

# The Tensor Cookbook

Thomas Dybdahl Ahle

May 12, 2024

# Chapter 1

## Introduction

**What is this?** These pages are a guide to tensors, using the visual language of “tensor diagrams”. For illustrating the generality of the approach, I’ve tried to closely follow the legendary “Matrix Cookbook”. As such, most of the presentation is a collection of facts (identities, approximations, inequalities, relations, ...) about tensors and matters relating to them. You won’t find many results not in the original cookbook, but hopefully the diagrams will give you a new way to understand and appreciate them.

**It’s ongoing:** The Matrix Cookbook is a long book, and not all the sections are equally amenable to diagrams. Hence I’ve opted to skip certain sections and shorten others. Perhaps in the future, I, or others, will expand the coverage further.

For example, while we cover all of the results on Expectation of Linear Combinations and Gaussian moments, we skip the section on general multi-variate distributions. I have also had to rearrange the material a bit, to avoid having to introduce all the notation up front.

**Complex Matrices and Covariance** Tensor diagrams (or networks) are currently most often seen in Quantum Physics. Here most values are complex numbers, which introduce some extra complexity. In particular transposing a matrix now involves taking the conjugate (flipping the sign of the imaginary part), which introduces the need for co- and contra-variant tensors. None of this complexity is present with standard real valued matrices, as is common e.g. in Machine Learning applications. For simplicity I have decided to not include these complexities.

**Tensorgrad** The symbolic nature of tensor diagrams make the well suited for symbolic computation.

**Advantages of Tensor Diagram Notation:** Tensor diagram notation has many benefits compared to other notations:

Various operations, such as a trace, tensor product, or tensor contraction can be expressed simply without extra notation. Names of indices and tensors can often be omitted. This saves time and lightens the notation, and is especially useful for internal indices which exist mainly to be summed over. The order of the tensor resulting from

a complicated network of contractions can be determined by inspection: it is just the number of unpaired lines. For example, a tensor network with all lines joined, no matter how complicated, must result in a scalar.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Notation and Nomenclature . . . . .	5
1.2	Basics . . . . .	5
1.2.1	Trace . . . . .	5
<b>2</b>	<b>Simple Derivatives</b>	<b>7</b>
2.1	Derivatives of Matrices, Vectors and Scalar Forms . . . . .	7
2.1.1	First Order . . . . .	7
2.1.2	Second Order . . . . .	8
2.1.3	Higher-order and non-linear . . . . .	8
2.1.4	Gradient and Hessian . . . . .	8
2.2	Derivatives of Traces . . . . .	9
2.2.1	First Order . . . . .	9
2.2.2	Second Order . . . . .	9
2.2.3	Higher Order . . . . .	9
2.2.4	Other . . . . .	9
<b>3</b>	<b>Statistics and Probability</b>	<b>10</b>
3.0.1	Definition of Moments . . . . .	10
3.0.2	Expectation of Linear Combinations . . . . .	10
3.0.3	Weighted Scalar Variable . . . . .	11
3.0.4	Gaussian Moments . . . . .	11
<b>4</b>	<b>Kronecker and Vec Operator</b>	<b>12</b>
4.1	Flattening . . . . .	12
4.2	The Kronecker Product . . . . .	13
4.3	The Vec Operator . . . . .	14
4.4	General Matrifcation . . . . .	15
4.4.1	The Lyapunov Equation . . . . .	15
4.4.2	Encapsulating Sum . . . . .	15
4.5	The Hadamard Product . . . . .	16
4.6	Khatri–Rao product . . . . .	16
<b>5</b>	<b>Determinant and Inverses</b>	<b>18</b>
5.1	Determinant . . . . .	18
5.2	Inverses . . . . .	19

<i>CONTENTS</i>	4
<b>6 Functions</b>	<b>20</b>
6.0.1 Pseudo-linear form . . . . .	20
6.1 Taylor Series . . . . .	20
<b>7 Advanced Derivatives</b>	<b>21</b>
7.1 Derivatives of vector norms . . . . .	21
7.1.1 Two-norm . . . . .	21
7.2 Derivatives of matrix norms . . . . .	21
7.3 Derivatives of Structured Matrices . . . . .	21
7.3.1 The Chain Rule . . . . .	21
7.3.2 The Hessian Chain Rule . . . . .	21
7.3.3 Symmetric . . . . .	22
7.3.4 Diagonal . . . . .	22
7.3.5 Toeplitz . . . . .	22
7.4 Derivatives of a Determinant . . . . .	22
7.5 General forms . . . . .	22
7.6 Linear forms . . . . .	22
7.7 Square forms . . . . .	22
7.8 Derivatives of an Inverse . . . . .	22
7.9 Derivatives of Eigenvalues . . . . .	22
<b>8 Special Matrices</b>	<b>23</b>
8.0.1 Block matrices . . . . .	23
8.0.2 The Discrete Fourier Transform Matrix . . . . .	23
8.0.3 Hermitian Matrices and skew-Hermitian . . . . .	23
8.0.4 Idempotent Matrices . . . . .	23
8.0.5 Orthogonal matrices . . . . .	23
8.0.6 Positive Definite and Semi-definite Matrices . . . . .	23
8.0.7 Singleentry Matrix, The . . . . .	23
8.0.8 Symmetric, Skew-symmetric/Antisymmetric . . . . .	24
8.0.9 Toeplitz Matrices . . . . .	24
8.0.10 Units, Permutation and Shift . . . . .	24
8.0.11 Vandermonde Matrices . . . . .	24
<b>9 Machine Learning Applications</b>	<b>25</b>
9.1 Least Squares . . . . .	25
9.2 Hessian of Cross Entropy Loss . . . . .	25
9.3 Convolutional Neural Networks . . . . .	25
9.4 Transformers / Attention . . . . .	25
9.5 Tensor Sketch . . . . .	25
<b>10 Tensorgrad</b>	<b>26</b>
10.1 Simplification Rules . . . . .	26
10.2 Functions . . . . .	26
10.3 Isomorphisms . . . . .	26
<b>11 Appendix</b>	<b>27</b>

## 1.1 Notation and Nomenclature

Dot product	$a-b$	$y = \sum_i a_i b_i$	$[\cdot \cdot \cdot \cdot] \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$	$= y$
Outer product	$-a \quad b-$	$Y_{i,j} = a_i b_j$	$\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} [\cdot \cdot \cdot \cdot]$	$= -Y-$
Matrix-Vector	$-A-b$	$y_i = \sum_j A_{i,j} b_j$	$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$	$= -y$
Matrix-Matrix	$-A-B-$	$Y_{i,k} = \sum_j A_{i,j} B_{j,k}$	$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$	$= -Y-$

The rank-4 copy tensor,  $C_{i,j,k,l} = [i = j = k = l]$ ,  $\succ\diagdown$ , differs from the outer product of two identity matrices (in the Cookbook denoted  $J$ ), which satisfies  $J_{i,j,k,l} = [i = k][j = l]$  and which we'd write as  $J = \begin{smallmatrix} \circ & \circ \\ \hline \circ & \circ \end{smallmatrix}$ , and satisfies, for example,  $\frac{dX}{dX} = J$ .

## 1.2 Basics

$-A- = -Q-\underset{\lambda}{\circ}-Q^{-1}-$  where  $\lambda_i$  is the  $i$ th eigenvalue of  $A$ .

More general principles:

1. You can contract edges in any order.
2. You can always contract connected subgraphs of Copy tensors.
3. distributive law of sums/products.
4. Sums and broadcasting.

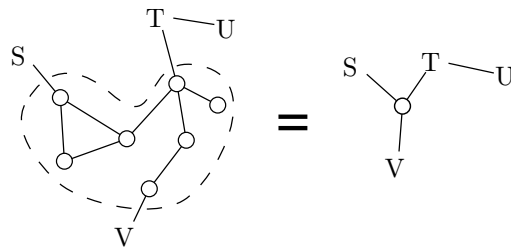


Figure 1.1: Contract spiders

### 1.2.1 Trace

The “trace” of a square matrix is defined  $\text{Tr}(A) = \sum_i A_{i,i}$ . In tensor diagram notation, that corresponds to a self-edge:  $\begin{smallmatrix} \circ \\ \hline \circ \end{smallmatrix} A$ .

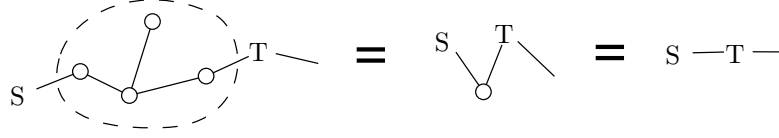


Figure 1.2: Eliminate Identity

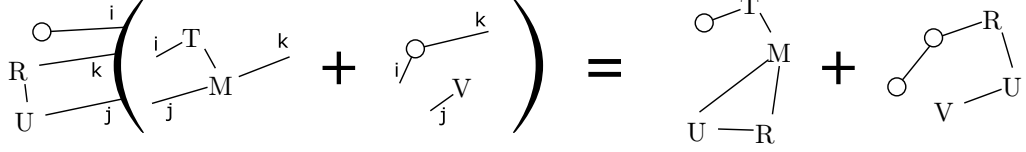


Figure 1.3: Distributive law

$$\sum_{i=1}^n A_{ii} = \text{Tr}(A) = \text{Tr}(AI) \quad \text{A} = \text{A} \text{---} \circ \quad (11)$$

$$\text{Tr}(A) = \sum_i \lambda_i = \langle 1, \lambda \rangle \quad \text{A} = \text{Q} \text{---} \circ \text{---} \text{Q}^{-1} \quad (12)$$

$$\text{Tr}(A) = \text{Tr}(A^T) \quad \text{A} = \text{A} \quad (13)$$

$$\text{Tr}(AB) = \text{Tr}(BA) \quad \text{A} \text{---} \text{B} = \text{B} \text{---} \text{A} \quad (14)$$

$$\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B) \quad \text{(A+B)---} \circ = \text{A---} \circ + \text{B---} \circ \quad (15)$$

$$\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB) \quad \text{A---B---C} = \text{B---C---A} = \text{C---A---B} \quad (16)$$

$$a^T a = \text{Tr}(aa^T) \quad a \text{---} a = \text{Tr}(-a \ a) = \text{a} \text{---} \text{a} \quad (17)$$

## Chapter 2

# Simple Derivatives

### 2.1 Derivatives of Matrices, Vectors and Scalar Forms

#### 2.1.1 First Order

The Matrix Cookbook defines the single-entry matrix  $J^{i,j} \in R^{n \times n}$  as the matrix which is zero everywhere except in the entry  $(i, j)$  in which it is 1. Alternatively we could write  $J_{n,m}^{i,j} = [i = n][j = m]$ .

$$\frac{\partial x^T a}{\partial x} = \frac{\partial a^T x}{\partial x} = a \quad (x - a) \begin{smallmatrix} \nearrow \\ \bullet \end{smallmatrix} = (x) \begin{smallmatrix} \nearrow \\ \bullet \end{smallmatrix} - a = \nearrow a \quad (69)$$

$$\frac{\partial a^T X b}{\partial X} = ab^T \quad (a - X - b) \begin{smallmatrix} \nearrow \\ \bullet \end{smallmatrix} = a - (X) \begin{smallmatrix} \nearrow \\ \bullet \end{smallmatrix} - b = a - \nearrow b \quad (70)$$

$$\frac{\partial X}{\partial X_{i,j}} = J^{i,j} \quad (-X - \begin{smallmatrix} \nearrow \\ \bullet \\ \nearrow \end{smallmatrix} \begin{smallmatrix} i \\ j \end{smallmatrix}) = -\begin{smallmatrix} \nearrow \\ \bullet \\ \nearrow \end{smallmatrix} \begin{smallmatrix} i \\ j \end{smallmatrix} \quad (73)$$

$$\begin{aligned} \frac{\partial (XA)_{i,j}}{\partial X_{m,n}} &= (J^{m,n} A)_{i,j} & (\begin{smallmatrix} i \\ - \end{smallmatrix} X - A \begin{smallmatrix} \nearrow \\ \bullet \\ \nearrow \end{smallmatrix} \begin{smallmatrix} m \\ n \end{smallmatrix}) &= - (X) \begin{smallmatrix} \nearrow \\ \bullet \end{smallmatrix} - A - \\ & & &= \begin{smallmatrix} i \\ \nearrow \end{smallmatrix} \begin{smallmatrix} m \\ \nearrow \end{smallmatrix} A \begin{smallmatrix} \nearrow \\ \bullet \end{smallmatrix} \end{aligned} \quad (74)$$



### 2.1.2 Second Order

$$\begin{aligned} \frac{\partial}{\partial X_{i,j}} \sum_{k,l,m,n} X_{k,l} X_{m,n} &= (\sum_{k,l} X_{k,l})^2 \\ &= 2 \sum_{k,l} X_{k,l} \end{aligned} \quad \begin{aligned} \left( \begin{array}{c} \circ - X - \circ \\ \circ - X - \circ \end{array} \right) \overset{i}{\nearrow}_j &= \begin{array}{c} \circ - X - \circ \\ \circ - (X) \end{array} + \begin{array}{c} \circ - (X) \\ \circ - X - \circ \end{array} \\ &= 2 \begin{array}{c} \circ - X - \circ \\ \circ - \overset{i}{\nearrow}_j \end{array} \end{aligned} \quad (76)$$

$$\begin{aligned} \frac{\partial b^T X^T X c}{\partial X} &= X(b c^T + c b^T) \\ (b - X^T - X - c) \overset{i}{\nearrow}_j &= b - X^T - (X) \overset{i}{\nearrow}_j - c \\ &+ b - (X^T) \overset{i}{\nearrow}_j - X - c \\ &= b - X^T - \overset{i}{\nearrow}_j c \\ &+ b - \overset{i}{\nearrow}_j X - c \\ &= -X - (-b c' + -c b') \end{aligned} \quad (77)$$

$$\begin{aligned} \frac{\partial}{\partial x} (Bx + b)^T C (Dx + d) &= B^T C (Dx + d) \\ &+ D^T C^T (Bx + b) \end{aligned} \quad (78)$$

$$\begin{aligned} \frac{\partial}{\partial X_{i,j}} (X^T B X)_{k,l} &= \delta_{l,j} (X^T B)_{k,i} \\ &+ \delta_{k,j} (B X)_{i,l} \end{aligned} \quad \begin{aligned} \dots \\ (\overset{k}{\leftarrow} X^T - B - X \overset{i}{\nearrow}_j) \overset{j}{\nwarrow}_i &= \overset{k}{\leftarrow} X^T - B \overset{j}{\nwarrow}_i \\ &+ \overset{k}{\leftarrow} \overset{j}{\nwarrow}_i B - X \overset{j}{\nwarrow}_i \end{aligned} \quad (79)$$

$$\frac{\partial}{\partial X_{i,j}} X^T B X = X^T B J^{i,j} + J^{j,i} B X \quad (\text{same as above}) \quad (80)$$

$$\begin{aligned} \frac{\partial}{\partial x} x^T B x &= (B + B^T) x \\ (x - B - x) \overset{i}{\nearrow}_j &= -B - x + x - B - \\ &= x \overset{i}{\nwarrow}_j \left( \overset{j}{\nwarrow}_i B \overset{i}{\nwarrow}_j \right) \end{aligned} \quad (81)$$

$$\frac{\partial}{\partial X} b^T X^T D X c = D^T X b c^T + D X c b^T \quad \dots \quad (82)$$

$$\frac{\partial}{\partial X} (Xb + c)^T D (Xb + c) = (D + D^T) (Xb + c) b^T \quad \dots \quad (83)$$

Assume  $W$  is symmetric, then... (84) - (88)

### 2.1.3 Higher-order and non-linear

...

### 2.1.4 Gradient and Hessian

...



## Chapter 3

# Statistics and Probability

### 3.0.1 Definition of Moments

Let  $x \in \mathbb{R}^n$  is a random variable. We write  $m = E[x] \in \mathbb{R}^n$  for the expectation and  $M = \text{Var}[x] = E[(x - m)(x - m)^T]$  for the covariance (when these quantities are defined.)

In tensor diagrams, we will use square brackets:

$$m = [-x] \quad \text{and} \quad M = [-(x \ominus m) \quad (x \div m) -]$$

Note we used the German minus,  $\div$ , to distinguish subtraction from contraction edges.

We can also define the third and fourth centralized moment tensors

$$M_3 = \begin{bmatrix} (x \div m) - \\ (x \div m) - \\ (x \div m) - \end{bmatrix} \quad \text{and} \quad M_4 = \begin{bmatrix} (x \div m) - \\ (x \div m) - \\ (x \div m) - \\ (x \div m) - \end{bmatrix}.$$

### 3.0.2 Expectation of Linear Combinations

General principle: The “linearity of expectation” lets you pull out all parts of the graph not involving  $X$ .

#### Linear Forms

$$E[AXB + C] = AE[X]B + C \quad \begin{bmatrix} -A-X-B- \\ + -C- \end{bmatrix} = \begin{bmatrix} -A-[X]-B- \\ + -C- \end{bmatrix} \quad (312)$$

$$\begin{aligned} \text{Var}[Ax] &= A\text{Var}[x]A^T \quad \begin{bmatrix} A-x \div [A-x] \\ A-x \div [A-x] \end{bmatrix} = \begin{bmatrix} A-(x \div m) \\ A-(x \div m) \end{bmatrix} \quad (313) \\ &= A \begin{bmatrix} (x \div m) \\ (x \div m) \end{bmatrix} \\ &= -A-M_2-A- \end{aligned}$$

**Quadratic Forms**

$$\begin{aligned}
E[x^T A x] &= \text{Tr}(A \Sigma) + \mu^T A \mu \\
[x - A - x] &= [(x \div \mu) - A - (x \div \mu) + \mu - A - \mu] \\
&= \left[ \begin{matrix} (x \div m) \\ (x \div m) \end{matrix} \right] A + \mu - A - \mu \\
&= \overbrace{\Sigma - A} + \mu - A - \mu
\end{aligned}$$

**Cubic Forms****3.0.3 Weighted Scalar Variable**

Let  $y = w^T x$ , and let  $m = E[y]$ , then

$$\begin{aligned}
E[y] &= m = w^T \mu \\
E[(y - m)^2] &= w - M_2 - w \\
E[(y - m)^3] &= w - \overset{w}{M_3} - w \\
E[(y - m)^4] &= w - \overset{w}{M_4} - w
\end{aligned}$$

For specific distributions, like  $x$  Gaussian, we can often reduce the moment tensors further. Khintchine's inequality also gives a way to bound all of these in terms of  $E[(y - m)^2]$ .

**3.0.4 Gaussian Moments**

**Mean and covariance of linear forms**

**Mean and variance of square forms**

**Cubic forms**

**Mean of Quartic Forms**

**Gaussian Integration by Parts**

General principle for Gaussian expectations.

## Chapter 4

# Kronecker and Vec Operator

### 4.1 Flattening

Flattening is a common operation for programmers. In the language of numpy, we may write `np.ones((2,3,4)).reshape(2, 12)` to flatten a shape (2,3,4) tensor into a shape (2,12) matrix. Similarly, in mathematical notation,  $\text{vec}(X)$  is commonly used to denote the flattening of a matrix into a vector.

Typically the main reason to do this is as a cludge for dealing with bad general notation for tensors. Hence, with tensor diagrams, we can avoid this operation entirely. However, it is still interesting to see how tensor diagrams can make a lot of properties of flattening much more transparent.

To begin with we note that flattening is a linear operation, and hence can be represented as a simple tensor. We'll use a triangle to denote this:

$$\triangleright_{i,j,k} = \begin{array}{c} i \\ \diagup \\ \diagdown \\ j \end{array} \triangleright^k = [i + jn = k].$$

Here  $n$  is the dimension of the  $i$  edge. Note we use a double line to denote the output of the flattening operation. This is simply a syntactic choice to remind ourselves that the output is a bundle of two edges.

Using this notation we can write

$$\text{vec}(X)_k = \sum_{i,j} \triangleright_{i,j,k} X_{i,j} = X \begin{array}{c} \text{---} \triangleright^k \\ \text{---} \end{array}.$$

The basic property of  $\triangleright$  is that opposing triangles cancel:

$$\begin{array}{c} \text{and} \end{array} \begin{array}{c} \text{---} \triangleright \text{---} \triangleleft \text{---} \\ \text{---} \triangleleft \text{---} \triangleright \text{---} \end{array} = \begin{array}{c} \text{---} \\ \text{---} \end{array}.$$

## 4.2 The Kronecker Product

The Kronecker product of an  $m \times n$  matrix  $A$  and an  $r \times q$  matrix  $B$ , is an  $mr \times nq$  matrix,  $A \otimes B$  defined as

$$A \otimes B = \begin{bmatrix} A_{1,1}B & A_{1,2}B & \cdots & A_{1,n}B \\ A_{2,1}B & A_{2,2}B & \cdots & A_{2,n}B \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1}B & A_{m,2}B & \cdots & A_{m,n}B \end{bmatrix}.$$

Using index notation we can also write this as  $(A \otimes B)_{p(r-1)+v, q(s-1)+w} = A_{rs}B_{vw}$ , but it's pretty hard to read.

In tensor notation the Kronecker Product is simply the outer product of two matrices, flattened “on both sides”:  $A \otimes B = \begin{array}{c} \text{A} \\ \text{B} \end{array}$ .

The Kronecker product has the following properties:

$$A \otimes (B + C) = A \otimes B + A \otimes C \quad \begin{array}{c} \text{A} \\ \text{(B+C)} \end{array} = \begin{array}{c} \text{A} \\ \text{B} \end{array} + \begin{array}{c} \text{A} \\ \text{C} \end{array} \quad (506)$$

$$A \otimes (B \otimes C) = (A \otimes B) \otimes C \quad \begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \end{array} = \begin{array}{c} \text{A} \\ \text{B} \end{array} \otimes \begin{array}{c} \text{C} \end{array} \quad (508)$$

$$aA \otimes bB = ab(A \otimes B) \quad \begin{array}{c} a \text{ A} \\ b \text{ B} \end{array} = \begin{array}{c} ab \\ \text{A} \\ \text{B} \end{array} \quad (509)$$

$$(A \otimes B)^T = A^T \otimes B^T \quad \dots \quad (510)$$

$$(A \otimes B)(C \otimes D) = AC \otimes BD \quad \begin{array}{c} \text{A} \\ \text{B} \end{array} \begin{array}{c} \text{C} \\ \text{D} \end{array} = \begin{array}{c} \text{A} \cdot \text{C} \\ \text{B} \cdot \text{D} \end{array} \quad (511)$$

$$(A \otimes I)(I \otimes B) = A \otimes B \quad \begin{array}{c} \text{A} \\ \text{B} \end{array} = \begin{array}{c} \text{A} \\ \text{B} \end{array} \quad (511b)$$

$$\text{Tr}(A \otimes B) = \text{Tr}(A)\text{Tr}(B) \quad \begin{array}{c} \text{A} \\ \text{B} \end{array} = \begin{array}{c} \text{A} \\ \text{B} \end{array} = \text{A} \text{ B} \quad (515)$$

$$\begin{aligned}
 \text{eig}(A \otimes B) = \text{eig}(A)\text{eig}(B) \quad & \begin{array}{c} \lambda_1 \\ \circ \\ Q_1 - \circ - Q_1^{-1} \\ \circ \\ Q_2 - \circ - Q_2^{-1} \\ \lambda_2 \end{array} \triangleright \triangleleft = \begin{array}{c} \lambda_1 \\ \circ \\ Q_1 \triangleright \times \triangleleft \circ \\ \circ \\ Q_2 \triangleright \times \triangleleft \circ \\ \lambda_2 \end{array} \begin{array}{c} Q_1^{-1} \\ \circ \\ Q_2^{-1} \end{array} \triangleright \triangleleft \\
 & (519) \\
 & = \begin{array}{c} Q_1 \\ \circ \\ Q_2 \end{array} \triangleright \triangleleft \circ \begin{array}{c} Q_1^{-1} \\ \circ \\ Q_2^{-1} \end{array} \triangleright \triangleleft \\
 & \quad \quad \quad \begin{array}{c} \lambda_1 \\ \circ \\ \lambda_2 \end{array}
 \end{aligned}$$

Here the last equation shows an interesting general property of  $\triangleright$ :

$$\begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \\ M \end{array} \triangleright \triangleleft = \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \\ M \end{array} = \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \\ M \end{array}.$$

While the right side is just  $\text{diag}(\text{vec}(M))$ , the left side is harder to write in classical notation. We can write it with index notation as

$$v_{k,\ell} = \sum_{i,j} \triangleright_{i,j,k} \triangleright_{i,j,\ell} M_{i,j} = \sum_{i,j} [i + jn = k][i + jn = \ell] M_{i,j},$$

but it is hardly enlightening.

### 4.3 The Vec Operator

The vec-operator applied on a matrix  $A$  stacks the columns into a vector, i.e. for a  $2 \times 2$  matrix

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \text{vec}(A) = \begin{bmatrix} A_{11} \\ A_{21} \\ A_{12} \\ A_{22} \end{bmatrix}$$

At the start of the chapter we showed how to represent the vec-operator using the flattening tensor:  $\text{vec}(X) = X \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \end{array}$ . The Matrix Cookbook gives the following properties of the vec-operator:

$$\text{vec}(A^T X B) = \text{vec}(X)^T (A \otimes B) \quad \begin{array}{c} A \\ \diagup \quad \diagdown \\ X \text{---} \text{---} \\ \diagdown \quad \diagup \\ B \end{array} \triangleright \triangleleft = X \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \end{array} \begin{array}{c} A \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \\ B \end{array} \triangleright \triangleleft \quad (520)$$

$$\text{Tr}(A^T B) = \text{vec}(A)^T \text{vec}(B) \quad A \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \end{array} B = A \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \end{array} \begin{array}{c} A \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \\ B \end{array} \triangleright \triangleleft \quad (521)$$

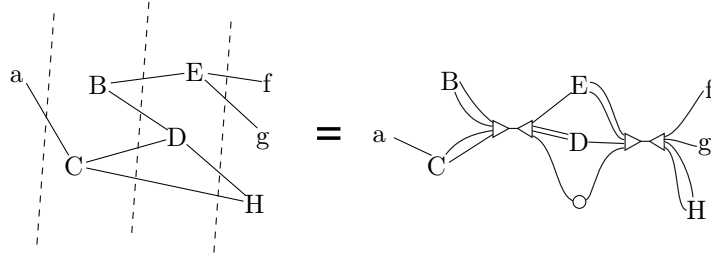
$$\text{vec}(A + B) = \text{vec}(A) + \text{vec}(B) \quad (A+B) \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \end{array} = A \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \end{array} + B \begin{array}{c} \circ \\ \diagup \quad \diagdown \\ \text{---} \quad \text{---} \\ \circ \end{array} \triangleright \triangleleft \quad (522)$$

$$\text{vec}(aA) = a \text{vec}(A) \quad aA \rhd = a A \rhd \quad (523)$$

$$\begin{aligned} a^T X B X^T c = \text{vec}(X)^T (B \otimes ca^T) \text{vec}(X) \quad a-X-B-X-c &= X \begin{array}{c} \text{B} \\ \swarrow \quad \searrow \\ a \quad b \end{array} X \quad (524) \\ &= X \rhd \rhd \rhd \begin{array}{c} \text{B} \\ \swarrow \quad \searrow \\ a \quad b \end{array} \rhd \rhd X \end{aligned}$$

## 4.4 General Matrifaction

The last equation is an example of a general idea: Any tensor network can be transformed into a series of matrix multiplications by applying the  $\text{vec}$ -operator to all tensors and the flattening tensor to all edges. For example, the following complicated graph:



Can be written as a simple vector-matrix-matrix-vector product,  $aM_1M_2b$ , where  $M_1 = \text{vec}(B) \otimes C'$ ,  $M_2 = E' \otimes D' \otimes I$  and  $b = f \otimes g \otimes \text{vec}(H)$ , where  $C'$ ,  $D'$  and  $E'$  are rank 3 tensors flattened on one side, and  $\text{vec}(B)$  is interpreted as a matrix with a single column.

### 4.4.1 The Lyapunov Equation

A nice application of Kronecker product rewritings is to solve equations like

$$AX + XB = C. \quad (272)$$

We use the rewriting  $\text{vec}(AX + XB) = (I \otimes A + B^T \otimes I)\text{vec}(X)$ , which follows from the tensor diagram massaging:

$$\left( \begin{array}{c} - A - X - \\ + - X - B - \end{array} \right) \rhd = X \begin{array}{c} \text{A} \\ \swarrow \quad \searrow \end{array} \rhd + X \begin{array}{c} \text{B} \\ \swarrow \quad \searrow \end{array} \rhd = X \rhd \left( \begin{array}{c} \text{A} \\ \swarrow \quad \searrow \\ \text{B} \end{array} \right) =$$

after which we can take the normal matrix inverse to get

$$\text{vec}(X) = (I \otimes A + B^T \otimes I)^{-1} \text{vec}(C). \quad (273)$$

### 4.4.2 Encapsulating Sum

This is a generalization of the previous equation.

$$\sum_n A_n X B_n = C \quad (274)$$

$$\text{vec}(X) = \left( \sum_n B_n^T \otimes A_n \right)^{-1} \text{vec}(C) \quad (275)$$



## 4.5 The Hadamard Product

The Hadamard product, also known as element-wise multiplication, is not described in the Matrix Cookbook. Yet, it is a very useful operation, and has some interesting properties in connection with the Kronecker product.

We define the Hadamard product of two  $2 \times 2$  matrices  $A$  and  $B$  as

$$A \circ B = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \circ \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} A_{11}B_{11} & A_{12}B_{12} \\ A_{21}B_{21} & A_{22}B_{22} \end{bmatrix}.$$

In tensor notation, the Hadamard product can be represented using two rank-3 copy tensors:  $- \circ \begin{smallmatrix} A \\ B \end{smallmatrix} \circ -$ . Some properties of the Hadamard product are:

$$x^T (A \circ B) y = \text{tr}(A^T D_x B D_y) \quad x - \circ \begin{smallmatrix} A \\ B \end{smallmatrix} \circ - y = \overbrace{A^T - \circ - B - \circ}^{\substack{x \\ y}}$$

$$(A \otimes B) \circ (C \otimes D) = (A \circ C) \otimes (B \circ D) \quad = \circ \begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \end{array} \circ = \begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \end{array}$$

To see why the last equation holds, it suffices to follow the double lines to see that  $A$  and  $C$  both use the “upper” part of the double edge, while  $B$  and  $D$  use the lower part.

## 4.6 Khatri–Rao product

Also known as the column-wise Kronecker, row-wise Kronecker or “Face-splitting Product”. We use the symbols  $*$  and  $\bullet$  for the column and row-wise Kronecker products, respectively.

$$A * B = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} * \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} A_{11}B_{11} & A_{12}B_{12} \\ A_{11}B_{21} & A_{22}B_{22} \\ A_{21}B_{11} & A_{22}B_{12} \\ A_{21}B_{21} & A_{22}B_{22} \end{bmatrix}$$

$$A \bullet B = \dots = \begin{bmatrix} A_{11}B_{11} & A_{11}B_{12} & A_{12}B_{11} & A_{12}B_{12} \\ A_{21}B_{21} & A_{21}B_{22} & A_{22}B_{21} & A_{22}B_{22} \end{bmatrix}$$

In terms of tensor diagrams, these products correspond simply to flattening the product on one side, and using a copy tensor on the other:

$$A * B = \begin{array}{c} \text{A} \\ \text{B} \end{array} \circ -$$

$$A \bullet B = - \circ \begin{array}{c} \text{A} \\ \text{B} \end{array}$$

Clearly the two are identical up to transpose. Indeed,  $(A*B)^T = B^T \bullet A^T$  and  $(A \bullet B)^T = B^T * A^T$ .

There are multiple “mixed product” identities:

$$\begin{aligned}
 (A \bullet B)(C \otimes D) &= (AC) \bullet (BD) & - \circ \begin{array}{c} A \\ \diagdown \end{array} \begin{array}{c} \diagup \\ B \end{array} \begin{array}{c} \diagdown \\ C \end{array} \begin{array}{c} \diagup \\ D \end{array} = - \circ \begin{array}{c} A - C \\ \diagdown \\ B - D \end{array} \begin{array}{c} \diagup \end{array} \\
 (Ax) \circ (By) &= (A \bullet C)(x \otimes y) & - \circ \begin{array}{c} A - x \\ \diagdown \\ B - y \end{array} = - \circ \begin{array}{c} A \\ \diagdown \end{array} \begin{array}{c} \diagup \\ B \end{array} \begin{array}{c} \diagdown \\ x \end{array} \begin{array}{c} \diagup \\ y \end{array}
 \end{aligned}$$

## Chapter 5

# Determinant and Inverses

### 5.1 Determinant

It's convenient to write the determinant in tensor notation as

$$\det(A) = \frac{1}{n!} \overline{\overline{A \cdots A}}$$

where  $\overline{\overline{i_1 \ i_2 \ \cdots \ i_n}} = \varepsilon_{i_1, \dots, i_n}$  is the rank- $n$  Levi-Civita tensor defined by

$$\varepsilon_{i_1, \dots, i_n} = \begin{cases} \text{sign}(\sigma) & \sigma = (i_1, \dots, i_n) \text{ is a permutation} \\ 0 & \text{otherwise.} \end{cases}$$

To see that the definition makes sense, let's first consider

$$\det(I) = \frac{1}{n!} \overline{\overline{\begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}}} = \frac{1}{n!} \sum_{i_1, \dots, i_n, j_1, \dots, j_n} \varepsilon_{i_1, \dots, i_n} \varepsilon_{j_1, \dots, j_n} [i = j] = \frac{1}{n!} \sum_{i_1, \dots, i_n} \varepsilon_{i_1, \dots, i_n}^2 = 1.$$

In general we get from the permutation definition of the determinant:

$$\begin{aligned} \overline{\overline{A \cdots A}} &= \sum_{i_1, \dots, i_n, j_1, \dots, j_n} \varepsilon_{i_1, \dots, i_n} \varepsilon_{j_1, \dots, j_n} A_{i_1, j_1} \cdots A_{i_n, j_n} \\ &= \sum_{\sigma, \tau} \text{sign}(\sigma) \text{sign}(\tau) A_{\sigma_1, \tau_1} \cdots A_{\sigma_n, \tau_n} \\ &= \sum_{\sigma} \text{sign}(\sigma) \sum_{\tau} \text{sign}(\tau) A_{\sigma_1, \tau_1} \cdots A_{\sigma_n, \tau_n} \\ &= \sum_{\sigma} \text{sign}(\sigma)^2 \det(A) \\ &= n! \det(A). \end{aligned}$$

The definition generalizes to Cayley's “hyper determinants” by . . . .

A curious property is that

$$\begin{array}{|c|} \hline A \\ \hline \end{array} \cdots \begin{array}{|c|} \hline A \\ \hline \end{array} = \begin{array}{|c|} \hline \cdots \\ \hline \end{array} \begin{array}{|c|} \hline A \\ \hline \end{array} \cdots \begin{array}{|c|} \hline A \\ \hline \end{array}$$

$$(18) \quad \det(A) = \prod_i \lambda_i \quad \dots$$

$$(19) \quad \det(cA) = c^n \det(A) \quad \begin{array}{|c|} \hline cA \\ \hline \end{array} \cdots \begin{array}{|c|} \hline cA \\ \hline \end{array} = c^n \begin{array}{|c|} \hline A \\ \hline \end{array} \cdots \begin{array}{|c|} \hline A \\ \hline \end{array}$$

$$(20) \quad \det(A) = \det(A^T) \quad \dots$$

$$(21) \quad \det(AB) = \det(A)\det(B) \quad \begin{array}{|c|} \hline A \\ \hline \end{array} \cdots \begin{array}{|c|} \hline A \\ \hline \end{array} \begin{array}{|c|} \hline B \\ \hline \end{array} \cdots \begin{array}{|c|} \hline B \\ \hline \end{array} = \begin{array}{|c|} \hline A \\ \hline \end{array} \cdots \begin{array}{|c|} \hline A \\ \hline \end{array} \begin{array}{|c|} \hline B \\ \hline \end{array} \cdots \begin{array}{|c|} \hline B \\ \hline \end{array}$$

$$(22) \quad \det(A^{-1}) = 1/\det(A) \quad \dots$$

$$(23) \quad \det(A^n) = \det(A)^n \quad \dots$$

$$(24) \quad \det(I + uv^T) = 1 + u^T v \quad \dots$$

## 5.2 Inverses

Might be reduced, unless cofactor matrices have a nice representation?

## Chapter 6

# Functions

Stuff about analytical matrix functions. Such as Exponential Matrix Function. I'd rather talk about my general function notation. And maybe about taylor series?

### 6.0.1 Pseudo-linear form

Maybe this should just be an example in a table?

Derivation of Peyman Milanfar's gradient

$$\begin{aligned}d[\mathbf{f}(\mathbf{x})] &= d[\mathbf{A}(\mathbf{x})\mathbf{x}] \\&= d[\mathbf{A}(\mathbf{x})]\mathbf{x} + \mathbf{A}(\mathbf{x})d\mathbf{x} \\&= \text{vec}\{d[\mathbf{A}(\mathbf{x})]\mathbf{x}\} + \mathbf{A}(\mathbf{x})d\mathbf{x} \\&= \text{vec}\{\mathbf{Id}[\mathbf{A}(\mathbf{x})]\mathbf{x}\} + \mathbf{A}(\mathbf{x})d\mathbf{x} \\&= (\mathbf{x}^T \otimes \mathbf{I}) \text{vec}\{d[\mathbf{A}(\mathbf{x})]\} + \mathbf{A}(\mathbf{x})d\mathbf{x} \\&= (\mathbf{x}^T \otimes \mathbf{I}) D \text{vec}[\mathbf{A}(\mathbf{x})]d\mathbf{x} + \mathbf{A}(\mathbf{x})d\mathbf{x} \\&= [(\mathbf{x}^T \otimes \mathbf{I}) D \text{vec}[\mathbf{A}(\mathbf{x})] + \mathbf{A}(\mathbf{x})] d\mathbf{x}\end{aligned}$$

### 6.1 Taylor Series

I've always wanted to write these out properly with tensors.

## Chapter 7

# Advanced Derivatives

### 7.1 Derivatives of vector norms

#### 7.1.1 Two-norm

$$\frac{d}{dx} \|x - a\|_2 = \frac{x - a}{\|x - a\|_2} \quad (7.1)$$

$$\frac{d}{dx} \frac{x - a}{\|x - a\|_2} = \frac{I}{\|x - a\|_2} - \frac{(x - a)(x - a)^T}{\|x - a\|_2^3} \quad (7.2)$$

$$\frac{d}{dx} \|x\|_2^2 = \frac{d}{dx} \|x^T x\|_2 = 2x \quad (7.3)$$

### 7.2 Derivatives of matrix norms

### 7.3 Derivatives of Structured Matrices

#### 7.3.1 The Chain Rule

Sometimes the objective is to find the derivative of a matrix which is a function of another matrix.

#### 7.3.2 The Hessian Chain Rule

See main.tex

7.3.3 Symmetric

7.3.4 Diagonal

7.3.5 Toeplitz

7.4 Derivatives of a Determinant

7.5 General forms

7.6 Linear forms

7.7 Square forms

7.8 Derivatives of an Inverse

7.9 Derivatives of Eigenvalues

## Chapter 8

# Special Matrices

### 8.0.1 Block matrices

Stuff like Schur complements is interesting. But can we say anything useful using tensor diagrams?

### 8.0.2 The Discrete Fourier Transform Matrix

I think FFT can be nicely described with diagrams

Let's start with the Hadamard matrix:  $H_n = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}^{\otimes n}$ . Hm. It's just a bunch of matrices below each other, kinda boring.

What about the FFT? Does that require a bit more?

### 8.0.3 Hermitian Matrices and skew-Hermitian

Complex. Skip

### 8.0.4 Idempotent Matrices

Skip

### 8.0.5 Orthogonal matrices

Skip

### 8.0.6 Positive Definite and Semi-definite Matrices

Skip

### 8.0.7 Singleentry Matrix, The

Describes the matrix  $J$ . All of this is trivial with diagrams.



**8.0.8 Symmetric, Skew-symmetric/Antisymmetric**

Could introduce Penrose's symmetric tensors here?

**8.0.9 Toeplitz Matrices**

Could talk about the convolution tensor here...

**8.0.10 Units, Permutation and Shift**

Not that interesting...

**8.0.11 Vandermonde Matrices**

Does this have a nice description? Not a lot of properties are given in the Cookbook.

## Chapter 9

# Machine Learning Applications

9.1 Least Squares

9.2 Hessian of Cross Entropy Loss

9.3 Convolutional Neural Networks

9.4 Transformers / Attention

9.5 Tensor Sketch

## Chapter 10

# Tensorgrad

Implementation details

### 10.1 Simplification Rules

### 10.2 Functions

### 10.3 Isomorphisms

## Chapter 11

# Appendix

Contains some proofs, such as of equation 524 or 571. They are pretty long and could be useful for contrasting with the diagram proofs.