

Differentiating the Singular Value Decomposition

James Townsend

August 10, 2016

1 The low rank case

Let \mathbf{A} be an $m \times n$ matrix of rank $k \leq \min(m, n)$. Then we may decompose \mathbf{A} as $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$, where \mathbf{U} is $m \times k$, \mathbf{S} is $k \times k$ diagonal, \mathbf{V} is $n \times k$ and the matrices \mathbf{U} and \mathbf{V} satisfy the relation

$$\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_k. \quad (1)$$

In this case the differential of \mathbf{A} may be expressed as

$$d\mathbf{A} = d\mathbf{U}\mathbf{S}\mathbf{V}^\top + \mathbf{U}d\mathbf{S}\mathbf{V}^\top + \mathbf{U}\mathbf{S}d\mathbf{V}^\top. \quad (2)$$

The constraint (1) implies that the differentials $d\mathbf{U}$ and $d\mathbf{V}$ are also constrained: focussing on \mathbf{U} for a moment, taking the differential of (1) gives

$$d\mathbf{U}^\top \mathbf{U} + \mathbf{U}^\top d\mathbf{U} = \mathbf{0}. \quad (3)$$

So the matrix $d\Omega_{\mathbf{U}} = \mathbf{U}^\top d\mathbf{U}$ is *skew-symmetric*. In fact, if we fix an $m \times (m - k)$ matrix \mathbf{U}_\perp such that $[\mathbf{U} \quad \mathbf{U}_\perp]$ is an orthogonal matrix (this could be computed using the Gram-Schmidt process) then we may expand $d\mathbf{U}$ as

$$d\mathbf{U} = \mathbf{U}d\Omega_{\mathbf{U}} + \mathbf{U}_\perp d\mathbf{K}_{\mathbf{U}} \quad (4)$$

where $d\mathbf{K}_{\mathbf{U}}$ is an unconstrained $(m - k) \times k$ matrix. Similarly we may expand $d\mathbf{V}$ as

$$d\mathbf{V} = \mathbf{V}d\Omega_{\mathbf{V}} + \mathbf{V}_\perp d\mathbf{K}_{\mathbf{V}} \quad (5)$$

where $d\Omega_{\mathbf{V}} = \mathbf{V}^\top d\mathbf{V}$ is $k \times k$ skew-symmetric and $d\mathbf{K}_{\mathbf{V}}$ is an $(n - k) \times k$ matrix. See [1] for more detail. Left-multiplying (2) by \mathbf{U}^\top and right-multiplying by \mathbf{V} gives

$$\mathbf{U}^\top d\mathbf{A}\mathbf{V} = d\Omega_{\mathbf{U}}\mathbf{S} + d\mathbf{S} + \mathbf{S}d\Omega_{\mathbf{V}}^\top. \quad (6)$$

Since $d\Omega_{\mathbf{U}}$ and $d\Omega_{\mathbf{V}}$ are skew-symmetric, they have zero diagonal and thus the products $d\Omega_{\mathbf{U}}\mathbf{S}$ and $\mathbf{S}d\Omega_{\mathbf{V}}^\top$ must also have zero diagonal. This means that we can split (6) into two components as follows. Letting $d\mathbf{P} := \mathbf{U}^\top d\mathbf{A}\mathbf{V}$ and using \circ to denote the Hadamard product, the diagonal component of (6) is

$$d\mathbf{S} = \mathbf{I}_k \circ d\mathbf{P} \quad (7)$$

and the off diagonal

$$\bar{\mathbf{I}}_k \circ d\mathbf{P} = d\Omega_{\mathbf{U}}\mathbf{S} - \mathbf{S}d\Omega_{\mathbf{V}} \quad (8)$$

where $\bar{\mathbf{I}}_k$ denotes the $k \times k$ matrix with zero diagonal and ones everywhere else.

Taking the transpose of (8) yields

$$\bar{\mathbf{I}}_k \circ d\mathbf{P}^\top = -\mathbf{S}d\Omega_{\mathbf{U}} + d\Omega_{\mathbf{V}}\mathbf{S}. \quad (9)$$

Now right multiply (8) by \mathbf{S} , left multiply (9) by \mathbf{S} and add. This gives

$$\bar{\mathbf{I}}_k \circ [d\mathbf{P}\mathbf{S} + \mathbf{S}d\mathbf{P}^\top] = d\Omega_{\mathbf{U}}\mathbf{S}^2 - \mathbf{S}^2d\Omega_{\mathbf{U}}, \quad (10)$$

which is solved by

$$d\Omega_{\mathbf{U}} = \mathbf{F} \circ [d\mathbf{P}\mathbf{S} + \mathbf{S}d\mathbf{P}^\top] \quad (11)$$

where $\mathbf{F}_{ij} = \begin{cases} \frac{1}{s_j^2 - s_i^2} & i \neq j \\ 0 & i = j \end{cases}$. By a similar process,

$$d\Omega_{\mathbf{V}} = \mathbf{F} \circ [\mathbf{S}d\mathbf{P} + d\mathbf{P}^\top\mathbf{S}]. \quad (12)$$

Finally, to find $d\mathbf{K}_{\mathbf{U}}$, we left multiply (2) by \mathbf{U}_\perp^\top , which yields

$$\mathbf{U}_\perp^\top d\mathbf{A} = d\mathbf{K}_{\mathbf{U}}\mathbf{S}\mathbf{V}^\top \quad (13)$$

which implies that

$$d\mathbf{K}_{\mathbf{U}} = \mathbf{U}_\perp^\top d\mathbf{A}\mathbf{V}\mathbf{S}^{-1}. \quad (14)$$

By a similar line of reasoning,

$$d\mathbf{K}_{\mathbf{V}} = \mathbf{V}_\perp^\top d\mathbf{A}^\top\mathbf{U}\mathbf{S}^{-1}. \quad (15)$$

All of this derivation can now be combined into formulae for the differentials $d\mathbf{U}$, $d\mathbf{S}$ and $d\mathbf{V}$ in terms of $d\mathbf{A}$, \mathbf{U} , \mathbf{S} and \mathbf{V} . We use the identity $\mathbf{U}_\perp\mathbf{U}_\perp^\top = \mathbf{I} - \mathbf{U}\mathbf{U}^\top$ to eliminate \mathbf{U}_\perp and \mathbf{V}_\perp .

$$d\mathbf{U} = \mathbf{U} \left(\mathbf{F} \circ [\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top\mathbf{U}] \right) + (\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top) d\mathbf{A}\mathbf{V}\mathbf{S}^{-1} \quad (16)$$

$$d\mathbf{S} = \mathbf{I}_k \circ [\mathbf{U}^\top d\mathbf{A}\mathbf{V}] \quad (17)$$

$$d\mathbf{V} = \mathbf{V} \left(\mathbf{F} \circ [\mathbf{S}\mathbf{U}^\top d\mathbf{A}\mathbf{V} + \mathbf{V}^\top d\mathbf{A}^\top\mathbf{U}\mathbf{S}] \right) + (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\top) d\mathbf{A}^\top\mathbf{U}\mathbf{S}^{-1} \quad (18)$$

1.1 Reverse mode AD updates

Suppose we have an objective function $f(\mathbf{x})$ whose gradient we wish to calculate. Use the shorthand $\bar{\cdot} = \nabla \cdot f$ to denote the grad of f with respect to \cdot , so the gradient we are looking for is $\bar{\mathbf{x}}$. Suppose that at some stage during the computation of f , we take a matrix $\mathbf{A}(\mathbf{x})$ and compute its svd $\mathbf{U}(\mathbf{x})\mathbf{S}(\mathbf{x})\mathbf{V}(\mathbf{x})^\top$

We may write

$$df = \text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U}) + \text{tr}(\bar{\mathbf{S}}^\top d\mathbf{S}) + \text{tr}(\bar{\mathbf{V}}^\top d\mathbf{V}). \quad (19)$$

To get the reverse mode AD update, we need to use the formulae (16), (17) and (18), and massage the right hand side into the form $\text{tr}(\bar{\mathbf{A}}^\top d\mathbf{A})$, then $\bar{\mathbf{A}}$ will be what we need for the update. Let us look first at the term $\text{tr}(\bar{\mathbf{S}}^\top d\mathbf{S})$. Using (17), this can be written as

$$\text{tr}(\bar{\mathbf{S}}^\top d\mathbf{S}) = \text{tr}\left(\bar{\mathbf{S}}^\top \left(\mathbf{I}_k \circ \left[\mathbf{U}^\top d\mathbf{A}\mathbf{V}\right]\right)\right) \quad (20)$$

$$= \text{tr}\left(\mathbf{U}^\top d\mathbf{A}\mathbf{V} \left(\mathbf{I}_k \circ \bar{\mathbf{S}}\right)\right) \quad (21)$$

$$= \text{tr}\left(\mathbf{V} \left(\mathbf{I}_k \circ \bar{\mathbf{S}}\right) \mathbf{U}^\top d\mathbf{A}\right) \quad (22)$$

using formula 65 of [2]. The expansion of $\text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U})$ is a little longer...

$$\text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U}) = \text{tr}\left(\bar{\mathbf{U}}^\top \left[\mathbf{U} \left(\mathbf{F} \circ \left[\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}\right]\right) + \left(\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top\right) d\mathbf{A}\mathbf{V}\mathbf{S}^{-1}\right]\right). \quad (23)$$

The right hand side is a sum of two terms. Again using formula 65 of [2] and the fact that $\mathbf{F}^\top = -\mathbf{F}$, the first term is

$$\text{tr}\left(\bar{\mathbf{U}}^\top \mathbf{U} \left(\mathbf{F} \circ \left[\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}\right]\right)\right) = \text{tr}\left(\left[\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}\right] \left(\mathbf{F} \circ \mathbf{U}^\top \bar{\mathbf{U}}\right)\right) \quad (24)$$

$$= \text{tr}\left(\mathbf{V}\mathbf{S} \left(\mathbf{F} \circ \mathbf{U}^\top \bar{\mathbf{U}}\right) \mathbf{U}^\top d\mathbf{A} - \mathbf{V}\mathbf{S} \left(\mathbf{F} \circ \bar{\mathbf{U}}^\top \mathbf{U}\right) \mathbf{U}^\top d\mathbf{A}\right) \quad (25)$$

$$= \text{tr}\left(\mathbf{V}\mathbf{S} \left(\mathbf{F} \circ \left[\mathbf{U}^\top \bar{\mathbf{U}} - \bar{\mathbf{U}}^\top \mathbf{U}\right]\right) \mathbf{U}^\top d\mathbf{A}\right) \quad (26)$$

The second term is more straightforward to deal with

$$\text{tr}\left(\bar{\mathbf{U}}^\top \left(\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top\right) d\mathbf{A}\mathbf{V}\mathbf{S}^{-1}\right) = \text{tr}\left(\mathbf{V}\mathbf{S}^{-1} \bar{\mathbf{U}}^\top \left(\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top\right) d\mathbf{A}\right) \quad (27)$$

and therefore

$$\text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U}) = \text{tr}\left(\mathbf{V} \left[\mathbf{S} \left(\mathbf{F} \circ \left[\mathbf{U}^\top \bar{\mathbf{U}} - \bar{\mathbf{U}}^\top \mathbf{U}\right]\right) \mathbf{U}^\top + \mathbf{S}^{-1} \bar{\mathbf{U}}^\top \left(\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top\right)\right] d\mathbf{A}\right). \quad (28)$$

A similar derivation leads to

$$\text{tr}(\bar{\mathbf{V}}^\top d\mathbf{V}) = \text{tr}\left(\left[\mathbf{V} \left(\mathbf{F} \circ \left[\mathbf{V}^\top \bar{\mathbf{V}} - \bar{\mathbf{V}}^\top \mathbf{V}\right]\right) \mathbf{S} + \left(\mathbf{I}_m - \mathbf{V}\mathbf{V}^\top\right) \bar{\mathbf{V}}\mathbf{S}^{-1}\right] \mathbf{U}^\top d\mathbf{A}\right). \quad (29)$$

Putting all of this together leads to the update equation

$$\bar{\mathbf{A}} = \left[\mathbf{U} \left(\mathbf{F} \circ \left[\mathbf{U}^\top \bar{\mathbf{U}} - \bar{\mathbf{U}}^\top \mathbf{U} \right] \right) \mathbf{S} + \left(\mathbf{I}_m - \mathbf{U} \mathbf{U}^\top \right) \bar{\mathbf{U}} \mathbf{S}^{-1} \right] \mathbf{V}^\top + \quad (30)$$

$$\mathbf{U} \left(\mathbf{I}_k \circ \bar{\mathbf{S}} \right) \mathbf{V}^\top + \mathbf{U} \left[\mathbf{S} \left(\mathbf{F} \circ \left[\mathbf{V}^\top \bar{\mathbf{V}} - \bar{\mathbf{V}}^\top \mathbf{V} \right] \right) \mathbf{V}^\top + \mathbf{S}^{-1} \bar{\mathbf{V}}^\top \left(\mathbf{I}_n - \mathbf{V} \mathbf{V}^\top \right) \right] \quad (31)$$

by taking the transposes of the expressions above and noting that the matrices $\mathbf{F} \circ \left[\mathbf{U}^\top \bar{\mathbf{U}} - \bar{\mathbf{U}}^\top \mathbf{U} \right]$ and $\mathbf{F} \circ \left[\mathbf{V}^\top \bar{\mathbf{V}} - \bar{\mathbf{V}}^\top \mathbf{V} \right]$ are symmetric.

References

- [1] Alan Edelman, Tomás A Arias, and Steven T Smith. The geometry of algorithms with orthogonality constraints. *SIAM journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.
- [2] Thomas P. Minka. Old and new matrix algebra useful for statistics. <http://research.microsoft.com/en-us/um/people/minka/papers/matrix/minka-matrix.pdf>.