithm 1 linearly convergent to  $\mathcal{A}^*$ .

$$\|\mathcal{A}^{(t)} - \mathcal{A}^*\|_F \leq \rho^t \|\mathcal{A}^{(0)} - \mathcal{A}^*\|_F$$

where  $\rho \in (0, 1)$  is a contraction parameter. If A4 holds, there exists a constant  $C$  such that

$$\|\mathcal{B}(\mathcal{A}^{(t)}) - \mathcal{B}(\mathcal{A}^*)\|_F \leq C \rho^t \|\mathcal{B}(\mathcal{A}^{(0)}) - \mathcal{B}(\mathcal{A}^*)\|_F.$$

Combining the Proposition and Theorem 4.1, we have the empirical performance of our estimate.

**Theorem 1** (Empirical Performance). Let  $\mathcal{A}^{(t)}$  be a sequence of estimates generated by Algorithm 1 with initial point  $\mathcal{A}^{(0)}$  and limiting point  $\mathcal{A}^*$ . Suppose  $\text{Loss}(\mathcal{B}^{\mathcal{A}^{(0)}}, \mathcal{A}_{true}) = \epsilon$  and  $\mathcal{L}(\mathcal{A}^*) > \mathcal{L}(\mathcal{A}_{true})$ . Suppose A1-A4 holds. With probability  $1 - \exp(-C_1 \sum_k p_k)$ , typo

$$\text{Loss}(\mathcal{B}^{\mathcal{A}^{(t)}}, \mathcal{A}_{true}) \leq C\rho^t \epsilon + C_2 \sum_k p_k,$$

where  $\rho \in (0, 1)$  is a contraction parameter and  $C_1, C_2 > 0$  are two constants.

## PROOF

grammar error

For Global Convergence, we need to show every sub-sequence of  $\mathcal{A}^{(t)}$  accumulates to the same stationary point. Let  $\mathcal{A}^*$  be any limiting point of sub-sequences of  $\mathcal{A}^{(t)}$ . Due to  $\mathcal{L}(\mathcal{A}^{(t)})$  monotonically increases along with  $t$ , every  $\mathcal{A}^*$  is a stationary point of  $\mathcal{L}(\mathcal{A})$  and  $\mathcal{A}^* \in \{\mathcal{A} : \mathcal{L}(\mathcal{A}) \geq \mathcal{L}(\mathcal{A}^{(0)})\}$ . The set of  $\mathcal{A}^*$  is also compact and connected because of A1. The isolation of stationary point also implies there are only finite  $\mathcal{A}^*$ s. Therefore, the set of  $\mathcal{A}^*$  becomes a single point. The Global Convergence holds.

To show the Local Convergence, define the differential mapping  $S : S(\mathcal{A}^{(t)}) = \mathcal{A}^{(t+1)}$ . Let  $H$  be the Hessian matrix of  $\mathcal{L}(\mathcal{A})$  at the local maximum  $\mathcal{A}^*$ . We partition the  $H$ :

$$d^2 \mathcal{L}(\mathcal{A}^*) = d^2 \mathcal{L}(C^*, M_1^*, \dots, M_K^*) = \begin{pmatrix} d_{CC}^2 \mathcal{L} & d_{CM_1}^2 \mathcal{L} & \dots & d_{CM_K}^2 \mathcal{L} \\ d_{M_1 C}^2 \mathcal{L} & d_{M_1 M_1}^2 \mathcal{L} & \dots & d_{M_1 M_K}^2 \mathcal{L} \\ \vdots & \vdots & \ddots & \vdots \\ d_{M_K C}^2 \mathcal{L} & d_{M_K M_1}^2 \mathcal{L} & \dots & d_{M_K M_K}^2 \mathcal{L} \end{pmatrix} = L + D + L^\top,$$

where  $L$  is strictly block lower triangle matrix and  $D$  is the diagonal part. Due to A2, every sub-matrix of Hessian is non-singular and negative definite. Therefore,  $L + D$  is invertible. According to Bezdek(2003), we have  $dS(\mathcal{A}^*) = -(L + d)^{-1}L$ . Let  $\rho$  denote the spectral radius of  $dS(\mathcal{A}^*)$  and  $\rho = \max_i |\lambda_i(-(L + d)^{-1}L)|$  is strictly smaller than 1. According to contraction principle, Algorithm 1 will be linearly convergent if its spectral radius of  $dS(\mathcal{A}^*)$  is smaller than 1. Therefore,

$$d(S(\mathcal{A}^{(t)}), S(\mathcal{A}^*)) \leq \rho^t d(\mathcal{A}^{(0)}, \mathcal{A}^*), \quad d(S(\mathcal{A}), S(\mathcal{A}')) = \|\mathcal{C} - \mathcal{C}'\|_F + \sum_k \|M_k - M'_k\|_F$$

where  $(\mathcal{A}, d)$  is a complete metric space,  $S(\mathcal{A}^*) = \mathcal{A}^*$  and  $\mathcal{A}^{(0)}$  is sufficiently near to  $\mathcal{A}^*$ . With sufficiently large  $t \in \mathbb{N}^+$  and A4, we have

$$\text{Loss}(\mathcal{B}(\mathcal{A}^{(t)}), \mathcal{B}^*) = \|t\mathcal{B}(\mathcal{A}^{(t)}) - \mathcal{B}^*\|_F \leq C\rho^t \|t\mathcal{B}(\mathcal{A}^{(0)}) - \mathcal{B}^*\|_F,$$

where  $C > 0$  is a constant.