

ML & LinAlg Math Cheat Sheet

October 24, 2017

Contents

1	Notation	1
2	Derivative	1
2.a	Vector Gradient	1
3	Determinant Operator	2
3.a	Random Properties	2
4	Trace Operator	2
4.a	Derivatives	2
4.a.i	$\nabla_{\mathbf{x}} \text{tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A}) = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$	2
4.b	Relation to Determinant	2
A	Proofs	2
A.a	Trace	2
A.a.i	$\nabla_{\mathbf{x}} \text{tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A}) = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$	2

1 Notation

Vectors are column vectors denoted by lower-case bolded variables, such that

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix}.$$

A row vector is denoted $\mathbf{x}^\top = [x_1 \dots x_N]$. A matrix is indicated by a bolded upper-case variable, such that an $N \times M$ matrix is

$$\mathbf{A} = \{a_{ij}\} = [\mathbf{a}_1 \dots \mathbf{a}_M] = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_N^\top \end{bmatrix} = \begin{bmatrix} a_{1,1}^\top & \dots & a_{1,M}^\top \\ \vdots & \ddots & \vdots \\ a_{N,1}^\top & \dots & a_{N,M}^\top \end{bmatrix}.$$

2 Derivative

2.a Vector Gradient

$$\nabla_{\mathbf{x}} \mathbf{y} = \left[\frac{\partial \mathbf{y}}{\partial x_1}, \dots, \frac{\partial \mathbf{y}}{\partial x_N} \right] \quad (1)$$

3 Determinant Operator

3.a Random Properties

For scalar c and $N \times N$ identity matrix I ,

$$\det(cI) = c^N.$$

4 Trace Operator

4.a Derivatives

$$\text{4.a.i} \quad \nabla_{\mathbf{x}} \text{tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A}) = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$$

For square matrix \mathbf{A} . Note that $\mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top) = 2\mathbf{x}^\top \mathbf{A}$ for symmetric \mathbf{A} .

See appendix [A.a.i](#) for proof.

4.b Relation to Determinant

A Proofs

A.a Trace

$$\text{A.a.i} \quad \nabla_{\mathbf{x}} \text{tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A}) = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$$

This proof can likely be generalized to non-square matrixes (and possibly some communicativeness, given the flexibility afforded by the trace), but the restricted case is presented here.

For square $N \times N$ matrix \mathbf{A} ,

$$\nabla_{\mathbf{x}} \text{tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A}) = \frac{d}{d\mathbf{x}} \text{tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A}) = \frac{d}{d\mathbf{x}} \sum_i^N \sum_k^N x_i x_k a_{ik}.$$

Recall eq. (1), and consider for any $j \in \{1, \dots, N\}$:

$$\begin{aligned} \frac{\partial}{\partial x_j} \sum_i^N \sum_k^N x_i x_k a_{ik} &= [x_1 a_{1,j} + x_2 a_{2,j} + \dots + x_{j-1} a_{j-1,j} + x_{j+1} a_{j+1,j} + \dots + x_N a_{N,j}] \\ &\quad + \frac{\partial}{\partial x_j} \sum_k^N x_j x_k a_{jk} \\ &= \left[\sum_i^N x_i a_{ij} - x_j a_{jj} \right] + \sum_k^N x_k a_{jk} - x_j a_{jj} + \frac{\partial}{\partial x_j} x_j x_j a_{jj} \\ &= \sum_i^N x_i a_{ij} + \sum_k^N x_k a_{jk} - 2x_j a_{jj} + 2x_j a_{jj} \\ &= \mathbf{x}^\top \mathbf{a}_j + \mathbf{x}^\top [\mathbf{a}^\top]_j, \end{aligned}$$

where $[\mathbf{a}^\top]_j$ is the j th column of \mathbf{A}^\top .

This equally applies for any j in $1 \dots N$, and so for the full gradient:

$$\begin{aligned} \nabla_{\mathbf{x}} \text{tr}(\mathbf{x}\mathbf{x}^\top \mathbf{A}) &= \frac{d}{d\mathbf{x}} \sum_i^N \sum_k^N x_i x_k a_{ik} = [\mathbf{x}^\top \mathbf{a}_1 \dots \mathbf{x}^\top \mathbf{a}_N] + [\mathbf{x}^\top [\mathbf{a}^\top]_1 \dots \mathbf{x}^\top [\mathbf{a}^\top]_N] \\ &= \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top). \end{aligned}$$