# ML & LinAlg Math Cheat Sheet

November 15, 2017

## Contents

## 1 Notation

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

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

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

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

For some random variable $x$, let $\mathbb{E}[x]$ denote its expected value.

## 2 Derivative

### 2.a Vector Gradient

$$\nabla_{\boldsymbol{x}}\boldsymbol{y} = [\frac{\partial \boldsymbol{y}}{\partial x_1}, \ldots, \frac{\partial \boldsymbol{y}}{\partial x_N}] \tag{1}$$

## 3 Determinant Operator

### 3.a Random Properties

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

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

## 4 Trace Operator

Defined for $N \times N$ square matrix $\boldsymbol{A}$ as

$$\mathrm{tr}(\boldsymbol{A}) \overset{\mathrm{def}}{=} \sum_i^N a_{ii} \tag{2}$$

### 4.a Properties

**4.a.i** $\mathrm{tr}(c\boldsymbol{A} + d\boldsymbol{B}) = c\,\mathrm{tr}(\boldsymbol{A}) + d\,\mathrm{tr}(\boldsymbol{B})$

For scalars $c$ and $d$, square matrices $\boldsymbol{A}$ and $\boldsymbol{B}$.

**4.a.ii** $\mathrm{tr}(\boldsymbol{AB}) = \mathrm{tr}(\boldsymbol{BA}) = \mathrm{tr}(\boldsymbol{A}^\top\boldsymbol{B}) = \mathrm{tr}(\boldsymbol{AB}^\top) = \sum_{i,j} a_{ij}b_{ij}$

And clearly, also $\mathrm{tr}(\boldsymbol{B}^\top\boldsymbol{A}) = \mathrm{tr}(\boldsymbol{BA}^\top) = \sum_{i,j} a_{ij}b_{ij} = \mathrm{tr}(\boldsymbol{AB})$.

### 4.b Derivatives

**4.b.i** $\nabla_{\boldsymbol{x}}\,\mathrm{tr}(\boldsymbol{x}\boldsymbol{x}^\top\boldsymbol{A}) = \boldsymbol{x}^\top(\boldsymbol{A} + \boldsymbol{A}^\top)$

For square matrix $\boldsymbol{A}$. Note that $\boldsymbol{x}^\top(\boldsymbol{A} + \boldsymbol{A}^\top) = 2\boldsymbol{x}^\top\boldsymbol{A}$ for symmetric $\boldsymbol{A}$.
See appendix A.a.i for proof.

### 4.c Relation to Determinant

## 5 Expected Values

For $\boldsymbol{x} \in \mathbb{R}^d$, with expected value $\mathbb{E}[\boldsymbol{x}] = \boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma} = \mathbb{E}[(\boldsymbol{x} - \mathbb{E}[\boldsymbol{x}])(\boldsymbol{x} - \mathbb{E}[\boldsymbol{x}])^\top]$,

$$\mathbb{E}[x_i^2] = \Sigma_{i,i} + \mu_i^2 \tag{3}$$

$$\mathbb{E}_x\left[(y - \boldsymbol{x}^\top\boldsymbol{w})^2\right] = (y - \boldsymbol{w}^\top\boldsymbol{\mu})^2 + \boldsymbol{w}^\top\boldsymbol{\Sigma}\boldsymbol{w}. \tag{4}$$

# A  Proofs

## A.a  Trace

### A.a.i  $\nabla_{\boldsymbol{x}} \operatorname{tr}(\boldsymbol{x}\boldsymbol{x}^\top \boldsymbol{A}) = \boldsymbol{x}^\top(\boldsymbol{A} + \boldsymbol{A}^\top)$

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 $\boldsymbol{A}$,

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

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

$$\frac{\partial}{\partial x_j} \sum_i^N \sum_k^N x_i x_k a_{ik} = [x_1 a_{1,j} + x_2 a_{2,j} + \cdots + x_{j-1} a_{j-1,j} + x_{j+1} a_{j+1,j} + \cdots x_N a_{N,j}]$$

$$+ \frac{\partial}{\partial x_j} \sum_k^N x_j x_k a_{jk}$$

$$= \left[ \sum_i^N x_i a_{ij} - x_j a_{jj} \right] + \sum_k^N x_k a_{jk} - x_j a_{jj} + \frac{\partial}{\partial x_j} x_j x_j a_{jj}$$

$$= \sum_i^N x_i a_{ij} + \sum_k^N x_k a_{jk} - 2 x_j a_{jj} + 2 x_j a_{jj}$$

$$= \boldsymbol{x}^\top \boldsymbol{a}_j + \boldsymbol{x}^\top [\boldsymbol{a}^\top]_j,$$

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

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

$$\nabla_{\boldsymbol{x}} \operatorname{tr}(\boldsymbol{x}\boldsymbol{x}^\top \boldsymbol{A}) = \frac{d}{d\boldsymbol{x}} \sum_i^N \sum_k^N x_i x_k a_{ik} = [\boldsymbol{x}^\top \boldsymbol{a}_1 \cdots \boldsymbol{x}^\top \boldsymbol{a}_N] + [\boldsymbol{x}^\top [\boldsymbol{a}^\top]_1 \cdots \boldsymbol{x}^\top [\boldsymbol{a}^\top]_N]$$

$$= \boldsymbol{x}^\top(\boldsymbol{A} + \boldsymbol{A}^\top).$$