

Proof of Theorem 1 in OP-TSOD

Theorem 1. Given $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$, $\mathbf{G} \in \mathbb{R}^{n_1 \times d}$, let $\mathbf{M} = \mathbf{A}^\top \mathbf{G} = \mathbf{Q}_T \mathbf{R}_T$ be the QR decomposition of $\mathbf{A}^\top \mathbf{G}$, where $\mathbf{Q}_T \in \mathbb{R}^{n_2 \times d}$ is an unitary matrix and is calculated by implementing Gram-Schmidt process on matrix $\mathbf{M} = \mathbf{A}^\top \mathbf{G}$, and $\mathbf{R}_T \in \mathbb{R}^{d \times d}$ is an upper triangle matrix, $\mathbf{u}_i \in \mathbb{R}^{n_2}$ and $\mathbf{w}_i \in \mathbb{R}^{n_2}$ denotes the i -th column vector of \mathbf{M} and \mathbf{Q}_T respectively. Then the gradient of $f(\mathbf{G}; \mathbf{A})$ is given by

$$\nabla f(\mathbf{G}; \mathbf{A}) = \mathbf{A} \mathbf{D} \mathbf{W}, \quad (1)$$

where

$$\begin{aligned} \mathbf{D} &= 2(\mathbf{H} \mathbf{A} + \mathbf{A}^\top \mathbf{H}^\top) \mathbf{Q}_T, \mathbf{D} \in \mathbb{R}^{n_2 \times d}; \\ \mathbf{H} &= \mathbf{Q}_T \mathbf{Q}_T^\top \mathbf{A}^\top - \mathbf{A}^\top, \mathbf{H} \in \mathbb{R}^{n_2 \times n_1}; \\ \mathbf{W} &= \left\{ \frac{\partial \mathbf{w}_i}{\partial \mathbf{u}_j} \right\}_{1 \leq i \leq d; 1 \leq j \leq d}, \mathbf{W} \in \mathbb{R}^{d \times d}; \\ \frac{\partial \mathbf{w}_i}{\partial \mathbf{u}_j} &= \begin{cases} 0, & i < j, \\ \frac{1}{\|\mathbf{u}_i - \sum_{k=1}^{i-1} (\mathbf{u}_i, \mathbf{w}_k) \mathbf{w}_k\|}, & i = j, \\ -\frac{\sum_{k=j+1}^i (\mathbf{u}_i, \mathbf{w}_{k-1}) \frac{\partial \mathbf{w}_{k-1}}{\partial \mathbf{u}_j}}{\|\mathbf{u}_i - \sum_{k=1}^i (\mathbf{u}_i, \mathbf{w}_k) \mathbf{w}_k\|}, & i > j. \end{cases} \end{aligned}$$

Eqn. (1) can be further simplified as,

$$\nabla f(\mathbf{G}; \mathbf{A}) = [(\mathbf{A} \mathbf{Q}_T)(\mathbf{A} \mathbf{Q}_T)^\top (\mathbf{A} \mathbf{Q}_T) - (\mathbf{A} \mathbf{A}^\top)(\mathbf{A} \mathbf{Q}_T)] \mathbf{W}. \quad (2)$$

Proof. The gradient $\frac{\partial f(\mathbf{G}; \mathbf{A})}{\partial \mathbf{G}}$ of the objective function $f(\mathbf{G}; \mathbf{A})$ w.r.t. \mathbf{G} can be calculated by using the chain rule,

$$\frac{\partial f}{\partial \mathbf{G}} = \frac{\partial f}{\partial \mathbf{Q}_T} \frac{\partial \mathbf{Q}_T}{\partial \mathbf{M}} \frac{\partial \mathbf{M}}{\partial \mathbf{G}}, \quad (3)$$

where $\mathbf{M} = \mathbf{A}^\top \mathbf{G}$. Thus the next step is to calculate $\frac{\partial f}{\partial \mathbf{Q}_T}$, $\frac{\partial \mathbf{Q}_T}{\partial \mathbf{M}}$ and $\frac{\partial \mathbf{M}}{\partial \mathbf{G}}$ respectively. For easy understanding of the derivation process, we denote $\mathbf{H} = \mathbf{Q}_T \mathbf{Q}_T^\top \mathbf{A}^\top - \mathbf{A}^\top$, and $\mathbf{H} \in \mathbb{R}^{n_2 \times n_1}$.

Step 1: Calculate $\frac{\partial f}{\partial \mathbf{Q}_T}$: Note that $\frac{\partial f}{\partial \mathbf{Q}_T}$ will produce a matrix. Following the trace trick in [4], we can derive the derivative as follow,

$$\begin{aligned}
df &= \text{tr}[df] \\
&= \text{tr}\left[\left(\frac{\partial f}{\partial \mathbf{H}}\right)^\top d\mathbf{H}\right] \\
&= \text{tr}[2\mathbf{H}^\top d[\mathbf{Q}_T \mathbf{Q}_T^\top \mathbf{A}^\top - \mathbf{A}^\top]] \\
&= \text{tr}[2\mathbf{H}^\top d[\mathbf{Q}_T \mathbf{Q}_T^\top \mathbf{A}^\top]] \\
&= \text{tr}[2\mathbf{H}^\top [(d\mathbf{Q}_T) \mathbf{Q}_T^\top \mathbf{A}^\top + \mathbf{Q}_T (d\mathbf{Q}_T) \mathbf{A}^\top]] \\
&= \text{tr}[2\mathbf{A}^\top \mathbf{H}^\top [(d\mathbf{Q}_T) \mathbf{Q}_T^\top + \mathbf{Q}_T (d\mathbf{Q}_T^\top)]] \\
&= \text{tr}[2[\mathbf{A}^\top \mathbf{H}^\top + \mathbf{H} \mathbf{A}](d\mathbf{Q}_T) \mathbf{Q}_T] \\
&= \text{tr}[2\mathbf{Q}_T^\top [\mathbf{A}^\top \mathbf{H}^\top + \mathbf{H} \mathbf{A}](d\mathbf{Q}_T)]
\end{aligned}$$

Then $\frac{\partial f}{\partial \mathbf{Q}_T} = 2(\mathbf{H} \mathbf{A} + \mathbf{A}^\top \mathbf{H}^\top) \mathbf{Q}_T$. We denote $\mathbf{D} = 2(\mathbf{H} \mathbf{A} + \mathbf{A}^\top \mathbf{H}^\top) \mathbf{Q}_T$, then $\frac{\partial f}{\partial \mathbf{Q}_T} = \mathbf{D}$.

Step 2: Calculate $\frac{\partial \mathbf{M}}{\partial \mathbf{G}}$: Note that $\frac{\partial \mathbf{M}}{\partial \mathbf{G}}$ is a matrix-matrix derivative and thus will produce a tensor, which can be calculated by using the vectorization approach [4, 5], i.e. $\frac{\partial \mathbf{M}}{\partial \mathbf{G}} = \frac{\partial \text{vec}(\mathbf{M})}{\partial \text{vec}(\mathbf{G})} = \mathbf{I}_d \otimes \mathbf{A} \in \mathbb{R}^{n_1 d \times n_2 d}$ following [2, 4, 5], where $\text{vec}(\mathbf{M})$ reshapes the matrix $\mathbf{M} \in \mathbb{R}^{n_2 \times d}$ to a $n_2 d \times 1$ vector, \otimes is the Kronecker product [5], $\mathbf{I}_d \in \mathbb{R}^{d \times d}$ is an identity matrix and $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$.

Step 3: Calculate $\frac{\partial \mathbf{Q}_T}{\partial \mathbf{M}}$: Denote $\mathbf{M} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_d]$; $\mathbf{Q}_T = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_d]$. Similar to the calculation of $\frac{\partial \mathbf{M}}{\partial \mathbf{G}}$, $\frac{\partial \mathbf{Q}_T}{\partial \mathbf{M}}$ is also a matrix-matrix derivative. We hence use the same vectorization approach. Since \mathbf{Q}_T is obtained through QR decomposition of \mathbf{M} by the Schmidt orthogonalization process [3], every column vector in \mathbf{Q}_T can be simplified as a combination of column vector in \mathbf{M} and given by the following process:

$$\begin{aligned}
\mathbf{w}_1 &= \frac{\mathbf{u}_1}{\|\mathbf{u}_1\|}; \\
\mathbf{w}_2 &= \frac{\mathbf{u}_2 - (\mathbf{u}_2, \mathbf{w}_1) \mathbf{w}_1}{\|\mathbf{u}_2 - (\mathbf{u}_2, \mathbf{w}_1) \mathbf{w}_1\|}; \\
&\dots
\end{aligned}$$

Then, we can calculate the vector-vector derivative of \mathbf{Q}_T w.r.t. \mathbf{G} as follows,

$$\begin{aligned}
\frac{\partial \mathbf{w}_1}{\partial \mathbf{u}_1} &= \frac{1}{\|\mathbf{u}_1\|}; \frac{\partial \mathbf{w}_1}{\partial \mathbf{u}_i} = 0, i = \{2, 3, \dots, d\}; \\
\frac{\partial \mathbf{w}_2}{\partial \mathbf{u}_1} &= \frac{-(\mathbf{u}_2, \mathbf{w}_1) \frac{\partial \mathbf{w}_1}{\partial \mathbf{u}_1}}{\|\mathbf{u}_2 - (\mathbf{u}_2, \mathbf{w}_1) \mathbf{w}_1\|} \\
\frac{\partial \mathbf{w}_2}{\partial \mathbf{u}_2} &= \frac{1}{\|\mathbf{u}_2 - (\mathbf{u}_2, \mathbf{w}_1) \mathbf{w}_1\|}; \frac{\partial \mathbf{w}_2}{\partial \mathbf{u}_i} = 0, i = \{3, 4, \dots, d\}; \\
\frac{\partial \mathbf{w}_3}{\partial \mathbf{u}_1} &= \frac{-(\mathbf{u}_3, \mathbf{w}_2) \frac{\partial \mathbf{w}_2}{\partial \mathbf{u}_1} - (\mathbf{u}_3, \mathbf{w}_1) \frac{\partial \mathbf{w}_2}{\partial \mathbf{u}_1}}{\|\mathbf{u}_3 - (\mathbf{u}_2, \mathbf{w}_2) \mathbf{w}_2 - (\mathbf{u}_3, \mathbf{w}_1) \mathbf{w}_1\|}; \\
\frac{\partial \mathbf{w}_3}{\partial \mathbf{u}_2} &= \frac{-(\mathbf{u}_3, \mathbf{w}_2) \frac{\partial \mathbf{w}_2}{\partial \mathbf{u}_2}}{\|\mathbf{u}_3 - (\mathbf{u}_2, \mathbf{w}_2) \mathbf{w}_2 - (\mathbf{u}_3, \mathbf{w}_1) \mathbf{w}_1\|}; \\
\frac{\partial \mathbf{w}_3}{\partial \mathbf{u}_3} &= \frac{1}{\|\mathbf{u}_3 - (\mathbf{u}_2, \mathbf{w}_2) \mathbf{w}_2 - (\mathbf{u}_3, \mathbf{w}_1) \mathbf{w}_1\|}; \\
\frac{\partial \mathbf{w}_3}{\partial \mathbf{u}_i} &= 0, i = \{4, 5, \dots, d\}. \\
&\dots\dots
\end{aligned}$$

Combining the above vector-vector derivatives will produce $\frac{\partial \mathbf{Q}_T}{\partial \mathbf{M}} = \frac{\partial \text{vec}(\mathbf{Q}_T)}{\partial \text{vec}(\mathbf{G})} = \mathbf{I}_{n_2} \otimes \mathbf{W} \in \mathbb{R}^{n_2 d \times n_2 d}$ following [1], where \mathbf{I}_{n_2} is an identity matrix, and

$$\begin{aligned}
\mathbf{W} &= \left\{ \frac{\partial \mathbf{w}_i}{\partial \mathbf{u}_j} \right\}_{1 \leq i \leq d; 1 \leq j \leq d}, \mathbf{W} \in \mathbb{R}^{d \times d}; \\
\frac{\partial \mathbf{w}_i}{\partial \mathbf{u}_j} &= \begin{cases} 0, & i < j, \\ \frac{1}{\|\mathbf{u}_i - \sum_{k=1}^{i-1} (\mathbf{u}_i, \mathbf{w}_k) \mathbf{w}_k\|}, & i = j, \\ -\frac{\sum_{k=j+1}^i (\mathbf{u}_i, \mathbf{w}_{k-1}) \frac{\partial \mathbf{w}_{k-1}}{\partial \mathbf{u}_j}}{\|\mathbf{u}_i - \sum_{k=1}^i (\mathbf{u}_i, \mathbf{w}_k) \mathbf{w}_k\|}, & i > j. \end{cases}
\end{aligned}$$

Finally, $\frac{\partial f(\mathbf{G}; \mathbf{A})}{\partial \text{vec}(\mathbf{G})} = (\mathbf{1} \otimes \mathbf{D})(\mathbf{I}_{n_2} \otimes \mathbf{W})(\mathbf{I}_d \otimes \mathbf{A}) \in \mathbb{R}^{1 \times n_1 d}$, where $\text{vec}(\mathbf{G})$ is the vector form of matrix \mathbf{G} . Finally, unstacking the vector form back to the matrix form will produce $\nabla f(\mathbf{G}; \mathbf{A}) = \mathbf{A} \mathbf{D} \mathbf{W} \in \mathbb{R}^{n_1 \times d}$ following [2, 4, 5].

References

1. Abadir, K.M., Magnus, J.R.: Matrix algebra, vol. 1. Cambridge University Press (2005)
2. Bodewig, E.: Matrix calculus. Elsevier (2014)
3. Gander, W.: Algorithms for the qr decomposition. Res. Rep **80**(02), 1251–1268 (1980)
4. Graham, A.: Kronecker products and matrix calculus with applications. Courier Dover Publications (2018)
5. Whitcomb, L.L.: Notes on kronecker products. Available: spray. me. jhu. edu/llw/courses/me530647/kron_1. pdf (2013)