

Matrix Calculus - Notes on the Derivative of a Trace

Johannes Traa

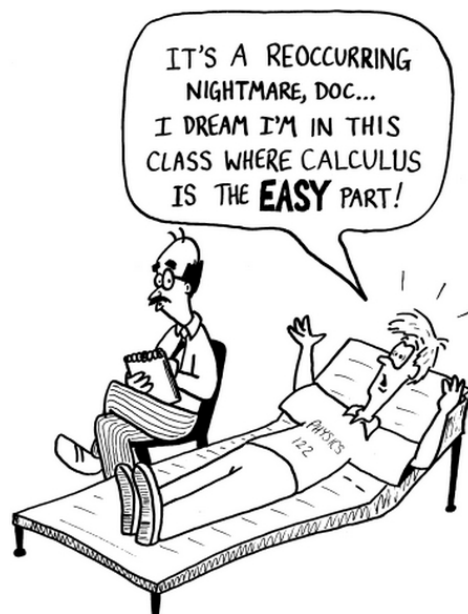
This write-up elucidates the rules of matrix calculus for expressions involving the trace of a function of a matrix \mathbf{X} :

$$f = \text{tr} [g(\mathbf{X})] . \quad (1)$$

We would like to take the derivative of f with respect to \mathbf{X} :

$$\frac{\partial f}{\partial \mathbf{X}} = ? . \quad (2)$$

One strategy is to write the trace expression as a scalar using index notation, take the derivative, and re-write in matrix form. An easier way is to reduce the problem to one or more smaller problems where the results for simpler derivatives can be applied. It's brute-force vs bottom-up.



MATRIX-VALUED DERIVATIVE

The derivative of a scalar f with respect to a matrix $\mathbf{X} \in \mathbb{R}^{M \times N}$ can be written as:

$$\frac{\partial f}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial f}{\partial X_{11}} & \frac{\partial f}{\partial X_{12}} & \cdots & \frac{\partial f}{\partial X_{1N}} \\ \frac{\partial f}{\partial X_{21}} & \frac{\partial f}{\partial X_{22}} & \cdots & \frac{\partial f}{\partial X_{2N}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial X_{M1}} & \frac{\partial f}{\partial X_{M2}} & \cdots & \frac{\partial f}{\partial X_{MN}} \end{bmatrix} \quad (3)$$

So, the result is the same size as \mathbf{X} .

MATRIX AND INDEX NOTATION

It is useful to be able to convert between matrix notation and index notation. For example, the product \mathbf{AB} has elements:

$$[\mathbf{AB}]_{ik} = \sum_j A_{ij} B_{jk} , \quad (4)$$

and the matrix product \mathbf{ABC}^T has elements:

$$[\mathbf{ABC}^T]_{il} = \sum_j A_{ij} [\mathbf{BC}^T]_{jl} = \sum_j A_{ij} \sum_k B_{jk} C_{lk} = \sum_j \sum_k A_{ij} B_{jk} C_{lk} . \quad (5)$$

FIRST-ORDER DERIVATIVES

EXAMPLE 1

Consider this example:

$$f = \text{tr} [\mathbf{AXB}] . \quad (6)$$

We can write this using index notation as:

$$f = \sum_i [\mathbf{AXB}]_{ii} = \sum_i \sum_j A_{ij} [\mathbf{XB}]_{ji} = \sum_i \sum_j A_{ij} \sum_k X_{jk} B_{ki} = \sum_i \sum_j \sum_k A_{ij} X_{jk} B_{ki} . \quad (7)$$

Taking the derivative with respect to X_{jk} , we get:

$$\frac{\partial f}{\partial X_{jk}} = \sum_i A_{ij} B_{ki} = [\mathbf{BA}]_{kj} . \quad (8)$$

The result has to be the same size as \mathbf{X} , so we know that the indices of the rows and columns must be j and k , respectively. This means we have to transpose the result above to write the derivative in matrix form as:

$$\frac{\partial \text{tr} [\mathbf{AXB}]}{\partial \mathbf{X}} = \mathbf{A}^T \mathbf{B}^T . \quad (9)$$

EXAMPLE 2

Similarly, we have:

$$f = \text{tr} [\mathbf{A}\mathbf{X}^T\mathbf{B}] = \sum_i \sum_j \sum_k A_{ij} X_{kj} B_{ki} , \quad (10)$$

so that the derivative is:

$$\frac{\partial f}{\partial X_{kj}} = \sum_i A_{ij} B_{ki} = [\mathbf{B}\mathbf{A}]_{kj} , \quad (11)$$

The \mathbf{X} term appears in (10) with indices kj , so we need to write the derivative in matrix form such that k is the row index and j is the column index. Thus, we have:

$$\frac{\partial \text{tr} [\mathbf{A}\mathbf{X}^T\mathbf{B}]}{\partial \mathbf{X}} = \mathbf{B}\mathbf{A} . \quad (12)$$

MULTIPLE-ORDER

Now consider a more complicated example:

$$f = \text{tr} [\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}\mathbf{C}^T] \quad (13)$$

$$= \sum_i \sum_j \sum_k \sum_l \sum_m A_{ij} X_{jk} B_{kl} X_{lm} C_{im} . \quad (14)$$

The derivative has contributions from both appearances of \mathbf{X} .

TAKE 1

In index notation:

$$\frac{\partial f}{\partial X_{jk}} = \sum_i \sum_l \sum_m A_{ij} B_{kl} X_{lm} C_{im} = [\mathbf{B}\mathbf{X}\mathbf{C}^T\mathbf{A}]_{kj} , \quad (15)$$

$$\frac{\partial f}{\partial X_{lm}} = \sum_i \sum_j \sum_k A_{ij} X_{jk} B_{kl} C_{im} = [\mathbf{C}^T\mathbf{A}\mathbf{X}\mathbf{B}]_{ml} . \quad (16)$$

Transposing appropriately and summing the terms together, we have:

$$\frac{\partial \text{tr} [\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}\mathbf{C}^T]}{\partial \mathbf{X}} = \mathbf{A}^T \mathbf{C}\mathbf{X}^T \mathbf{B}^T + \mathbf{B}^T \mathbf{X}^T \mathbf{A}^T \mathbf{C} . \quad (17)$$

TAKE 2

We can skip this tedious process by applying (9) for each appearance of \mathbf{X} :

$$\frac{\partial \text{tr} [\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}\mathbf{C}^T]}{\partial \mathbf{X}} = \frac{\partial \text{tr} [\mathbf{A}\mathbf{X}\mathbf{D}]}{\partial \mathbf{X}} + \frac{\partial \text{tr} [\mathbf{E}\mathbf{X}\mathbf{C}^T]}{\partial \mathbf{X}} = \mathbf{A}^T \mathbf{D}^T + \mathbf{E}^T \mathbf{C} . \quad (18)$$

where $\mathbf{D} = \mathbf{BXC}^T$ and $\mathbf{E} = \mathbf{AXB}$. So we just evaluate the matrix derivative for each appearance of \mathbf{X} assuming that everything else is a constant (including other \mathbf{X} 's). To see why this rule is useful, consider the following beast:

$$f = \text{tr} [\mathbf{AXX}^T \mathbf{BCX}^T \mathbf{XC}] . \quad (19)$$

We can immediately write down the derivative using (9) and (12):

$$\frac{\partial \text{tr} [\mathbf{AXX}^T \mathbf{BCX}^T \mathbf{XC}]}{\partial \mathbf{X}} = (\mathbf{A})^T (\mathbf{X}^T \mathbf{BCX}^T \mathbf{XC})^T + (\mathbf{BCX}^T \mathbf{XC}) (\mathbf{AX}) + (\mathbf{XC}) (\mathbf{AXX}^T \mathbf{BC}) + (\mathbf{AXX}^T \mathbf{BCX}^T)^T (\mathbf{C})^T \quad (20)$$

$$= \mathbf{AC}^T \mathbf{X}^T \mathbf{XC}^T \mathbf{B}^T \mathbf{X} + \mathbf{BCX}^T \mathbf{XCAX} + \mathbf{XCAXX}^T \mathbf{BC} + \mathbf{XC}^T \mathbf{B}^T \mathbf{XX}^T \mathbf{A}^T \mathbf{C}^T . \quad (21)$$

FROBENIUS NORM

The Frobenius norm shows up when we have an optimization problem involving a matrix factorization and we want to minimize a sum-of-squares error criterion:

$$f = \sum_i \sum_k \left(X_{ik} - \sum_j W_{ij} H_{jk} \right)^2 = \|\mathbf{X} - \mathbf{WH}\|_F^2 = \text{tr} [(\mathbf{X} - \mathbf{WH})(\mathbf{X} - \mathbf{WH})^T] . \quad (22)$$

We can work with the expression in index notation, but it's easier to work directly with matrices and apply the results derived earlier. Suppose we want to find the derivative with respect to \mathbf{W} . Expanding the matrix outer product, we have:

$$f = \text{tr} [\mathbf{XX}^T] - \text{tr} [\mathbf{XH}^T \mathbf{W}^T] - \text{tr} [\mathbf{WHX}^T] + \text{tr} [\mathbf{WHH}^T \mathbf{W}^T] . \quad (23)$$

Applying (9) and (12), we easily deduce that:

$$\frac{\partial \text{tr} [(\mathbf{X} - \mathbf{WH})(\mathbf{X} - \mathbf{WH})^T]}{\partial \mathbf{W}} = -2\mathbf{XH}^T + 2\mathbf{WHH}^T . \quad (24)$$

LOG

Consider this trace expression:

$$f = \text{tr} [\mathbf{V}^T \log(\mathbf{AXB})] = \sum_i \sum_j V_{ij} \log \left(\sum_m \sum_n A_{im} X_{mn} B_{nj} \right) . \quad (25)$$

Taking the derivative with respect to X_{mn} , we get:

$$\frac{\partial f}{\partial X_{mn}} = \sum_i \sum_j \left(\frac{V_{ij}}{\sum_m \sum_n A_{im} X_{mn} B_{nj}} \right) A_{im} B_{nj} = \sum_i \sum_j A_{im} \left(\frac{\mathbf{V}}{\mathbf{AXB}} \right)_{ij} B_{nj} . \quad (26)$$

Thus:

$$\frac{\partial \operatorname{tr} [\mathbf{V}^T \log(\mathbf{A}\mathbf{X}\mathbf{B})]}{\partial \mathbf{X}} = \mathbf{A}^T \left(\frac{\mathbf{V}}{\mathbf{A}\mathbf{X}\mathbf{B}} \right) \mathbf{B}^T . \quad (27)$$

Similarly:

$$\frac{\partial \operatorname{tr} [\mathbf{V}^T \log(\mathbf{A}\mathbf{X}^T\mathbf{B})]}{\partial \mathbf{X}} = \mathbf{B} \left(\frac{\mathbf{V}}{\mathbf{A}\mathbf{X}^T\mathbf{B}} \right)^T \mathbf{A} . \quad (28)$$

These bare a spooky resemblance to (9) and (12).

DIAG

EXAMPLE 1

Consider the tricky case of a $\operatorname{diag}(-)$ operator:

$$f = \operatorname{tr} [\mathbf{A} \operatorname{diag}(\mathbf{x}) \mathbf{B}] = \sum_i \sum_j A_{ij} x_j B_{ji} . \quad (29)$$

Taking the derivative, we have:

$$\frac{\partial f}{\partial x_j} = \sum_i A_{ij} B_{ji} = [(\mathbf{A}^T \odot \mathbf{B}) \mathbf{1}]_j . \quad (30)$$

So we can write:

$$\frac{\partial \operatorname{tr} [\mathbf{A} \operatorname{diag}(\mathbf{x}) \mathbf{B}]}{\partial \mathbf{x}} = (\mathbf{A}^T \odot \mathbf{B}) \mathbf{1} . \quad (31)$$

EXAMPLE 2

Consider the following take on the last example:

$$f = \operatorname{tr} [J \mathbf{A} \operatorname{diag}(\mathbf{x}) \mathbf{B}] = \sum_i \sum_j \sum_k A_{ij} x_j B_{jk} , \quad (32)$$

where J is the matrix of ones.

TAKE 1

Taking the derivative, we have:

$$\frac{\partial f}{\partial x_j} = \sum_i \sum_k A_{ij} B_{jk} = \left(\sum_i A_{ij} \right) \left(\sum_k B_{jk} \right) = [\mathbf{A}^T \mathbf{1} \odot \mathbf{B} \mathbf{1}]_j . \quad (33)$$

So we can write:

$$\frac{\partial \operatorname{tr} [J \mathbf{A} \operatorname{diag}(\mathbf{x}) \mathbf{B}]}{\partial \mathbf{x}} = \mathbf{A}^T \mathbf{1} \odot \mathbf{B} \mathbf{1} . \quad (34)$$

TAKE 2

We could have derived this result from the previous example using the rotation property of the trace operator:

$$f = \text{tr} [\mathbf{J} \mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B}] = \text{tr} [\mathbf{1} \mathbf{1}^T \mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B}] = \text{tr} [\mathbf{1}^T \mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B} \mathbf{1}] = \text{tr} [\mathbf{a}^T \text{diag}(\mathbf{x}) \mathbf{b}] , \quad (35)$$

where we have defined $\mathbf{a} = \mathbf{A}^T \mathbf{1}$ and $\mathbf{b} = \mathbf{B} \mathbf{1}$. Applying (31), we have:

$$\frac{\partial \text{tr} [\mathbf{a}^T \text{diag}(\mathbf{x}) \mathbf{b}]}{\partial \mathbf{x}} = (\mathbf{a} \odot \mathbf{b}) \mathbf{1} = \mathbf{A}^T \mathbf{1} \odot \mathbf{B} \mathbf{1} . \quad (36)$$

EXAMPLE 3

Consider a more complicated example:

$$f = \text{tr} [\mathbf{V}^T \log(\mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B})] = \sum_i \sum_j V_{ij} \log \left(\sum_m A_{im} x_m B_{mj} \right) . \quad (37)$$

Taking the derivative with respect to x_m , we have:

$$\frac{\partial f}{\partial x_m} = \sum_i \sum_j \left(\frac{V_{ij}}{\sum_m A_{im} x_m B_{mj}} \right) A_{im} B_{mj} \quad (38)$$

$$= \sum_j \left[\sum_i A_{im} \left(\frac{\mathbf{V}}{\mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B}} \right)_{ij} \right] B_{mj} \quad (39)$$

$$= \left[\left(\mathbf{A}^T \left(\frac{\mathbf{V}}{\mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B}} \right) \odot \mathbf{B} \right) \mathbf{1} \right]_m . \quad (40)$$

The final result is:

$$\frac{\partial \text{tr} [\mathbf{V}^T \log(\mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B})]}{\partial \mathbf{x}} = \left(\mathbf{A}^T \left(\frac{\mathbf{V}}{\mathbf{A} \text{diag}(\mathbf{x}) \mathbf{B}} \right) \odot \mathbf{B} \right) \mathbf{1} . \quad (41)$$

EXAMPLE 4

How about when we have a trace composed of a sum of expressions, each of which depends on what row of a matrix \mathbf{B} is chosen:

$$f = \text{tr} \left[\sum_k \mathbf{V}^T \log(\mathbf{A} \text{diag}(\mathbf{B}_k; \mathbf{X}) \mathbf{C}) \right] = \sum_k \sum_i \sum_j V_{ij} \log \left(\sum_m A_{im} \left(\sum_n B_{kn} X_{nm} \right) C_{mj} \right) . \quad (42)$$

Taking the derivative, we get:

$$\frac{\partial f}{\partial X_{nm}} = \sum_k \sum_i \sum_j \left(\frac{V_{ij}}{\sum_m A_{im} \left(\sum_n B_{kn} X_{nm} \right) C_{mj}} \right) A_{im} B_{kn} C_{mj} \quad (43)$$

$$= \sum_k B_{kn} \sum_i \sum_j \left(\frac{\mathbf{V}}{\mathbf{A} \operatorname{diag}(\mathbf{B}_k; \mathbf{X}) \mathbf{C}} \right)_{ij} A_{im} C_{mj} \quad (44)$$

$$= \sum_k B_{kn} \left[\mathbf{1}^T \left(\left(\frac{\mathbf{V}}{\mathbf{A} \operatorname{diag}(\mathbf{B}_k; \mathbf{X}) \mathbf{C}} \right)^T \mathbf{A} \odot \mathbf{C}^T \right) \right]_{km}, \quad (45)$$

so we can write:

$$\frac{\partial \operatorname{tr} \left[\sum_k \mathbf{V}^T \log(\mathbf{A} \operatorname{diag}(\mathbf{B}_k; \mathbf{X}) \mathbf{C}) \right]}{\partial \mathbf{X}} = \mathbf{B}^T \begin{bmatrix} \mathbf{1}^T \left(\left(\frac{\mathbf{V}}{\mathbf{A} \operatorname{diag}(\mathbf{B}_1; \mathbf{X}) \mathbf{C}} \right)^T \mathbf{A} \odot \mathbf{C}^T \right) \\ \vdots \\ \mathbf{1}^T \left(\left(\frac{\mathbf{V}}{\mathbf{A} \operatorname{diag}(\mathbf{B}_K; \mathbf{X}) \mathbf{C}} \right)^T \mathbf{A} \odot \mathbf{C}^T \right) \end{bmatrix}. \quad (46)$$

CONCLUSION

We've learned two things. First, if we don't know how to find the derivative of an expression using matrix calculus directly, we can always fall back on index notation and convert back to matrices at the end. This reduces a potentially unintuitive matrix-valued problem into one involving scalars, which we are used to. Second, it's less painful to massage an expression into a familiar form and apply previously-derived identities.