# The Tensor Cookbook

Thomas Dybdahl Ahle

May 28, 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 amendable 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

## 1.1 Notation and Nomenclature

Dot product $\qquad a\!-\!b \qquad y = \sum_i a_i b_i \qquad [\,\cdot\ \cdot\ \cdot\ \cdot\,]\begin{bmatrix}\cdot\\\cdot\\\cdot\\\cdot\end{bmatrix} \qquad = y$

Outer product $\quad -a\ \ b- \qquad Y_{i,j} = a_i b_j \qquad \begin{bmatrix}\cdot\\\cdot\\\cdot\\\cdot\end{bmatrix}[\,\cdot\ \cdot\ \cdot\ \cdot\,] \qquad = -Y-$

Matrix-Vector $\quad -A\!-\!b \qquad y_i = \sum_j A_{i,j} b_j \qquad \begin{bmatrix}\cdot\ \cdot\ \cdot\ \cdot\\\cdot\ \cdot\ \cdot\ \cdot\\\cdot\ \cdot\ \cdot\ \cdot\end{bmatrix}\begin{bmatrix}\cdot\\\cdot\\\cdot\\\cdot\end{bmatrix} \qquad = -y$

Matrix-Matrix $\quad -A\!-\!B- \quad Y_{i,k} = \sum_j A_{i,j} B_{j,k} \quad \begin{bmatrix}\cdot\ \cdot\ \cdot\ \cdot\\\cdot\ \cdot\ \cdot\ \cdot\\\cdot\ \cdot\ \cdot\ \cdot\end{bmatrix}\begin{bmatrix}\cdot\ \cdot\ \cdot\ \cdot\\\cdot\ \cdot\ \cdot\ \cdot\\\cdot\ \cdot\ \cdot\ \cdot\end{bmatrix} \quad = -Y-$

### 1.1.1 The Copy Tensor

We define $\circ\!-$ to be the all-ones vector. That is $\circ_i = 1$. We generalize $\circ$ to rank-n tensors by $\circ_{i,j,k,\ldots} = [i = j = k = \ldots]$. That is, the tensor with 1 on the diagonal, and 0 everywhere else. This is also known as the "copy" or "spider" tensor, or "generalized Kronecker delta". For rank-2 tensors, $-\circ- = I$, the identity matrix.

The rank-4 copy tensor, $C_{i,j,k,l} = [i = j = k = l]$, $\times\!\!\!\times$, 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 $J = \begin{smallmatrix}-\circ-\\-\circ-\end{smallmatrix}$, 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.
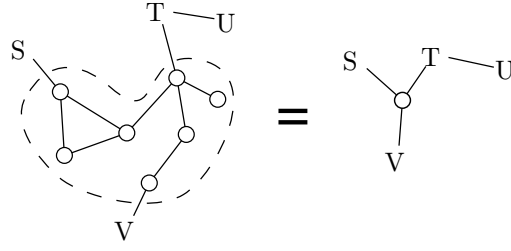
4. Sums and broadcasting.
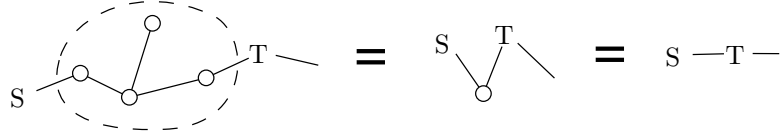
Figure 1.1: Contract spiders



Figure 1.2: Eliminate Identity

### 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: $\overset{\frown}{A}$ .

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

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

$$\text{Tr}(A) = \text{Tr}(A^T) \qquad \overset{\frown}{A} = \overset{\frown}{A} \qquad (13)$$

$$\text{Tr}(AB) = \text{Tr}(BA) \qquad \overset{\frown}{A-B} = \overset{\frown}{B-A} \qquad (14)$$

$$\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B) \qquad \overset{\frown}{(A+B)-\circ} = \overset{\frown}{A-\circ} \qquad (15)$$
$$+ \overset{\frown}{B-\circ}$$

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

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

Figure 1.3: Distributive law

# 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 $J^{i,j}_{n,m} = [i = n][j = m]$.

$$\frac{\partial x^T a}{\partial x} = \frac{\partial a^T x}{\partial x} = a \qquad\qquad (x\!-\!a) = (x)\!-\!a = \frown\!a \qquad (69)$$

$$\frac{\partial a^T X b}{\partial X} = ab^T \qquad\qquad (a\!-\!X\!-\!b) = a\!-\!(X)\!-\!b = a\!-\!\!\!/\!-\!b \qquad (70)$$

$$\frac{\partial X}{\partial X_{i,j}} = J^{i,j} \qquad\qquad (\!-\!X\!-\!)_j^i = \frac{_i}{_j} \qquad (73)$$

$$\frac{\partial (XA)_{i,j}}{\partial X_{m,n}} = (J^{m,n}A)_{i,j} \qquad (\underset{j}{\overset{i}{-}}\!X\!-\!A\!\underset{j}{\overset{m}{-}})^n = -(X)\!-\!A\!- \qquad (74)$$

$$= \frac{_i{}^m}{_n}\!A\!\underset{j}{-}$$

8

### 2.1.2 Second Order

$$\frac{\partial}{\partial X_{i,j}} \sum_{k,l,m,n} X_{k,l} X_{m,n} = (\sum_{k,l} X_{k,l})^2 \tag{76}$$

$$= 2 \sum_{k,l} X_{k,l}$$

$$\frac{\partial b^T X^T X c}{\partial X} = X(bc^T + cb^T) \tag{77}$$

$$\frac{\partial}{\partial x}(Bx + b)^T C(Dx + d) = B^T C(Dx + d) \tag{78}$$
$$+ D^T C^T (Bx + b)$$

$$\frac{\partial}{\partial X_{i,j}}(X^T B X)_{k,l} = \delta_{l,j}(X^T B)_{k,i} \tag{79}$$
$$+ \delta_{k,j}(BX)_{i,l}$$

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

$$\frac{\partial}{\partial x} x^T B x = (B + B^T)x \tag{81}$$

$$\frac{\partial}{\partial X} b^T X^T D X c = D^T X b c^T + D X c b^T \qquad \dots \tag{82}$$

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

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

### 2.1.3 Higher-order and non-linear

$$\frac{\partial (\mathbf{X}^n)_{kl}}{\partial X_{ij}} = \sum_{r=0}^{n-1} \left(\mathbf{X}^r \mathbf{J}^{ij} \mathbf{X}^{n-1-r}\right)_{kl}$$

For proof of the above, see B.1.3.

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{a}^T \mathbf{X}^n \mathbf{b} = \sum_{r=0}^{n-1} \left(\mathbf{X}^r\right)^T \mathbf{a}\mathbf{b}^T \left(\mathbf{X}^{n-1-r}\right)^T$$

## 2.2  Derivatives of Traces

Assume F(X) to be a differentiable function of each of the elements of X. It then holds that

$$\frac{d\mathrm{Tr}(F(x))}{dX} = f(X)^T,$$

where $f(\cdot)$ is the scalar derivative of $F(\cdot)$.

TODO: To show this with tensor diagrams, we first need to introduce our notation for functions.

### 2.2.1  First Order

$$\frac{\partial}{\partial X}\mathrm{Tr}(X) = I \qquad\qquad (99)$$

$$\frac{\partial}{\partial X}\mathrm{Tr}(XA) = A^T \qquad\qquad (100)$$

$$\frac{\partial}{\partial X}\mathrm{Tr}(AXB) = A^T B^T \qquad\qquad (101)$$

Continues for (102-105). The last one uses the Kronecker product, which we may have to introduce first.

### 2.2.2  Second Order

$$\frac{\partial}{\partial X}\mathrm{Tr}(X^2) = 2X^T$$

More:

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}^2 \mathbf{B}\right) = (\mathbf{X}\mathbf{B} + \mathbf{B}\mathbf{X})^T$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}^T \mathbf{B} \mathbf{X}\right) = \mathbf{B}\mathbf{X} + \mathbf{B}^T \mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{B}\mathbf{X}\mathbf{X}^T\right) = \mathbf{B}\mathbf{X} + \mathbf{B}^T \mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}\mathbf{X}^T \mathbf{B}\right) = \mathbf{B}\mathbf{X} + \mathbf{B}^T \mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}\mathbf{B}\mathbf{X}^T\right) = \mathbf{X}\mathbf{B}^T + \mathbf{X}\mathbf{B}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{B}\mathbf{X}^T \mathbf{X}\right) = \mathbf{X}\mathbf{B}^T + \mathbf{X}\mathbf{B}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}^T \mathbf{X}\mathbf{B}\right) = \mathbf{X}\mathbf{B}^T + \mathbf{X}\mathbf{B}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}(\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}) = \mathbf{A}^T \mathbf{X}^T \mathbf{B}^T + \mathbf{B}^T \mathbf{X}^T \mathbf{A}^T$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}^T \mathbf{X}\right) = \frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}\mathbf{X}^T\right) = 2\mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{B}^T \mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{B}\right) = \mathbf{C}^T \mathbf{X}\mathbf{B}\mathbf{B}^T + \mathbf{C}\mathbf{X}\mathbf{B}\mathbf{B}^T$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left[\mathbf{X}^T \mathbf{B}\mathbf{X}\mathbf{C}\right] = \mathbf{B}\mathbf{X}\mathbf{C} + \mathbf{B}^T \mathbf{X}\mathbf{C}^T$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{A}\mathbf{X}\mathbf{B}\mathbf{X}^T \mathbf{C}\right) = \mathbf{A}^T \mathbf{C}^T \mathbf{X}\mathbf{B}^T + \mathbf{C}\mathbf{A}\mathbf{X}\mathbf{B}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left[(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})^T\right] = 2\mathbf{A}^T (\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C})\mathbf{B}^T$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}(\mathbf{X} \otimes \mathbf{X}) = \frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}(\mathbf{X}) \operatorname{Tr}(\mathbf{X}) = 2\operatorname{Tr}(\mathbf{X})\mathbf{I}$$

### 2.2.3 Higher Order

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{X}^k\right) = k\left(\mathbf{X}^{k-1}\right)^T$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left(\mathbf{A}\mathbf{X}^k\right) = \sum_{r=0}^{k-1} \left(\mathbf{X}^r \mathbf{A}\mathbf{X}^{k-r-1}\right)^T$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr}\left[\mathbf{B}^T \mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{B}\right] = \mathbf{C}\mathbf{X}\mathbf{X}^T \mathbf{C}\mathbf{X}\mathbf{B}\mathbf{B}^T$$
$$+ \mathbf{C}^T \mathbf{X}\mathbf{B}\mathbf{B}^T \mathbf{X}^T \mathbf{C}^T \mathbf{X}$$
$$+ \mathbf{C}\mathbf{X}\mathbf{B}\mathbf{B}^T \mathbf{X}^T \mathbf{C}\mathbf{X}$$
$$+ \mathbf{C}^T \mathbf{X}\mathbf{X}^T \mathbf{C}^T \mathbf{X}\mathbf{B}\mathbf{B}$$

### 2.2.4 Other

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr} \left( \mathbf{A} \mathbf{X}^{-1} \mathbf{B} \right) = - \left( \mathbf{X}^{-1} \mathbf{B} \mathbf{A} \mathbf{X}^{-1} \right)^T = -\mathbf{X}^{-T} \mathbf{A}^T \mathbf{B}^T \mathbf{X}^{-T}$$

Assume $\mathbf{B}$ and $\mathbf{C}$ to be symmetric, then

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr} \left[ \left( \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \mathbf{A} \right] = - \left( \mathbf{C} \mathbf{X} \left( \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \right) \left( \mathbf{A} + \mathbf{A}^T \right) \left( \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1}$$

$$\frac{\partial}{\partial \mathbf{X}} \operatorname{Tr} \left[ \left( \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \left( \mathbf{X}^T \mathbf{B} \mathbf{X} \right) \right] = -2 \mathbf{C} \mathbf{X} \left( \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{B} \mathbf{X} \left( \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1}$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{X}} \operatorname{Tr} \left[ \left( \mathbf{A} + \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \left( \mathbf{X}^T \mathbf{B} \mathbf{X} \right) \right] = {}&-2 \mathbf{B} \mathbf{X} \left( \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \\ &- 2 \mathbf{C} \mathbf{X} \left( \mathbf{A} + \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{B} \mathbf{X} \left( \mathbf{A} + \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \\ &+ 2 \mathbf{B} \mathbf{X} \left( \mathbf{A} + \mathbf{X}^T \mathbf{C} \mathbf{X} \right)^{-1} \end{aligned}$$

See [7].

$$\frac{\partial \operatorname{Tr}(\sin(\mathbf{X}))}{\partial \mathbf{X}} = \cos(\mathbf{X})^T$$

# 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:

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

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**

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

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

$$= A\!-\!\!\begin{bmatrix} (x \div m) \\ (x \div m) \end{bmatrix}$$

$$= -A-M_2-A-$$

**Quadratic Forms**

$$\mathrm{E}[x^T A x] = \mathrm{Tr}(A\Sigma) + \mu^T A \mu$$
$$[x - A - x] = [(x \div \mu) - A - (x \div \mu) + \mu - A - \mu]$$
$$= \begin{bmatrix} (x \div m) \\ (x \div m) \end{bmatrix} A \\ + \mu - A - \mu$$
$$= \overset{\frown}{\Sigma - A} + \mu - A - \mu$$

**Cubic Forms**

### 3.0.3 Weighted Scalar Variable

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

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

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, $\mathrm{vec}(X)$ 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:

$$\rhd_{i,j,k} = \ {}_{j}\!\!\underset{}{\overset{i}{\diagdown}}\!\!\rhd\!=\!\!=^{k} = [i + jn = k].$$

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

$$\mathrm{vec}(X)_k = \sum_{i,j} \rhd_{i,j,k} X_{i,j} = \ \mathrm{X}\ \bigcirc\!\!\!\rhd\!\!=^{k}.$$

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

$$\rhd\!\!=\!\!\lhd\ \bigcirc \quad = \quad \bigcirc$$

$$\text{and} \quad =\!\!\lhd\bigcirc\rhd\!\!= \quad = \quad =\!\!= \quad .$$

## 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 = $  .

The Kronecker product has the following properties:

$$A \otimes (B + C) = A \otimes B + A \otimes C \qquad\qquad (506)$$

$$A \otimes (B \otimes C) = (A \otimes B) \otimes C \qquad\qquad (508)$$

$$aA \otimes bB = ab(A \otimes B) \qquad\qquad (509)$$

$$(A \otimes B)^T = A^T \otimes B^T \qquad\qquad (510)$$

$$(A \otimes B)(C \otimes D) = AC \otimes BD \qquad\qquad (511)$$

$$(A \otimes I)(I \otimes B) = A \otimes B \qquad\qquad (511b)$$

$$\operatorname{Tr}(A \otimes B) = \operatorname{Tr}(A)\operatorname{Tr}(B) \qquad\qquad (515)$$

$$\text{eig}(A \otimes B) = \text{eig}(A)\text{eig}(B)$$

(519)

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

(4.1)

This is easier to see when we consider that $V = $ represents the tensor where $V_{i,j,i,j} = M_{i,j}$ and 0 otherwise. So flattening $V$ on both sides is the same as $\text{diag}(\text{vec}(M))$, which is the rhs of (4.1).

## 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 . The Matrix Cookbook gives the following properties of the vec-operator:

$$\text{vec}(A^T X B) = \text{vec}(X)^T (A \otimes B)$$

(520)

$$\text{Tr}(A^T B) = \text{vec}(A)^T \text{vec}(B)$$

(521)

$$\text{vec}(A + B) = \text{vec}(A) + \text{vec}(B)$$

(522)

$$\text{vec}(aA) = a\,\text{vec}(A) \qquad\qquad aA \rangle\!\!=\; =\; a\; A \rangle\!\!= \qquad (523)$$

$$a^T X B X^T c = \text{vec}(X)^T (B \otimes ca^T)\text{vec}(X) \quad a\text{-}X\text{-}B\text{-}X\text{-}c = \; X \overset{B}{\underset{a\quad b}{\longrightarrow}} X \qquad (524)$$

$$= \; X \rangle\!\!\prec\!\!\!\ll \overset{B}{\underset{a\quad b}{\longrightarrow}} \rangle\!\!\prec\!\!\!\ll X$$

## 4.4   General Matrification

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 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, $aM_1 M_2 b$, 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

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

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:



after which we can take the normal matrix inverse to get

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

### 4.4.2   Encapsulating Sum

This is a generalization of the previous equation.

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

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

## 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: $- \circ \!\! \stackrel{A}{\underset{B}{\diamondsuit}} \!\! \circ -$. Some properties of the Hadamard product are:

$$x^T(A \circ B)y = \operatorname{tr}(A^T D_x B D_y) \qquad \mathrm{x} - \circ \!\! \stackrel{A}{\underset{B}{\diamondsuit}} \!\! \circ - \mathrm{y} \ = \ \overparen{A^T - \circ - B - \circ}_{\underset{\mathrm{x}}{\vert} \quad \underset{\mathrm{y}}{\vert}}$$

$$(A \otimes B) \circ (C \otimes D) = (A \circ C) \otimes (B \circ D) \qquad = \circ \!\! \stackrel{A}{\underset{D}{\overset{B}{\underset{C}{\bowtie}}}} \!\! \circ = = \ \stackrel{A}{\underset{D}{\overset{B}{\underset{C}{\bowtie}}}}$$

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.

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

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

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

$$A * B = \ = \!\! \stackrel{A}{\underset{B}{\diamondsuit}} \!\! \circ -$$

$$A \bullet B = \ - \circ \!\! \stackrel{A}{\underset{B}{\diamondsuit}} \!\! =$$

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:

$$(A \bullet B)(C \otimes D) = (AC) \bullet (BD)$$

$$(Ax) \circ (By) = (A \bullet C)(x \otimes y)$$

### 4.6.1 Stacking

Can be part of Kronecker section

From Josh: Proposition 2.5. For any field $\mathbb{F}$, integers $d_1, d_2, d_3, d_4$ and matrices $X_1 \in \mathbb{F}^{d_1 \times d_2}, X_2 \in \mathbb{F}^{d_2 \times d_3}, X_3 \in \mathbb{F}^{d_1 \times d_4}$, and $X_4 \in \mathbb{F}^{d_4 \times d_3}$, we have

$$X_1 \times X_2 + X_3 \times X_4 = (X_1 \mid X_3) \times \left( \frac{X_2}{X_4} \right),$$

where we are writing '|'to denote matrix concatenation.

With tensor diagrams we can write stacking along a new axis $i$ as

$$\text{stack}_i(X, Y) = \frac{e^{(0)} \dashv i}{-X-} + \frac{e^{(1)} \dashv i}{-Y-}$$

where $e_i^{(i)} = 1$ and 0 elsewhere.

Fro this we easily get the identity

$$(A \mid C) \left( \frac{B}{D} \right) = \text{stack}_i(A, C) \, \text{stack}_i(B, D) \tag{4.2}$$

$$= \left( \frac{e^{(0)} \dashv i}{k - A - j} + \frac{e^{(1)} \dashv i}{k - C - j} \right) \left( \frac{e^{(0)} \dashv i}{j - B - m} + \frac{e^{(1)} \dashv i}{j - D - m} \right) \tag{4.3}$$

$$= \frac{e^{(0)} - - e^{(0)}}{k - A - - - B - m} + \frac{e^{(1)} - - e^{(1)}}{k - C - - - D - m} \tag{4.4}$$

$$= AB + CD \tag{4.5}$$

TODO: Relation to direct sum, which is basically stacking + flattening. Also have a bunch of properties like $A \otimes (B \oplus C) = A \otimes B \oplus A \otimes C$.

# Chapter 5

# Functions

| | | | |
|---|---|---|---|
| Function to scalar | $f : \mathbb{R}^n \to \mathbb{R}$ | $f(x) \in \mathbb{R}$ | $\{f \leftarrow x\}$ |
| Function to vector | $g : \mathbb{R}^n \to \mathbb{R}^m$ | $g(x) \in \mathbb{R}^m$ | $\{g \overset{/}{\leftarrow} x\}$ |
| Element-wise function | $h : \mathbb{R}^n \to \mathbb{R}$ | $h(x) \in \mathbb{R}^m$ | $\{h \quad \overset{/}{x}\}$ |
| Vector times vector function | $u : \mathbb{R}^n \to \mathbb{R}^m, v \in \mathbb{R}^m$ | $v^T u(x) \in \mathbb{R}$ | $v \ \overset{\frown}{\{u} \leftarrow x\}$ |
| Vector times matrix function | $A : \mathbb{R}^n \to \mathbb{R}^{m \times n}, v \in \mathbb{R}^n$ | $A(x)v \in \mathbb{R}^m$ | $\overset{m}{\{A} \leftarrow x\} \overset{\frown}{v}$ |
| Batch function application | $f : \mathbb{R}^d \to \mathbb{R}, X \in \mathbb{R}^{b \times d}$ | $f(X) \in \mathbb{R}^b$ | $\{f \leftarrow \overset{\flat}{X}\}$ |

As an example of a more complicated function, let $\exp : \mathbb{R} \to \mathbb{R}$ be the element-wise exponential function, and $\mathrm{pow}^{-1} : \mathbb{R} \to \mathbb{R}$ be the element-wise inverse power function. Then we can write $\mathrm{softmax} : \mathbb{R}^d \to \mathbb{R}^d$ as the tensor diagram:

$$\mathrm{softmax}(x) = \{\exp \quad \overset{/}{x}\} \quad \{\mathrm{pow}^{-1} \quad \{\exp \quad \overset{\frown}{x}\} \ \circ\} \,.$$

Note in particular how $\{\exp \quad \overset{\frown}{x}\} \ \circ$ is the diagram way of writing $\sum_i \exp(x_i)$. Alternatively we could have used a function $\mathrm{sum} : \mathbb{R}^n \to \mathbb{R}$, but the more we can express in terms of tensor diagrams, the more we can use the powerful tensor diagram calculus.

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

## 5.1   The Chain Rule

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

E.g. $f : \mathbb{R}^n -> \mathbb{R}^n$

Standard chain rule. Here we let $f \in \mathbb{R}^d \to \mathbb{R}$ be a scalar function, and $v \in \mathbb{R}^d \to \mathbb{R}^d$ be a vector function as used in backprop.

Visualization of the Chain Rule: $J_{f \circ v}(x) = \nabla_f(v(x)) J_v(x).$

$$( \quad \{ \ \widehat{f} \ \{ v \leftarrow x \}\} \quad \}' \ = \{ \ \widehat{f_\bullet \{v \leftarrow x}\}\} \quad \{ v \overset{\bullet}{\leftarrow} x \}$$

### 5.1.1   The Chain Rule

Let $f \in \mathbb{R}^d \to \mathbb{R}$ be a scalar function, and $v \in \mathbb{R}^d \to \mathbb{R}^d$ be a vector function, as used in backprop. Then we can write the chain rule as:



Using standard notation: $J_{f \circ v}(x) = \nabla_f(v(x)) J_v(x).$

The second derivative, (or Hessian Chain rule):



Using standard notation: $H_{f \circ v}(x) = Dv(x)^T \cdot D^2 f(v(x)) \cdot Dv(x) + \sum_{k=1}^{d} \frac{\partial f}{\partial u_k}(v(x)) \frac{\partial^2 v_k}{\partial x \partial x^T}(x).$

$$\frac{\partial A(x)x}{\partial x} = (x^T \otimes I)\frac{\partial}{\partial x}\text{vec}[A(x)] + A(x) \quad ( \ \widehat{x \ \{ \ A \leftarrow x \ \}}' \ = \widehat{x \ (\{ \ A \leftarrow x \ \}}' + ( \ \widehat{x \ \}' \{ \ A \leftarrow x \ \}}$$

$$= \widehat{x \ \{ \ A \overset{\bullet}{\leftarrow} x \ \}} + \{ \ A \overset{\bullet}{\leftarrow} x \ \}$$



are common.  All pixel-adaptive filters like non-local means, bilateral, etc, and the so-called attention mechanism in transformers can be written this way

Gradient of this f(x) is important & has a form worth remembering. . .

**Pseudo-linear form**

Maybe this should just be an example in a table?

Derivation of Peyman Milanfar's gradient

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

## 5.1.2 Taylor

For an n-times differentiable function $v : \mathbb{R}^d \to \mathbb{R}^d$ we can write the Taylor expansion:

$$
v(x+\varepsilon) \approx v(x) + \left[\frac{\partial}{\partial x}v(x)\right]\varepsilon + \frac{1}{2}\left[\frac{\partial}{\partial x}\left[\frac{\partial}{\partial x}v(x)\right]\varepsilon\right]\varepsilon + \frac{1}{6}\left[\frac{\partial}{\partial x}\left[\frac{\partial}{\partial x}\left[\frac{\partial}{\partial x}v(x)\right]\varepsilon\right]\varepsilon\right]\varepsilon + \dots
$$

$$
= v(x) + \left[\frac{\partial}{\partial x}v(x)\right]\varepsilon + \frac{1}{2}(I \otimes \varepsilon)\left[\frac{\partial\text{vec}}{\partial x}\left[\frac{\partial v(x)}{\partial x}\right]\right]\varepsilon + \frac{1}{6}(I \otimes \varepsilon \otimes \varepsilon)\left[\frac{\partial\text{vec}}{\partial x}\left[\frac{\partial\text{vec}}{\partial x}\left[\frac{\partial v(x)}{\partial x}\right]\right]\right]\varepsilon + \dots
$$

Or with indices:

$$
v_i(x+\varepsilon) \approx v_i(x) + \sum_j \frac{\partial v_i(x)}{\partial x_j}\varepsilon_j + \frac{1}{2}\sum_{j,k}\frac{\partial v_i(x)}{\partial x_j \partial x_k}\varepsilon_j\varepsilon_k + \frac{1}{6}\sum_{j,k,\ell}\frac{\partial v_i(x)}{\partial x_j \partial x_k \partial x_\ell}\varepsilon_j\varepsilon_k\varepsilon_\ell
$$

Or diagrams:



TODO: Examples based on idempotent matrices etc.

# Chapter 6

# Determinant and Inverses

## 6.1 Determinant

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

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

where $\underline{\quad^{i_1} \quad^{i_2} \quad \cdots \quad^{i_n}\quad} = \varepsilon_{i_1,\ldots,i_n}$ is the rank-$n$ Levi-Civita tensor defined by

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

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

$$\det(I) = \frac{1}{n!} \; \boxed{\quad \cdots \quad} = \frac{1}{n!} \sum_{i_1,\ldots,i_n,j_1,\ldots,j_n} \varepsilon_{i_1,\ldots,i_n} \varepsilon_{j_1,\ldots,j_n}[i = j] = \frac{1}{n!} \sum_{i_1,\ldots,i_n} \varepsilon^2_{i_1,\ldots,i_n} = 1.$$

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

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

The definition generalizes to Cayley's "hyper determinants" by ....

A curious property is that

$$\underline{A} \;\cdots\; \underline{A} \;=\; \overline{\underline{A} \;\cdots\; \underline{A}}$$

(18)     $\det(A) = \prod_i \lambda_i$                     $\cdots$

(19)     $\det(cA) = c^n \det(A)$     $\overline{cA \;\cdots\; cA} = c^n \overline{A \;\cdots\; A}$

(20)     $\det(A) = \det(A^T)$                     $\cdots$

(21)   $\det(AB) = \det(A)\det(B)$     $\overline{\underline{\begin{array}{c}A \;\cdots\; A\\ B \;\cdots\; B\end{array}}} = \overline{\underline{\begin{array}{c}A \;\cdots\; A\\ B \;\cdots\; B\end{array}}}$

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

(23)     $\det(A^n) = \det(A)^n$                     $\cdots$

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

## 6.2   Inverses

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

# Chapter 7

# Advanced Derivatives

## 7.1 Derivatives of vector norms

### 7.1.1 Two-norm

$$\frac{d}{dx}\|x - a\|_2 = \frac{x - a}{\|x - a\|_2} \tag{7.1}$$

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

$$\frac{d}{dx}\|x\|_2^2 = \frac{d}{dx}\|x^T x\|_2 = 2x \tag{7.3}$$

## 7.2   Derivatives of matrix norms

## 7.3   Derivatives of Structured Matrices

### 7.3.1   Symmetric

### 7.3.2   Diagonal

### 7.3.3   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: $H_n = \left[\begin{smallmatrix} 1 & 1 \\ 1 & -1 \end{smallmatrix}\right]^{\otimes n}$. 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  Fast Kronecker Multiplication

Say we want to compute $(A_1 \otimes A_2 \cdots A_n)x$, where $A_i$ is a $a_i \times a_i$ matrix, and $x \in \mathbb{R}^{a_1 a_2 \cdots a_n}$. If we first compute the Kronecker product, and then the matrix-vector multiplication, this would take $(a_1 \cdots a_n)^2$ time.

Instead we can reshape $x$ into a $a_1 \times \cdots a_n$ tensor and perform the multiplication



by contracting the $a_i$ edges one by one. This takes time

$$a_1^2(a_2 \cdots a_n) + a_2^2(a_1 a_3 \cdots a_n) + \cdots + a_n^2(a_1 \cdots a_{n-1}) = (a_1 + \cdots + a_n)(a_1 \cdots a_n),$$

which is the basis of many fast algorithm as we will see.

### Hadamard

The Hadamard matrix is defined as $H_{2^n} = H_2^{\otimes n} = H_2 \otimes \cdots \otimes H_2$ where $H_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$.

For example

$$H_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}.$$

This gives a very simple tensor diagram representation as:



The Fast Hadamard Transform (FHT) transform is usually described recursively by:

$$H_{2^n} x = \begin{bmatrix} H_{2^{n-1}} & H_{2^{n-1}} \\ H_{2^{n-1}} & -H_{2^{n-1}} \end{bmatrix} \begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix},$$

where $\begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix}$ is the first and second half of $x$. Because of the redundancy in the matrix multiplication (it only depends on $H_{2^{n-1}} x^{(1)}$ and $H_{2^{n-1}} x^{(2)}$, the algorithm computes $H_N x$ in $O(N \log N)$ time.

Alternatively we could just use the general fact, as described above, where $a_i = 2$ for all $i$. Then the "fast Kronecker multiplication" method takes time $(a_1 a_2 \cdots a_n)(a_1 + a_2 + \cdots a_n) = 2n \log_2 n$.

**Fourier**

A more common matrix is the Discrete Fourier Matrix. It is defined by $(F_N)_{i,j} = \omega^{ij}$, where $\omega = e^{-2\pi i/N}$. We can write it out as

$$
F_N = \begin{bmatrix}
1 & 1 & 1 & 1 & \cdots & 1 \\
1 & \omega & \omega^2 & \omega^3 & \cdots & \omega^{N-1} \\
1 & \omega^2 & \omega^4 & \omega^6 & \cdots & \omega^{2(N-1)} \\
1 & \omega^3 & \omega^6 & \omega^9 & \cdots & \omega^{3(N-1)} \\
\vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\
1 & \omega^{N-1} & \omega^{2(N-1)} & \omega^{3(N-1)} & \cdots & \omega^{(N-1)(N-1)}
\end{bmatrix}.
$$

The Good-Thomas Fast Fourier Transformer (FFT) uses a decomposition based on the Chinese Remainder Theorem:



where $N = p_1^{i_1} p_2^{i_2} \cdots p_n^{i_n}$ is the prime factorisation of $N$, and $P_1$ and $P_2$ are some permutation matrices.

Using fast Kronecker multiplication, the algorithm this takes $(p_1^{i_1} + \cdots + p_n^{i_n})N$ time. By padding $x$ with zeros, we can increase $N$ by a constant factor to get a string of $n = O(\log(N)/\log\log(N))$ primes, the sum of which is $\sim n^2/\log n = O(\log(N)^2)$. The complete algorithm thus takes time $O(N\log(N)^2)$. Next we will see how to reduce this to $O(N\log N)$.

The classical Cooley-Tukey FFT algorithm uses a recursion:

$$
F_N = \begin{bmatrix} I & I \\ I & -I \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & D_{N/2} \end{bmatrix} \begin{bmatrix} F_{N/2} & 0 \\ 0 & F_{N/2} \end{bmatrix} \begin{bmatrix} \text{even-odd} \\ \text{permutation} \end{bmatrix},
$$

where $D_N = [1, w^N, w^{2N}, \dots]$. The even-odd permutation moves all the even values to the start. If we reshape $I_{2^n}$ as $I_2 \otimes \cdots \otimes I_2$, this permutation is just $P_N = \times$, or in pytorch: `x.permute([3,0,1,2])`. Also note that $\begin{bmatrix} I & I \\ I & -I \end{bmatrix} = H_2 \otimes I$ and $\begin{bmatrix} F_{N/2} & 0 \\ 0 & F_{N/2} \end{bmatrix} = I_2 \otimes F_{N/2}$. So we can write in tensor diagram notation:

Since one can multiply with the permutation and diagonal matrices in linear time, the $O(n \log n)$ time complexity follows from the same argument as for Hadamard.

TODO: We can also talk about the row/column version of Cooley-Tukey (Bailey's FFT algorithm). The nice thing is that we can compose arbitrarily (not necessarily in $H_2$ chunks), and the $D_N$ matrices have a nice form (they are just fourier matrices!) So they could in principle be decomposed too, by the same algorithm...

Note that this figure may look different from some FFT diagrams you have seen. These typically look like this:



and have $2^n$ rows. The tensor diagram only has $n$ rows (or $\log_2 N$).

## Multi-dimensional Fourier Transform

This is just taking the Fourier transform along each axis.

### 8.0.4   Hermitian Matrices and skew-Hermitian

Complex. Skip

### 8.0.5   Idempotent Matrices

Skip

### 8.0.6   Orthogonal matrices

Skip

### 8.0.7   Positive Definite and Semi-definite Matrices

Skip

### 8.0.8   Singleentry Matrix, The

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

### 8.0.9   Symmetric, Skew-symmetric/Antisymmetric

Could introduce Penrose's symmetric tensors here?

### 8.0.10   Toeplitz Matrices

Could talk about the convolution tensor here...

### 8.0.11   Units, Permutation and Shift

Not that interesting...

### 8.0.12   Vandermonde Matrices

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

# Chapter 9

# Decompositions

## 9.1 Higher-order singular value decomposition

Say we have an order $n$ tensor $A$. We "unfold" $A$ along each dimension. This means pulling the edge $i$ to the left, and flattening the rest to the right. Then we compute the SVD, $USV$. Here $U$ is a square matrix, which we keep. We multiply the $i$th edge of $A$ by $U^T$ (which is also the inverse of $U$). The result is a "core" tensor as well as a sequence of $U$ tensors. If we want a more compact SVD, we can make each $U$ low rank, like normal SVD. There is also the "Interlacing computation" where we multiply the $U^T$s onto $A$ as we go along.

For order 3 tensors, this method is called a "Tucker decomposition".

If the "core matrix" is diagonal, this is called tensor rank decomposition. If we were good at that, we could use it to factor $I^{\otimes 3}$ to get better matrix multiplication algorithms. Unfortunately tensor rank decomposition is NP hard.

I guess HOSVD gives a rank decomposition if we diagonalize the core tensor. It just won't be an efficient one.

## 9.2 Rank Decomposition

## 9.3 Fast Matrix Multiplication

Strassen defines 3 tensors of shape $7 \times 2 \times 2$:

$$S_A = \left[ \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} -1 & 0 \\ 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 0 & -1 \end{bmatrix} \right]$$

$$S_B = \left[ \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 0 & -1 \end{bmatrix}, \begin{bmatrix} -1 & 0 \\ 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix} \right]$$

$$W = \left[ \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 0 & -1 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} -1 & 0 \\ 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \right]$$

These tensors have the neat property that they factor $I_2 \otimes I_2 \otimes I_2$:



To multiply two matrices, $A$ and $B$, faster than the normal $n^3$ time, we reshape them as block matrices, shape $(2, \frac{n}{2}, 2, \frac{n}{2})$ and use Strassen's tensor:



Contracting the edges in the right order, uses only $7/8n^3 + O(n^2)$ operations.

If we instead reshape to $(2, 2, \ldots, 2)$,



and using Strassen's tensor along each axis reduces the work by $(7/8)^{\log_2(n)}$, giving us matrix multiplication in time $n^{3+\log_2(7/8)} = n^{2.80735}$.

Contracting the double edges, $S_A - A$ and $S_B - B$, is both $O(n^2)$ time.

It remains to verify that this is actually faster than the naive matrix multiplication: Contracting $S_A - A$ takes $7 \cdot 2^2 (n/2)^2$ operations, and likewise $S_B - B$. Next we contract $S_A A - S_B B$ which takes $7(n/2)^3$ time. And finally we contract the edge with $W$ which takes $2^2 \cdot 7(n/2)^2$. The important term is the cubic $7/8 n^3$, which if instead done recursively, leads to the classical $O(n^{\log_2 7})$ algorithm.

FIXME: What "edge with $W$"? I think we have to/want to contract the hyperedge with $W$ immediately?

### Other

If we instead wrote $A$ and $B$ using $(n, m)$ and $(m, p)$ shaped blocks, we could factor $I_n \otimes I_m \otimes I_p$ and get a matrix multiplication algorithm using the same approach as the Strassen $(2, 2, 2)$ tensors above. Lots of papers have investigated this problem, which has led to the best algorithms by Josh Alman and others. For example, Deep Mind found a rank 47 factorization of $I_3 \otimes I_4 \otimes I_5$.

Maybe a more interesting example is the $(4, 4, 4)$ tensor, for which they find a rank 47 factorization. This an easy way to create a rank 49 is to take Strassen and double it. Would this be a nice thing to show? Maybe too messy? Well, actually their rank 47 construction only works in the "modular" case. Then $(3, 4, 5)$ is general.

# Chapter 10

# Machine Learning Applications

# Chapter 11

# Tensorgrad

Implementation details

## 11.1   Simplification Rules

## 11.2   Functions

## 11.3   Isomorphisms

# Chapter 12

# 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.