

Gradient of Multiview MDS stress function

1 Some matrix calculus

Let X be an $n \times p$ matrix variable and let $y = f(X)$, where $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$ is a function of X . The derivative of y with respect to X is the $p \times n$ matrix given by

$$\frac{dy}{dX} := \left[\frac{\partial f}{\partial X_{ji}} \right]_{ij},$$

where X_{ij} is the (i, j) entry of X . The gradient of y with respect to X is the $n \times p$ matrix given by

$$\nabla_X y := \left[\frac{\partial f}{\partial X_{ij}} \right]_{ij} = \left(\frac{dy}{dX} \right)^T.$$

The first notation is useful when deriving differentiation rules, the second will be used for optimization.

Note that

$$\begin{aligned} y(X_0 + \Delta X) &= y(X_0) + \sum_{i,j} (\nabla_X y|_{X_0})_{ij} (\Delta X)_{ij} + \text{higher order terms} \\ &= y(X_0) + \text{tr} \left(\frac{dy}{dX}(X_0) \Delta X \right) + \text{higher order terms} \end{aligned},$$

which we can write in terms of differentials as

$$dy = \text{tr} \left(\frac{dy}{dX}(X_0) dX \right).$$

We can then use properties of differentials and the trace function to derive many differentiability rules for matrix calculus.

If $g : \mathbb{R} \rightarrow \mathbb{R}$ and $z = g(f(X))$, the chain rule says that

$$\frac{dz}{dX}(X) = g'(f(X)) \frac{df}{dX}(X).$$

If P is a fixed $p \times p$ matrix and $y(X) = f(XP^T)$, we have

$$\begin{aligned} y(X + \Delta X) &= f((X + \Delta X)P^T) \\ &= f(XP^T + \Delta X P^T) \\ &= f(XP^T) + \text{tr} \left(\frac{df}{dX}(XP^T) \Delta X P^T \right) + \mathcal{O}(\|\Delta X P^T\|^2), \\ &= f(XP^T) + \text{tr} \left(P^T \frac{df}{dX}(X_0 P^T) \Delta X \right) + \mathcal{O}(\|\Delta X\|^2) \end{aligned}$$

so that

$$\frac{d}{dX} (f(XP^T)) = P^T \frac{df}{dX} (XP^T)$$

and

$$\nabla_X (f(XP^T)) = \nabla f(XP^T) P.$$

If we differentiate with respect to P instead, set $y(P) = f(XP^T)$ and note that

$$\begin{aligned} y(P + \Delta P) &= f(X(P + \Delta P)^T) \\ &= f(XP^T) + \text{tr} \left(\frac{df}{dX} (XP^T) X \Delta P^T \right) + \mathcal{O}(\|X \Delta P^T\|^2) \\ &= f(XP^T) + \text{tr} \left(X^T \left(\frac{df}{dX} (XP^T) \right)^T \Delta P \right) + \mathcal{O}(\|\Delta P^T\|^2), \end{aligned}$$

and therefore

$$\frac{d}{dP} (f(XP^T)) = X^T \left(\frac{df}{dX} (XP^T) \right)^T$$

and

$$\nabla_P (f(XP^T)) = (\nabla f(XP^T))^T X.$$

If we write $P = QQ^T$ and differentiate with respect to Q , we have $y(Q) = f(XQQ^T)$ and

$$\begin{aligned} y(Q + \Delta Q) &= f(X(Q + \Delta Q)(Q + \Delta Q)^T) \\ &= f(XQQ^T + X(Q\Delta Q^T + \Delta QQ^T + \Delta Q\Delta Q^T)) \\ &= f(XQQ^T) + \text{tr} \left(\frac{df}{dX} (XQQ^T) X (Q\Delta Q^T + \Delta QQ^T + \Delta Q\Delta Q^T) \right) + \mathcal{O}(\|\Delta Q\|^2) \\ &= f(XQQ^T) + \text{tr} \left(\frac{df}{dX} (XQQ^T) X Q \Delta Q^T \right) + \text{tr} \left(\frac{df}{dX} (XQQ^T) X \Delta Q Q^T \right) + \mathcal{O}(\|\Delta Q\|^2) \\ &= f(XQQ^T) + \text{tr} \left(Q^T X^T \left(\frac{df}{dX} (XQQ^T) \right)^T \Delta Q \right) + \text{tr} \left(Q^T \frac{df}{dX} (XQQ^T) X \Delta Q \right) + \mathcal{O}(\|\Delta Q\|^2) \\ &= f(XQQ^T) + \text{tr} \left(Q^T \left(X^T \left(\frac{df}{dX} (XQQ^T) \right)^T + \frac{df}{dX} (XQQ^T) X \right) \Delta Q \right) + \mathcal{O}(\|\Delta Q\|^2) \end{aligned}$$

so it follows that

$$\frac{d}{dQ} (f(XQQ^T)) = Q^T \left(X^T \left(\frac{df}{dX} (XQQ^T) \right)^T + \frac{df}{dX} (XQQ^T) X \right)$$

and

$$\nabla_Q (f(XQQ^T)) = \left((\nabla f(XQQ^T))^T X + X^T \nabla f(XQQ^T) \right) Q$$

2 Multiview MDS using gradient descent

2.1 Gradients of distance function

Given an $n \times p$ matrix X , containing the coordinates of n points in \mathbb{R}^p , the distance between points i and j is

$$d_{ij}(X) = \|X^T e_i - X^T e_j\|_2 = \|X^T (e_i - e_j)\|_2,$$

where $e_i, e_j \in \mathbb{R}^n$ are the i th and j th (column) basis vectors.

The square distance can be written as

$$\begin{aligned} d_{ij}^2(X) &= \|X^T(e_i - e_j)\|_2^2 \\ &= \text{tr}(X^T(e_i - e_j)(e_i - e_j)^T X) , \\ &= \text{tr}(X^T A_{ij} X) \end{aligned}$$

where

$$A_{ij} := (e_i - e_j)(e_i - e_j)^T.$$

Note that A_{ij} is symmetric. The square distance can be written more compactly as

$$d_{ij}^2(X) = (e_i - e_j)^T X X^T (e_i - e_j),$$

but the first form is easier to work with. Note that

$$\begin{aligned} d\text{tr}(X^T A_{ij} X) &= \text{tr}(d(X^T A_{ij} X)) \\ &= \text{tr}(dX^T A_{ij} X + X^T A_{ij} dX) \\ &= \text{tr}(X^T A_{ij}^T dX + X^T A_{ij} dX) , \\ &= \text{tr}((2X^T A_{ij}) dX) \end{aligned}$$

and so

$$\frac{dd_{ij}^2}{dX}(X) = 2X^T A_{ij}.$$

It then follows that

$$\begin{aligned} \frac{dd_{ij}}{dX}(X) &= \frac{d}{dX} \sqrt{d_{ij}^2(X)} \\ &= \frac{1}{2\sqrt{d_{ij}^2(X)}} \frac{dd_{ij}^2}{dX}(X) \\ &= \frac{1}{d_{ij}(X)} X^T A_{ij} \end{aligned}$$

and that

$$\nabla d_{ij}(X) = \frac{1}{d_{ij}(X)} A_{ij} X$$

If P is a $p \times p$ matrix, we have

$$\begin{aligned} \nabla_X (d_{ij}(X P^T)) &= \nabla_X d_{ij}(X P^T) P \\ &= \left(\frac{1}{d_{ij}(X P^T)} A_{ij} (X P^T) \right) P . \\ &= \frac{1}{d_{ij}(X P^T)} A_{ij} X P^T P \end{aligned}$$

If we differentiate with respect to P instead, we obtain

$$\begin{aligned} \nabla_P (d_{ij}(X P^T)) &= (\nabla d_{ij}(X P^T))^T X \\ &= \left(\frac{1}{d_{ij}(X P^T)} A_{ij} (X P^T) \right)^T X . \\ &= \frac{1}{d_{ij}(X P^T)} P X^T A_{ij} X \end{aligned}$$

Finally, if we set $P = Q Q^T$ and differentiate with respect to Q ,

$$\begin{aligned} \nabla_Q (d_{ij}(X Q Q^T)) &= \left((\nabla d_{ij}(X Q Q^T))^T X + X^T \nabla d_{ij}(X Q Q^T) \right) Q \\ &= \left(\left(\frac{1}{d_{ij}(X Q Q^T)} A_{ij} X Q Q^T \right)^T X + X^T \frac{1}{d_{ij}(X Q Q^T)} A_{ij} X Q Q^T \right) Q . \\ &= \frac{1}{d_{ij}(X Q Q^T)} (Q Q^T X^T A_{ij} X + X^T A_{ij} X Q Q^T) Q \end{aligned}$$

2.2 Gradient of MDS stress

For a fixed $n \times n$ distance matrix D , the MDS stress is defined by

$$\sigma^2(X; D) = \sum_{i < j} (d_{ij}(X) - D_{ij})^2.$$

Its gradient is

$$\begin{aligned} \nabla \sigma^2(X; D) &= \nabla_X \sum_{i < j} (d_{ij}(X) - D_{ij})^2 \\ &= \sum_{i < j} 2(d_{ij}(X) - D_{ij}) \nabla_X (d_{ij}(X) - D_{ij}) \\ &= \sum_{i < j} 2(d_{ij}(X) - D_{ij}) \nabla_X d_{ij}(X) \\ &= \sum_{i < j} 2(d_{ij}(X) - D_{ij}) \frac{1}{d_{ij}(X)} A_{ij} X \\ &= \left(2 \sum_{i < j} \frac{(d_{ij}(X) - D_{ij})}{d_{ij}(X)} A_{ij} \right) X \\ &:= B(X; D) X \end{aligned}$$

where

$$B(X; D) := 2 \sum_{i < j} \frac{(d_{ij}(X) - D_{ij})}{d_{ij}(X)} A_{ij}.$$

2.3 Gradient of MDS stress with fixed projections

If P is a $p \times p$ matrix (such as a projection matrix), then the action of P on the rows of X is the matrix XP^T . This is an $n \times p$ matrix giving the new coordinates (e.g. after projecting). The gradient of the MDS stress function w.r. to X is

$$\begin{aligned} \nabla_X (\sigma^2(XP^T; D)) &= \nabla \sigma^2(XP^T; D) P \\ &= B(XP^T; D) XP^T P. \end{aligned}$$

If $\{(P_k, D_k)\}_{k=1}^K$ are k pairs of $p \times p$ transformations and $n \times n$ distance matrices, then the multiview MDS stress function is

$$\begin{aligned} \sigma_m^2(X; \{(P_k, D_k)\}_{k=1}^K) &:= \sum_{k=1}^K \sigma^2(XP_k^T; D_k) \\ &= \sum_k \sum_{i < j} (d_{ij}(XP_k^T) - (D_k)_{ij})^2, \end{aligned}$$

and its gradient is

$$\begin{aligned} \nabla_X \sigma_m^2(X; \{(P_k, D_k)\}_{k=1}^K) &= \sum_k \nabla_X \sigma^2(XP_k^T; D_k) \\ &= \sum_k B(XP_k^T; D_k) XP_k^T P_k. \end{aligned}$$

2.4 Gradient with respect to transformations

The gradient of the MDS stress function w.r. to P is

$$\begin{aligned} \nabla_P \sigma^2(XP^T; D) &= (\nabla \sigma^2(XP^T; D))^T X \\ &= (B(XP^T; D) XP^T)^T X \\ &= P X^T B(XP^T; D) X \end{aligned}$$

The gradient of the multiview MDS stress function with respect to one of the transformations is

$$\begin{aligned}\nabla_{P_k} \sigma_m^2 \left(X; \{(P_k, D_k)\}_{k=1}^K \right) &= \nabla_{P_k} \sigma^2(X P_k^T; D_k) \\ &= P_k X^T B(X P_k^T; D_k) X.\end{aligned}$$

2.5 Gradient with respect to orthogonal projections

A rank- q , $p \times p$ orthogonal projection matrix is one of the form $P_A = Q Q^T$, where Q is an $p \times q$ orthogonal matrix (that is, its q columns are orthonormal). We want to restrict optimization of multi-MDS stress to this type of transformations. The gradient of the MDS stress function with respect to Q is

$$\begin{aligned}\nabla_Q \sigma^2(X Q Q^T; D) &= \left((\nabla \sigma^2(X Q Q^T; D))^T X + X^T \nabla \sigma^2(X Q Q^T; D) \right) Q \\ &= \left((B(X Q Q^T; D) X Q Q^T)^T X + X^T B(X Q Q^T; D) X Q Q^T \right) Q \\ &= \left(Q Q^T X^T (B(X Q Q^T; D))^T X + X^T B(X Q Q^T; D) X Q Q^T \right) Q\end{aligned}$$

The gradient for the multiview MDS stress function with respect to one of the Q matrices is

$$\nabla_{Q_k} \sigma_m^2 \left(X; \{(Q_k Q_k^T, D_k)\}_{k=1}^K \right) = \left(Q_k Q_k^T X^T (B(X Q_k Q_k^T; D_k))^T X + X^T B(X Q_k Q_k^T; D_k) X Q_k Q_k^T \right) Q_k.$$

The projection of a $p \times p$ matrix B into the subspace of rank- q orthogonal matrices is given by $U I_q V^T$, where $U \Sigma V^T$ is the singular-value decomposition of B and $I_q = [e_1 \cdots e_q 0 \cdots 0]$. That is, the matrix $U I_q V^T$ minimizes $\|U \Sigma V^T - C\|_F^2$ over all q -rank orthogonal matrices C .

Since $\tilde{Q}_k^{(i+1)} = Q_k^{(i)} + \alpha \nabla_{Q_k} \sigma_m^2 \left(X; \left\{ \left(Q_k^{(i)} Q_k^{(i)T}, D_k \right) \right\}_{k=1}^K \right)$ is not guaranteed to be a rank- q orthogonal matrix, We set $Q_k^{(i+1)} = \mathcal{P}_q \left(\tilde{Q}_k^{(i+1)} \right)$, where \mathcal{P}_q is the projection described above.