

# Named Tensor Notation

David Chiang and Sasha Rush

December 11, 2020

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Informal Overview</b>	<b>3</b>
2.1	Named tensors . . . . .	3
2.2	Named tensor operations . . . . .	4
<b>3</b>	<b>Examples</b>	<b>8</b>
3.1	Building blocks . . . . .	8
3.2	Discrete random variables . . . . .	11
3.3	Transformer . . . . .	11
3.4	LeNet . . . . .	13
3.5	Other examples . . . . .	14
<b>4</b>	<b><math>\LaTeX</math> Macros</b>	<b>16</b>
<b>5</b>	<b>Formal Definitions</b>	<b>17</b>
5.1	Records . . . . .	17
5.2	Named tensors . . . . .	17
5.3	Extending functions to named tensors . . . . .	18
<b>6</b>	<b>Extensions</b>	<b>19</b>
6.1	Index types . . . . .	19
6.2	Duality . . . . .	20
6.3	Indexing with a tensor of indices . . . . .	23

# 1 Introduction

Most papers about neural networks use the notation of vectors and matrices from applied linear algebra. This notation is very well-suited to talking about vector spaces, but less well-suited to talking about neural networks. Consider the following equation (Vaswani et al., 2017):

$$\text{Att}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{\sqrt{d_k}} \right) V.$$

where  $Q$ ,  $K$ , and  $V$  are sequences of query, key, and value vectors packed into matrices. Does the product  $QK^\top$  sum over the sequence, or over the query/key features? We would need to know the sizes of  $Q$ ,  $K$ , and  $V$  to know that it's taken over the query/key features. Is the softmax taken over the query sequence or the key sequence? The usual notation doesn't even offer a way to answer this question. With multiple attention heads, the notation becomes more complicated and leaves more questions unanswered. With multiple sentences in a minibatch, the notation becomes more complicated still, and most papers wisely leave this detail out.

Libraries for programming with neural networks (Harris et al., 2020; Paszke et al., 2019) provide multidimensional arrays, called tensors (although usually without the theory associated with tensors in linear algebra and physics), and a rich array of operations on tensors. But they inherit from math the convention of identifying indices by *position*, making code bug-prone. Quite a few libraries have been developed to identify indices by *name* instead: Nexus (Chen, 2017), tsalib (Sinha, 2018), NamedTensor (Rush, 2019), named tensors in PyTorch (Torch Contributors, 2019), and Dex (Maclaurin et al., 2019). (Some of these libraries also add types to indices, but here we are only interested in adding names.)

Back in the realm of mathematical notation, then, we want two things: first, the flexibility of working with multidimensional arrays, and second, the perspicuity of identifying indices by name instead of by position. This document describes our proposal to do both.

As a preview, the above equation becomes

$$\text{Att}(Q, K, V) = \underset{\text{seq}}{\text{softmax}} \left( \frac{Q \underset{\text{key}}{\cdot} K}{\sqrt{d_k}} \right) \underset{\text{seq}}{\cdot} V$$

making it unambiguous which index each operation applies to. The same equation works with multiple heads and with minibatching.

More examples of the notation are given in §3.

The source code for this document can be found at <https://github.com/namedtensor/notation/>. We invite anyone to make comments on this proposal by submitting issues or pull requests on this repository.

## 2 Informal Overview

Let's think first about the usual notions of vectors, matrices, and tensors, without named indices.

Define  $[n] = \{1, \dots, n\}$ . We can think of a size- $n$  real vector  $v$  as a function from  $[n]$  to  $\mathbb{R}$ . We get the  $i$ th element of  $v$  by applying  $v$  to  $i$ , but we normally write this as  $v_i$  (instead of  $v(i)$ ).

Similarly, we can think of an  $m \times n$  real matrix as a function from  $[m] \times [n]$  to  $\mathbb{R}$ , and an  $l \times m \times n$  real tensor as a function from  $[l] \times [m] \times [n]$  to  $\mathbb{R}$ . In general, then, real tensors are functions from *tuples of natural numbers* to reals.

### 2.1 Named tensors

We want to make tensors into functions, no longer on tuples, but on *records*, which look like this:

$$\{\text{foo}(1), \text{bar}(3)\}$$

where `foo` and `bar` are *names* (written in sans-serif font), mapped to 1 and 3, respectively. The pairs `foo(1)` and `bar(3)` are called *named indices*. Their order doesn't matter:  $\{\text{foo}(1), \text{bar}(3)\}$  and  $\{\text{bar}(3), \text{foo}(1)\}$  are the same record.

The set of records that can be used to index a named tensor is defined by a *named shape*, which looks like this:

$$\text{foo}[2] \times \text{bar}[3]$$

which stands for records where index `foo` ranges from 1 to 2, and index `bar` ranges from 1 to 3. Again, the order of the elements, called *named index sets*, doesn't matter:  $\text{foo}[2] \times \text{bar}[3]$  and  $\text{bar}[3] \times \text{foo}[2]$  are the same shape.

Then, a real *named tensor* is a function from (the set of records defined by) a named shape to the real numbers. For example, here is a tensor with shape  $\text{foo}[2] \times \text{bar}[3]$ .

$$A = \text{foo} \begin{matrix} & \text{bar} \\ \begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 & 9 \end{bmatrix} \end{matrix}.$$

We access elements of  $A$  using subscripts:  $A_{\text{foo}(1), \text{bar}(3)} = 4$ . We also allow partial indexing:

$$A_{\text{foo}(1)} = \begin{matrix} & \text{bar} \\ \begin{bmatrix} 3 & 1 & 4 \end{bmatrix} \end{matrix} \qquad A_{\text{bar}(3)} = \begin{matrix} \text{foo} \\ \begin{bmatrix} 4 & 9 \end{bmatrix} \end{matrix}.$$

We use uppercase italic letters for variables standing for named tensors. We don't mind if you use another convention, but urge you not to use different styles for tensors and their elements. For example, if  $\mathbf{A}$  is a tensor, then an element of  $\mathbf{A}$  is written as  $\mathbf{A}_{\text{foo}(2), \text{bar}(3)}$  – not  $A_{\text{foo}(2), \text{bar}(3)}$  or  $a_{\text{foo}(2), \text{bar}(3)}$ .

Just as the set of all size- $n$  real vectors is written  $\mathbb{R}^n$ , and the set of all  $m \times n$  real matrices is often written  $\mathbb{R}^{m \times n}$  (which makes sense because one sometimes writes  $Y^X$  for the set of all functions from  $X$  to  $Y$ ), we write  $\mathbb{R}^{\text{foo}[2] \times \text{bar}[3]}$  for the set of all tensors with shape `foo[2] × bar[3]`.

It’s very common for an index name to be used with only one size. If you write

$$\begin{aligned}\text{foo} &\stackrel{\text{ind}}{=} [2] \\ \text{bar} &\stackrel{\text{ind}}{=} [3]\end{aligned}$$

then you can simply write  $\mathbb{R}^{\text{foo} \times \text{bar}}$  for the set of tensors with index `foo` ranging from 1 to 2 and `bar` ranging from 1 to 3. For more information on what  $\stackrel{\text{ind}}{=}$  does, please see §5.

What are good choices for index names? We recommend meaningful *words* instead of single letters, and we recommend words that describe a *whole* rather than its parts. For example, a minibatch of sentences, each of which is a sequence of one-hot vectors, would be represented by a tensor with three indices, which we might name `batch`, `seq`, and `vocab`. Please see §3 for more examples.

## 2.2 Named tensor operations

### 2.2.1 Elementwise operations

Any function from scalars to scalars can be applied elementwise to a named tensor:

$$\text{exp } A = \text{foo} \begin{array}{c} \text{bar} \\ \left[ \begin{array}{ccc} \text{exp } 3 & \text{exp } 1 & \text{exp } 4 \\ \text{exp } 1 & \text{exp } 5 & \text{exp } 9 \end{array} \right] \end{array}.$$

More elementwise unary operations:

$kA$	scalar multiplication by $k$
$-A$	negation
$A^k$	elementwise exponentiation
$\sqrt{A}$	elementwise square root
$\text{exp } A$	elementwise exponential function
$\tanh A$	hyperbolic tangent
$\sigma(A)$	logistic sigmoid
$\text{ReLU}(A)$	rectified linear unit

Any function or operator that takes two scalar arguments can be applied elementwise to two named tensors with the same shape. If  $A$  is as above and

$$B = \text{foo} \begin{array}{c} \text{bar} \\ \left[ \begin{array}{ccc} 2 & 7 & 1 \\ 8 & 2 & 8 \end{array} \right] \end{array}$$

then

$$A + B = \text{foo} \begin{matrix} & \text{bar} \\ \begin{bmatrix} 3 + 2 & 1 + 7 & 4 + 1 \\ 1 + 8 & 5 + 2 & 9 + 8 \end{bmatrix} \end{matrix}.$$

But things get more complicated when  $A$  and  $B$  don't have the same shape. If  $A$  and  $B$  each have an index with the same name (and size), the two indices are *aligned*, as above. But if  $A$  has an index named  $i$  and  $B$  doesn't, then we do *broadcasting*, which means effectively that we replace  $B$  with a new tensor  $B'$  that contains a copy of  $B$  for every value of index  $i$ .

$$\begin{aligned} A + 1 &= \text{foo} \begin{matrix} & \text{bar} \\ \begin{bmatrix} 3 + 1 & 1 + 1 & 4 + 1 \\ 1 + 1 & 5 + 1 & 9 + 1 \end{bmatrix} \end{matrix} \\ A + B_{\text{foo}(1)} &= \text{foo} \begin{matrix} & \text{bar} \\ \begin{bmatrix} 3 + 2 & 1 + 7 & 4 + 1 \\ 1 + 2 & 5 + 7 & 9 + 1 \end{bmatrix} \end{matrix} \\ A + B_{\text{bar}(3)} &= \text{foo} \begin{matrix} & \text{bar} \\ \begin{bmatrix} 3 + 1 & 1 + 1 & 4 + 1 \\ 1 + 8 & 5 + 8 & 9 + 8 \end{bmatrix} \end{matrix}. \end{aligned}$$

Similarly, if  $B$  has an index named  $i$  and  $A$  doesn't, then we effectively replace  $A$  with a new tensor  $A'$  that contains a copy of  $A$  for every value of index  $i$ . If you've programmed with NumPy or any of its derivatives, this should be unsurprising to you.

More elementwise binary operations:

$A + B$	addition
$A - B$	subtraction
$A \odot B$	elementwise (Hadamard) product
$\frac{A}{B}$	elementwise division
$\max\{A, B\}$	elementwise maximum
$\min\{A, B\}$	elementwise minimum

### 2.2.2 Reductions

The same rules for alignment and broadcasting apply to functions that take tensor as arguments or return tensors. The gory details are in §5.3, but we present the most important subcases here. The first is *reductions*, which are functions from vectors to scalars. Unlike with functions on scalars, we always have to specify which index these functions apply to, using a subscript. (This is equivalent to the `axis` argument in NumPy and `dim` in PyTorch.)

For example, using the same example tensor  $A$  from above,

$$\sum_{\text{foo}} A = \begin{matrix} & \text{bar} \\ [3+1 & 1+5 & 4+9] \end{matrix}$$

$$\sum_{\text{bar}} A = \begin{matrix} & \text{foo} \\ [3+1+4 & 1+5+9] \end{matrix}.$$

More reductions: If  $A$  has shape  $\text{foo}[I] \times \dots$ , then

$$\sum_{\text{foo}} A = \sum_{i \in I} A_{\text{foo}(i)} = \begin{matrix} & \text{bar} \\ [4 & 6 & 13] \end{matrix}$$

$$\text{norm}_{\text{foo}} A = \sqrt{\sum_{\text{foo}} A^2} = \begin{matrix} & \text{bar} \\ [\sqrt{10} & \sqrt{26} & \sqrt{97}] \end{matrix}$$

$$\min_{\text{foo}} A = \min_{i \in I} A_{\text{foo}(i)} = \begin{matrix} & \text{bar} \\ [1 & 1 & 4] \end{matrix}$$

$$\max_{\text{foo}} A = \max_{i \in I} A_{\text{foo}(i)} = \begin{matrix} & \text{bar} \\ [3 & 5 & 9] \end{matrix}$$

$$\text{mean}_{\text{foo}} A = \frac{1}{|I|} A = \begin{matrix} & \text{bar} \\ [2 & 3 & 6.5] \end{matrix}$$

$$\text{var}_{\text{foo}} A = \frac{1}{|I|} \sum_i (A - \text{mean}_{\text{foo}} A)^2 = \begin{matrix} & \text{bar} \\ [1 & 4 & 6.25] \end{matrix}$$

(Note that  $\max$  and  $\min$  are overloaded; with multiple arguments and no subscript, they are elementwise, and with a single argument and a subscript, they are reductions.)

You can also write multiple names to perform the reduction over multiple indices at once.

### 2.2.3 Contraction

The vector dot product (inner product) is a function from *two* vectors to a scalar, which generalizes to named tensors to give the ubiquitous *contraction* operator, which performs elementwise multiplication, then sums along an index. It can

be used, for example, for matrix multiplication:

$$C = \text{bar} \begin{bmatrix} 1 & -1 \\ 2 & -2 \\ 3 & -3 \end{bmatrix}$$

$$A \underset{\text{bar}}{\cdot} C = \text{foo} \begin{bmatrix} 17 & -17 \\ 53 & -53 \end{bmatrix}$$

However, note that (like vector dot-product, but unlike matrix multiplication) this operator is commutative, but not associative! Specifically, if

$$A \in \mathbb{R}^{\text{foo}[m]}$$

$$B \in \mathbb{R}^{\text{foo}[m] \times \text{bar}[n]}$$

$$C \in \mathbb{R}^{\text{foo}[m] \times \text{bar}[n]}$$

then  $(A \underset{\text{foo}}{\cdot} B) \underset{\text{bar}}{\cdot} C$  and  $A \underset{\text{foo}}{\cdot} (B \underset{\text{bar}}{\cdot} C)$  don't even have the same shape.

#### 2.2.4 Vectors to vectors

A very common example of a function from vectors to vectors is the softmax:

$$\text{softmax}_{\text{foo}} A = \frac{\exp A}{\sum_{\text{foo}} \exp A} \approx \text{foo} \begin{array}{ccc} \text{bar} & & \\ \begin{bmatrix} 0.731 & 0.002 & 0.953 \\ 0.269 & 0.998 & 0.047 \end{bmatrix} \end{array}$$

And it's also very handy to have a function that renames an index:

$$[A]_{\text{bar} \rightarrow \text{baz}} = \text{foo} \begin{array}{ccc} \text{baz} & & \\ \begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 & 9 \end{bmatrix} \end{array}$$

Concatenation combines two vectors into one:

$$A \underset{\text{foo}}{\oplus} B = \text{foo} \begin{array}{ccc} \text{bar} & & \\ \begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 & 9 \\ 2 & 7 & 1 \\ 8 & 2 & 8 \end{bmatrix} \end{array}$$

$$A \underset{\text{bar}}{\oplus} B = \text{foo} \begin{array}{cccccc} \text{bar} & & & & & \\ \begin{bmatrix} 3 & 1 & 4 & 2 & 7 & 1 \\ 1 & 5 & 9 & 8 & 2 & 8 \end{bmatrix} \end{array}$$

### 2.2.5 Matrices

Finally, we briefly consider functions on matrices, for which you have to give *two* index names (and the order in general matters). Let  $A$  be a named tensor with shape  $\{\text{foo}[2] \times \text{bar}[2] \times \text{baz}[2]\}$ :

$$\begin{aligned}
 A_{\text{foo}(1)} &= \text{bar} \begin{matrix} \text{baz} \\ \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \end{matrix} \\
 A_{\text{foo}(2)} &= \text{bar} \begin{matrix} \text{baz} \\ \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \end{matrix} \\
 \det_{\text{bar}, \text{baz}} A &= \begin{matrix} \text{foo} \\ \begin{bmatrix} \det \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} & \det \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \end{bmatrix} \end{matrix} \\
 \det_{\text{baz}, \text{bar}} A &= \begin{matrix} \text{foo} \\ \begin{bmatrix} \det \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} & \det \begin{bmatrix} 5 & 7 \\ 6 & 8 \end{bmatrix} \end{bmatrix} \end{matrix} \\
 \det_{\text{foo}, \text{bar}} A &= \begin{matrix} \text{baz} \\ \begin{bmatrix} \det \begin{bmatrix} 1 & 3 \\ 5 & 7 \end{bmatrix} & \det \begin{bmatrix} 2 & 4 \\ 6 & 8 \end{bmatrix} \end{bmatrix} \end{matrix}
 \end{aligned}$$

For matrix inverses, there's no easy way to put a subscript under  $\cdot^{-1}$ , so we recommend writing  $\text{inv}_{\text{foo}, \text{bar}}$ .

## 3 Examples

In this section we give a series of examples illustrating how to use named tensors in various situations, mostly related to machine learning.

### 3.1 Building blocks

#### 3.1.1 Fully-connected layers

We use rectified linear units (ReLUs) because some of the larger models we define later require it; a logistic function or tanh would also work here.

$$\begin{aligned}
 \text{Full}: \mathbb{R}^{\text{layer}} &\rightarrow \mathbb{R}^{\text{layer}'} \\
 \text{Full}(x; W, b) &= \max \left\{ 0, W_{\text{layer}} \cdot x + b \right\}
 \end{aligned}$$

where

$$\begin{aligned}
 W &\in \mathbb{R}^{\text{layer}' \times \text{layer}} \\
 b &\in \mathbb{R}^{\text{layer}'}
 \end{aligned}$$



### 3.1.2 Recurrent neural networks

As a second example, let's define a simple (Elman) RNN. Let  $d$  be a positive integer.

$$\begin{aligned}
\text{state}, \text{state}' &\stackrel{\text{ind}}{=} [d] \\
x^{(t)} &\in \mathbb{R}^{\text{emb}} & t = 1, \dots, n \\
h^{(t)} &\in \mathbb{R}^{\text{state}} & t = 0, \dots, n \\
A &\in \mathbb{R}^{\text{state} \times \text{state}'} \\
B &\in \mathbb{R}^{\text{emb} \times \text{state}'} \\
c &\in \mathbb{R}^{\text{state}'} \\
h^{(t+1)} &= \left[ \tanh \left( A \underset{\text{state}}{\cdot} h^{(t)} + B \underset{\text{emb}}{\cdot} x^{(t)} + c \right) \right]_{\text{state}' \rightarrow \text{state}}
\end{aligned}$$

Here the index name **state** has to stay the same across time, and the renaming is necessary because our notation doesn't provide a one-step way to apply a linear transformation ( $A$ ) to one index and put the result in the same index. For possible solutions, see §6.2.

### 3.1.3 Attention

Let  $d_k$  and  $d_v$  be positive integers, and let  $n$  be the input sequence length.

$$\begin{aligned}
\text{key} &\stackrel{\text{ind}}{=} [d_k] \\
\text{val} &\stackrel{\text{ind}}{=} [d_v] \\
\text{Att}_n &: \mathbb{R}^{\text{key}} \times \mathbb{R}^{\text{seq}[n] \times \text{key}} \times \mathbb{R}^{\text{seq}[n] \times \text{val}} \rightarrow \mathbb{R}^{\text{val}} \\
\text{Att}_n(Q, K, V) &= \underset{\text{seq}}{\text{softmax}} \left( \frac{Q \underset{\text{key}}{\cdot} K}{\sqrt{d_k}} \right) \underset{\text{seq}}{\cdot} V
\end{aligned}$$

Sometimes we need to apply a mask to keep from attending to certain positions.

$$\begin{aligned}
\text{Att}_n &: \mathbb{R}^{\text{key}} \times \mathbb{R}^{\text{seq}[n] \times \text{key}} \times \mathbb{R}^{\text{seq}[n] \times \text{val}} \times \mathbb{R}^{\text{seq}} \rightarrow \mathbb{R}^{\text{val}} \\
\text{Att}_n(Q, K, V, M) &= \underset{\text{seq}}{\text{softmax}} \left( \frac{Q \underset{\text{key}}{\cdot} K}{\sqrt{d_k}} + M \right) \underset{\text{seq}}{\cdot} V
\end{aligned}$$

If  $Q$ ,  $K$ , and  $V$  have a **head** index for multiple attention heads, then the above equations will compute multi-head attention without modification.

### 3.1.4 Convolution

A 1-dimensional convolution can be easily written by unrolling a tensor and then applying a standard dot product.

$$\begin{aligned}
\text{kernel} &\stackrel{\text{ind}}{=} [k] \\
X &\in \mathbb{R}^{\text{channels} \times \text{seq}[n]} \\
W &\in \mathbb{R}^{\text{channels} \times \text{kernel}} & b \in \mathbb{R} \\
\text{conv1d}(X; W, b) &= W \underset{\text{channels, kernel}}{\cdot} U + b \\
U &\in \mathbb{R}^{\text{channels} \times \text{seq}[n-k+1] \times \text{kernel}} \\
U_{\text{seq}(i), \text{kernel}(j)} &= X_{\text{seq}(i+j-1)}
\end{aligned}$$

A 2-dimensional convolution:

$$\begin{aligned}
kh &\stackrel{\text{ind}}{=} [kh] \\
kw &\stackrel{\text{ind}}{=} [kw] \\
X &\in \mathbb{R}^{\text{channels} \times \text{height}[h] \times \text{width}[w]} \\
W &\in \mathbb{R}^{\text{channels} \times kh \times kw} & b \in \mathbb{R} \\
\text{conv2d}(X; W, b) &= W \underset{\text{channels, kh, kw}}{\cdot} U + b \\
U &\in \mathbb{R}^{\text{channels} \times \text{height}[h-kh+1] \times \text{width}[w-kw+1] \times kh \times kw} \\
U_{\text{height}(i), \text{width}(j), kh(ki), kw(kj)} &= X_{\text{height}(i+ki-1), \text{width}(j+kj-1)}
\end{aligned}$$

### 3.1.5 Max pooling

$$\begin{aligned}
X &\in \mathbb{R}^{\text{height}[h] \times \text{width}[w]} \\
\text{maxpool2d}(X, kh, kw) &= \max_{kh, kw} U \\
U &\in \mathbb{R}^{\text{height}[h/kh] \times \text{width}[w/kw] \times kh[kh] \times kw[kw]} \\
U_{\text{height}(i), \text{width}(j), kh(di), kw(dj)} &= X_{\text{height}(i \times kh + di - 1), \text{width}(j \times kw + dj - 1)}
\end{aligned}$$

### 3.1.6 Normalization layers

Batch, instance, and layer normalization are often informally described using the same equation, but they each correspond to very different functions. They differ by which axes are normalized.

We can define a single generic normalization layer, parameterized by a named index set  $\mathbf{d}$ :

$$\begin{aligned}
&\text{xnorm}: \mathbb{R}^{\mathbf{d}} \rightarrow \mathbb{R}^{\mathbf{d}} \\
&\text{xnorm}_{\mathbf{d}}(X; \gamma, \beta, \epsilon) = \frac{X - \text{mean}_{\mathbf{d}}(X)}{\sqrt{\text{var}_{\mathbf{d}}(X) + \epsilon}} \odot \gamma + \beta
\end{aligned}$$

where

$$\begin{aligned}\gamma, \beta &\in \mathbb{R}^d \\ \epsilon &> 0\end{aligned}$$

Now, suppose that the input has three indices:

$$X \in \mathbb{R}^{\text{batch} \times \text{channels} \times \text{layer}}$$

Then the three kinds of normalization layers can be written as:

$$\begin{aligned}Y &= \underset{\text{batch}}{\text{xnorm}}(X) && \text{batch normalization} \\ Y &= \underset{\text{layer}}{\text{xnorm}}(X) && \text{instance normalization} \\ Y &= \underset{\text{layer, channels}}{\text{xnorm}}(X) && \text{layer normalization}\end{aligned}$$

### 3.2 Discrete random variables

Named indices are very helpful for working with discrete random variables, because each random variable can be represented by an index with the same name. For instance, if  $A$  and  $B$  are random variables, we can treat  $p(B | A)$  and  $p(A)$  as tensors:

$$\begin{aligned}A &\stackrel{\text{ind}}{=} [a] \\ B &\stackrel{\text{ind}}{=} [b] \\ p(B | A) &\in [0, 1]^{A \times B} && \sum_B p(B | A) = 1 \\ p(A) &\in [0, 1]^A && \sum_A p(A) = 1\end{aligned}$$

Then Bayes' rule is just:

$$p(A | B) = \frac{p(B | A) \odot p(A)}{p(B | A) \cdot_A p(A)}.$$

### 3.3 Transformer

We define a Transformer used autoregressively as a language model. The input is a sequence of one-hot vectors  $I$ , from which we compute word embeddings

and positional encodings. Let  $d_{\text{model}}$  be a positive integer.

$$\begin{aligned}
\text{layer} &\stackrel{\text{ind}}{=} [d_{\text{model}}] \\
I &\in \{0, 1\}^{\text{seq}[n] \times \text{vocab}} & \sum_{\text{vocab}} I &= 1 \\
E &\in \mathbb{R}^{\text{vocab} \times \text{layer}} \\
W &= (E \underset{\text{vocab}}{\cdot} I) \sqrt{d_{\text{model}}} \\
P &\in \mathbb{R}^{\text{seq}[n] \times \text{layer}} \\
P_{\text{seq}(p), \text{layer}(i)} &= \begin{cases} \sin((p-1)/10000^{(i-1)/d_{\text{model}}}) & i \text{ odd} \\ \cos((p-1)/10000^{(i-2)/d_{\text{model}}}) & i \text{ even} \end{cases}
\end{aligned}$$

Then we use  $L$  layers of self-attention and feed-forward neural networks.

$$\begin{aligned}
X^0 &= W + P \\
T^1 &= \underset{\text{layer}}{\text{xnorm}}(\text{SelfAtt}(X^0)) + X^0 \\
X^1 &= \underset{\text{layer}}{\text{xnorm}}(\text{FFN}(T^1)) + T^1 \\
&\vdots \\
T^L &= \underset{\text{layer}}{\text{xnorm}}(\text{SelfAtt}(X^{L-1})) + X^{L-1} \\
X^L &= \underset{\text{layer}}{\text{xnorm}}(\text{FFN}(T^L)) + T^L \\
O &= \underset{\text{vocab}}{\text{softmax}}(E \underset{\text{layer}}{\cdot} X^L)
\end{aligned}$$

where SelfAtt and FFNN are defined below.

The feedforward neural networks are:

$$\begin{aligned}
W^1 &\in \mathbb{R}^{\text{hidden}, \text{layer}} & b^1 &\in \mathbb{R}^{\text{hidden}} \\
W^2 &\in \mathbb{R}^{\text{layer}, \text{hidden}} & b^2 &\in \mathbb{R}^{\text{layer}} \\
\text{FFN}: \mathbb{R}^{\text{layer}} &\rightarrow \mathbb{R}^{\text{layer}} \\
\text{FFN}(x; W^1, b^1, W^2, b^2) &= \text{ReLU}(W^2 \underset{\text{hidden}}{\cdot} y + b^2) \\
y &= \text{ReLU}(W^1 \underset{\text{layer}}{\cdot} x + b^1)
\end{aligned}$$

We defined attention above (§3.1.3); the Transformer uses multi-head self-attention, in which  $Q$ ,  $K$ , and  $V$  are all computed from the same sequence.

The parameters are:

$$\begin{aligned}
\text{heads} &\stackrel{\text{ind}}{=} [h] \\
\text{key, val} &\stackrel{\text{ind}}{=} [d_{\text{model}}/h] \\
W^Q &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\
W^K &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\
W^V &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{val}} \\
W^O &\in \mathbb{R}^{\text{heads} \times \text{val} \times \text{layer}}
\end{aligned}$$

Then define

$$\begin{aligned}
&\text{SelfAtt}_n: \mathbb{R}^{\text{seq}[n] \times \text{layer}} \rightarrow \mathbb{R}^{\text{seq}[n] \times \text{layer}} \\
&\text{SelfAtt}_n(X; W^Q, W^K, W^V, W^O) = \sum_{\text{heads}} W^O \cdot_{\text{val}} [\text{Att}(Q, K, V, M)]_{\text{seq}' \rightarrow \text{seq}}
\end{aligned}$$

where

$$\begin{aligned}
Q &= \left[ W^Q \cdot_{\text{layer}} X \right]_{\text{seq} \rightarrow \text{seq}'} \\
K &= W^K \cdot_{\text{layer}} X \\
V &= W^V \cdot_{\text{layer}} X \\
M &\in \mathbb{R}^{\text{seq}[n] \times \text{seq}'[n]} \\
M_{\text{seq}(i), \text{seq}'(j)} &= \begin{cases} 0 & i \leq j \\ -\infty & \text{otherwise.} \end{cases}
\end{aligned}$$

### 3.4 LeNet

Parameters:

$$\begin{aligned}
W^1 &\in \mathbb{R}^{\text{channels}'[c_2] \times \text{channels}[c_1] \times \text{kh}[kh_1] \times \text{kw}[kw_1]} & b^1 &\in \mathbb{R}^{\text{channels}'[c_2]} \\
W^2 &\in \mathbb{R}^{\text{channels}'[c_3] \times \text{channels}[c_2] \times \text{kh}[kh_2] \times \text{kw}[kw_2]} & b^2 &\in \mathbb{R}^{\text{channels}'[c_3]}
\end{aligned}$$

Model:

$$\begin{aligned}
X^0 &\in \mathbb{R}^{\text{batch} \times \text{channels}[c_1] \times \text{height} \times \text{width}} \\
X^1 &= [\text{ReLU}(\text{conv2d}(X^0; W^1, b^1))]_{\text{channels}' \rightarrow \text{channels}} \\
X^2 &= \text{maxpool2d}(X^1, ph_1, ph_2) \\
X^3 &= [\text{ReLU}(\text{conv2d}(X^2; W^2, b^2))]_{\text{channels}' \rightarrow \text{channels}} \\
X^4 &= \text{maxpool2d}(X^3, ph_2, ph_2) \\
X^5 &= [X^4]_{\text{height, width, channels} \rightarrow \text{layer}} \\
H &= \text{ReLU}(W^3_{\text{layer}} \cdot X^5 + b^3) \\
O &= \text{softmax}(\text{ReLU}(W^4_{\text{hidden}} \cdot H + b^4))_{\text{class}}
\end{aligned}$$

The flattening operation in the equation for  $X^5$  is defined in §??.

### 3.5 Other examples

#### 3.5.1 Continuous bag of words

A continuous bag-of-words model classifies by summing up the embeddings of a sequence of words  $X$  and then projecting them to the space of classes.

$$\begin{aligned}
\text{vocab} &\stackrel{\text{ind}}{=} [v] \\
\text{hidden} &\stackrel{\text{ind}}{=} [h] \\
\text{classes} &\stackrel{\text{ind}}{=} [c] \\
X &\in \{0, 1\}^{\text{seq}[n] \times \text{vocab}} & \sum_{\text{vocab}} X &= 1 \\
E &\in \mathbb{R}^{\text{vocab} \times \text{hidden}} \\
W &\in \mathbb{R}^{\text{classes} \times \text{hidden}} \\
\text{cbow}(X; E, W) &= \text{softmax}(W_{\text{class}} \cdot E_{\text{hidden}} \cdot X_{\text{vocab}})
\end{aligned}$$

Here, the two contractions can be done in either order, so we leave the parentheses off.

#### 3.5.2 Sudoku ILP

Sudoku puzzles can be represented as binary tiled tensors. Given a grid we can check that it is valid by converting it to a grid of grids. Constraints then ensure that there is one digit per row, per column and per sub-box

$$\begin{aligned}
X &\in \{0, 1\}^{\text{height}[9] \times \text{width}[9] \times \text{assign}[9]} \\
\text{check}(X) &= \left( \sum_{\text{assign}} Y = \sum_{\text{Height}, \text{height}} Y = \sum_{\text{Width}, \text{width}} Y = \sum_{\text{height}, \text{width}} Y = 1 \right) \\
Y &\in \{0, 1\}^{\text{Height}[3] \times \text{Width}[3] \times \text{height}[3] \times \text{width}[3] \times \text{assign}[9]} \\
Y_{\text{Height}(r), \text{height}(r'), \text{Width}(c), \text{width}(c')} &= X_{\text{height}(r \times 3 + r' - 1), \text{width}(c \times 3 + c' - 1)}
\end{aligned}$$

### 3.5.3 *K*-means clustering

The following equations define one step of *k*-means clustering. Given a set of points  $X$  and an initial set of cluster centers  $C$ ,

$$\begin{aligned}
\text{batch} &\stackrel{\text{ind}}{=} [b] \\
\text{space} &\stackrel{\text{ind}}{=} [k] \\
\text{clusters} &\stackrel{\text{ind}}{=} [c] \\
X &\in \mathbb{R}^{\text{batch} \times \text{space}} \\
C &\in \mathbb{R}^{\text{clusters} \times \text{space}}
\end{aligned}$$

we compute cluster assignments

$$\begin{aligned}
Q &= \underset{\text{clusters}}{\text{argmin}} \underset{\text{space}}{\text{norm}}(C - X) \\
&= \lim_{\alpha \rightarrow -\infty} \underset{\text{clusters}}{\text{softmax}} \left( \alpha \underset{\text{space}}{\text{norm}}(C - X) \right)
\end{aligned}$$

then we recompute the cluster centers:

$$C \leftarrow \sum_{\text{batch}} \frac{Q \odot X}{Q}.$$

### 3.5.4 Beam search

Beam search is a commonly used approach for approximate discrete search. Here  $H$  is the score of each element in the beam,  $S$  is the state of each element in the beam, and  $f$  is an update function that returns the score of each state transition. Beam step returns the new  $H$  tensor.

$$\begin{aligned}
\text{batch} &\stackrel{\text{ind}}{=} [b] \\
\text{beam} &\stackrel{\text{ind}}{=} [k] \\
\text{state}, \text{state}' &\stackrel{\text{ind}}{=} [s] \\
H &\in \mathbb{R}^{\text{batch} \times \text{beam}} \\
S &\in \{0, 1\}^{\text{batch} \times \text{beam} \times \text{state}} \\
f &: \{0, 1\}^{\text{state}} \rightarrow \mathbb{R}^{\text{state}'} \\
\text{beamstep}(H, S) &= \max_{\text{beam}, \text{state}'} \left( \text{softmax}_{\text{state}'}(f(S)) \odot H \right)
\end{aligned}
\qquad \sum_{\text{state}} S = 1$$

### 3.5.5 Multivariate normal distribution

In our notation, the application of a bilinear form is more verbose than the standard notation  $((X - \mu)^\top \Sigma^{-1} (X - \mu))$ , but also makes it look more like a function of two arguments (and would generalize to three or more arguments).

$$\begin{aligned}
\text{batch} &\stackrel{\text{ind}}{=} [b] \\
\text{d}, \text{d1}, \text{d2} &\stackrel{\text{ind}}{=} [k] \\
X &\in \mathbb{R}^{\text{batch} \times \text{d}} \\
\mu &\in \mathbb{R}^{\text{d}} \\
\Sigma &\in \mathbb{R}^{\text{d1} \times \text{d2}} \\
\mathcal{N}(X; \mu, \Sigma) &= \frac{\exp \left( -\frac{1}{2} \left( \text{inv}_{\text{d1}, \text{d2}}(\Sigma) \cdot_{\text{d1}} [X - \mu]_{\text{d} \rightarrow \text{d1}} \right) \cdot_{\text{d2}} [X - \mu]_{\text{d} \rightarrow \text{d2}} \right)}{\sqrt{(2\pi)^k \det_{\text{d1}, \text{d2}}(\Sigma)}}
\end{aligned}$$

## 4 L<sup>A</sup>T<sub>E</sub>X Macros

Many of the L<sup>A</sup>T<sub>E</sub>X macros used in this document are available in the style file <https://namedtensor.github.io/namedtensor.sty>. To use it, put

`\usepackage{namedtensor}`

in the preamble of your L<sup>A</sup>T<sub>E</sub>X source file (after `\documentclass{article}` but before `\begin{document}`).

The style file contains a small number of macros:

- Use `\name{foo}` to write an index name: `foo`.
- Use `\ndef` to define a named index set: `foo`  $\stackrel{\text{ind}}{=} [n]$ .



- Use  $\backslash\mathrm{ndot}\{\mathrm{foo}\}$   $B$  for contraction:  $A \cdot_{\mathrm{foo}} B$ . Similarly, use  $\backslash\mathrm{ncat}\{\mathrm{foo}\}$   $B$  for concatenation.
- Use  $\backslash\mathrm{nsum}\{\mathrm{foo}\}$   $A$  for summation:  $\sum_{\mathrm{foo}} A$ .
- Use  $\backslash\mathrm{nfun}\{\mathrm{foo}\}\{\mathrm{qux}\}$   $A$  for a function named *qux* with a name under it:  $\mathrm{qux}_{\mathrm{foo}} A$ .
- Use  $\backslash\mathrm{nmov}\{\mathrm{foo}\}\{\mathrm{bar}\}\{A\}$  for renaming:  $[A]_{\mathrm{foo} \rightarrow \mathrm{bar}}$ .

## 5 Formal Definitions

### 5.1 Records

A *named index* is a pair, written  $i(x)$ , where  $i$  is a *name* and  $x$  is usually a positive integer. We write both names and variables ranging over names using sans-serif font.

A *record* is a set of named indices  $\{i_1(x_1), \dots, i_r(x_r)\}$ , where  $i_1, \dots, i_r$  are pairwise distinct names.

If  $i$  is a name and  $X$  is a set, define the *named index set*  $i[X]$  to be the set  $\{i(x) \mid x \in X\}$ . We deal with named index sets of the form  $i[[n]]$  ( $= i[\{1, \dots, n\}]$ ) so frequently that we abbreviate this as  $i[n]$ .

A *named shape* is a set of named index sets  $\{i_1[X_1], \dots, i_r[X_r]\}$ , where  $i_1, \dots, i_r$  are pairwise distinct names. A named shape defines a set of records:

$$\mathrm{ind}\{i_1[X_1], \dots, i_r[X_r]\} = \{\{i_1(x_1), \dots, i_r(x_r)\} \mid x_1 \in X_1, \dots, x_r \in X_r\},$$

and for this reason we often write a shape with the alternative notation  $i_1[X_1] \times \dots \times i_r[X_r]$ , because it's like an “unordered Cartesian product” of the index sets.

We say two shapes  $\mathcal{S}$  and  $\mathcal{T}$  are *compatible* if whenever  $i(X) \in \mathcal{S}$  and  $i(Y) \in \mathcal{T}$ , then  $X = Y$ . We say that  $\mathcal{S}$  and  $\mathcal{T}$  are *orthogonal* if there is no  $i$  such that  $i(X) \in \mathcal{S}$  and  $i(Y) \in \mathcal{T}$  for any  $X, Y$ .

If  $t \in \mathrm{ind}\mathcal{T}$  and  $\mathcal{S} \subseteq \mathcal{T}$ , then we write  $t|_{\mathcal{S}}$  for the unique tuple in  $\mathrm{ind}\mathcal{S}$  such that  $t|_{\mathcal{S}} \subseteq t$ .

Because it's very common for a name to be always used with the same index set, we can write  $i \stackrel{\mathrm{ind}}{=} X$  to create a named index set  $i[X]$ , and this named index set is *also* called  $i$ . The context always makes it clear whether the name or the named index set is meant.

### 5.2 Named tensors

Let  $F$  be a field and let  $\mathcal{T}$  be a set of records. Then a *named tensor over  $F$  with shape  $\mathcal{T}$*  is a mapping from  $\mathcal{T}$  to  $F$ . We write the set of all named tensors

with shape  $\mathcal{T}$  as  $F^{\mathcal{T}}$ .

We don't make any distinction between a scalar (an element of  $F$ ) and a named tensor with empty shape (an element of  $F^{\emptyset}$ ).

If  $A \in F^{\mathcal{T}}$ , then we access an element of  $A$  by applying it to a record  $t \in \text{ind } \mathcal{T}$ ; but we write this using the usual subscript notation:  $A_t$  rather than  $A(t)$ . To avoid clutter, in place of  $A_{\{i_1(x_1), \dots, i_r(x_r)\}}$ , we usually write  $A_{i_1(x_1), \dots, i_r(x_r)}$ . When a named tensor is an expression like  $(A+B)$ , we surround it with square brackets like this:  $[A+B]_{i_1(x_1), \dots, i_r(x_r)}$ .

We also allow partial indices. If  $A$  is a tensor with shape  $\mathcal{T}$  and  $s \in \text{ind } \mathcal{S}$  where  $\mathcal{S} \subseteq \mathcal{T}$ , then we define  $A_s$  to be the named tensor with shape  $\mathcal{T} \setminus \mathcal{S}$  such that, for any  $t \in \text{ind}(\mathcal{T} \setminus \mathcal{S})$ ,

$$[A_s]_t = A_{s \cup t}.$$

(For the edge case  $\mathcal{S} = \mathcal{U}$  and  $\mathcal{T} = \emptyset$ , our definitions for indexing and partial indexing coincide: one gives a scalar and the other gives a tensor with empty shape, but we don't distinguish between the two.)

### 5.3 Extending functions to named tensors

In §2, we described several classes of functions that can be extended to named tensors. Here, we define how to do this for general functions.

Let  $f: F^{\mathcal{S}} \rightarrow G^{\mathcal{T}}$  be a function from tensors to tensors. For any shape  $\mathcal{U}$  orthogonal to both  $\mathcal{S}$  and  $\mathcal{T}$ , we can extend  $f$  to:

$$\begin{aligned} f: F^{\mathcal{S} \cup \mathcal{U}} &\rightarrow G^{\mathcal{T} \cup \mathcal{U}} \\ [f(A)]_u &= f(A_u) \quad \text{for all } u \in \text{ind } \mathcal{U}. \end{aligned}$$

If  $f$  is a multary function, we can extend its arguments to larger shapes, and we don't have to extend all the arguments with the same names. We consider just the case of two arguments; three or more arguments are analogous. Let  $f: F^{\mathcal{S}} \times G^{\mathcal{T}} \rightarrow H^{\mathcal{U}}$  be a binary function from tensors to tensors. For any shapes  $\mathcal{S}'$  and  $\mathcal{T}'$  that are compatible with each other and orthogonal to  $\mathcal{S}$  and  $\mathcal{T}$ , respectively, and  $\mathcal{U}' = \mathcal{S}' \cup \mathcal{T}'$  is orthogonal to  $\mathcal{U}$ , we can extend  $f$  to:

$$\begin{aligned} f: F^{\mathcal{S} \cup \mathcal{S}'} \times G^{\mathcal{T} \cup \mathcal{T}'} &\rightarrow H^{\mathcal{U} \cup \mathcal{U}'} \\ [f(A, B)]_u &= f(A_{u|_{\mathcal{S}'}}, B_{u|_{\mathcal{T}'}}) \quad \text{for all } u \in \text{ind } \mathcal{U}'. \end{aligned}$$

All of the tensor operations described in §2.2 can be defined in this way. For example, the contraction operator extends the following “named dot-product”:

$$\begin{aligned} \cdot_i: F^{i[n]} \times F^{i[n]} &\rightarrow F \\ A \cdot_i B &= \sum_{i=1}^n A_{i(i)} B_{i(i)}. \end{aligned}$$

## 6 Extensions

### 6.1 Index types

We have defined a named index set as a pair  $i[X]$ , where  $i$  is a name and  $X$  is a set, usually  $[n]$  for some  $n$ . In this section, we consider some other possibilities for  $X$ .

#### 6.1.1 Non-integral types

The sets  $X$  don't have to contain integers. For example, if  $V$  is the vocabulary of a natural language ( $V = \{\text{cat}, \text{dog}, \dots\}$ ), we could define a matrix of word embeddings:

$$E \in \mathbb{R}^{\text{vocab}[V] \times \text{emb}[d]}.$$

#### 6.1.2 Integers with units

If  $u$  is a symbol and  $n > 0$ , define  $[n]u = \{1u, 2u, \dots, nu\}$ . You could think of  $u$  as analogous to a physical unit, like kilograms. The elements of  $[n]u$  can be added and subtracted like integers ( $au + bu = (a + b)u$ ) or multiplied by unitless integers ( $c \cdot au = (c \cdot a)u$ ), but numbers with different units are different ( $au \neq av$ ).

Then the set  $[n]u$  could be used as an index set, which would prevent the index from being aligned with another index that uses different units. For example, if we want to define a tensor representing an image, we might write

$$A \in \mathbb{R}^{\text{height}[[h]\text{pixels}] \times \text{width}[[w]\text{pixels}]}.$$

If we have another tensor representing a go board, we might write

$$B \in \mathbb{R}^{\text{height}[[n]\text{points}] \times \text{width}[[n]\text{points}]},$$

and even if it happens that  $h = w = n$ , it would be incorrect to write  $A + B$  because the units do not match.

#### 6.1.3 Tuples of integers

An index set could also be  $[m] \times [n]$ , which would be a way of sneaking ordered indices into named tensors, useful for matrix operations. For example, instead of defining an  $\text{inv}$  operator that takes two subscripts, we could write

$$\begin{aligned} A &\in \mathbb{R}^{d[m \times n]} = \mathbb{R}^{d[[m] \times [n]]} \\ B &= \text{inv}_d A. \end{aligned}$$

We could also define an operator  $\circ$  for matrix-matrix and matrix-vector multiplication:

$$\begin{aligned} c &\in \mathbb{R}^{d[n]} \\ D &= A \underset{d}{\circ} B \underset{d}{\circ} c. \end{aligned}$$

## 6.2 Duality

In applied linear algebra, we distinguish between column and row vectors; in pure linear algebra, vector spaces and dual vector spaces; in tensor algebra, contravariant and covariant indices; in quantum mechanics, bras and kets. Do we need something like this?

In §3.1.2 we saw that defining an RNN requires renaming of indices, because a linear transformation must map one index to another index; if we want to map an index to itself, we need to use renaming.

In this section, we describe three possible solutions to this problem, and welcome comments about which (if any) would be best.

### 6.2.1 Contracting two names

We define a version of the contraction operator that can contract two indices with different names (and the same size):

$$\begin{aligned} \cdot_{ij} : F^{i[n]} \times F^{j[n]} &\rightarrow F \\ A \cdot_{ij} B &= \sum_{k=1}^n A_{i(k)} B_{j(k)}. \end{aligned}$$

For example, the RNN would look like this.

$$\begin{aligned} x^{(t)} &\in \mathbb{R}^{\text{emb}[d]} \\ h^{(t)} &\in \mathbb{R}^{\text{state}[d]} \\ A &\in \mathbb{R}^{\text{state}[d] \times \text{state}'[d]} \\ B &\in \mathbb{R}^{\text{emb}[d] \times \text{state}[d]} \\ c &\in \mathbb{R}^{\text{state}[d]} \\ h^{(t+1)} &= \tanh \left( A \cdot_{\text{state}'|\text{state}} h^{(t)} + B \cdot_{\text{emb}} x^{(t)} + c \right) \end{aligned}$$

### 6.2.2 Starred index names

If  $i$  is a name, we also allow a tensor to have an index  $i^*$  (alternatively: superscript  $i$ ). Multiplication contracts starred indices in the left operand with

non-starred indices in the right operand.

$$\begin{aligned}
x^{(t)} &\in \mathbb{R}^{\text{emb}[d]} \\
h^{(t)} &\in \mathbb{R}^{\text{state}[d]} \\
A &\in \mathbb{R}^{\text{state}*[d] \times \text{state}[d]} \\
B &\in \mathbb{R}^{\text{emb}*[d] \times \text{state}[d]} \\
c &\in \mathbb{R}^{\text{state}[d]} \\
h^{(t+1)} &= \tanh \left( A \underset{\text{state}}{\cdot} h^{(t)} + B \underset{\text{emb}}{\cdot} x^{(t)} + c \right)
\end{aligned}$$

The contraction operator can be defined as:

$$\begin{aligned}
\underset{i}{\cdot} : F^{i*[n]} \times F^{i[n]} &\rightarrow F \\
A \underset{i}{\cdot} B &= \sum_{i=1}^n A_{i*(i)} B_{i(i)}.
\end{aligned}$$

There are a few variants of this idea that have been floated:

1.  $\cdot$  (no subscript) contracts every starred index in its left operand with every corresponding unstarred index in its right operand. Rejected.
2.  $\underset{i}{\cdot}$  contracts  $i$  with  $i$ , and we need another notation like  $\underset{i(*)}{\cdot}$  or  $\underset{i}{\times}$  for contracting  $i*$  with  $i$ .
3.  $\underset{i}{\cdot}$  always contracts  $i*$  with  $i$ ; there's no way to contract  $i$  with  $i$ .

### 6.2.3 Named and numbered indices

We allow indices to have names that are natural numbers  $1, 2, \dots$ , and we define “numbering” and “naming” operators:

$$\begin{aligned}
A_i &\quad \text{rename index } i \text{ to } 1 \\
A_{i,j} &\quad \text{rename index } i \text{ to } 1 \text{ and } j \text{ to } 2 \\
A_{\rightarrow i} &\quad \text{rename index } 1 \text{ to } i \\
A_{\rightarrow i,j} &\quad \text{rename index } 1 \text{ to } i \text{ and } 2 \text{ to } j
\end{aligned}$$

The numbering operators are only defined on tensors that have no numbered indices.

Then we adopt the convention that standard vector/matrix operations operate on the numbered indices. For example, vector dot-product always uses index 1 of both its operands, so that we can write

$$C = A_i \cdot B_i$$

equivalent to  $C = A \underset{i}{\cdot} B$ .

Previously, we had to define a new version of every operation; most of the time, it looked similar to the standard version (e.g.,  $\max$  vs  $\max_i$ ), but occasionally it looked quite different (e.g., matrix inversion). With numbered indices, we can use standard notation for everything. (This also suggests a clean way to integrate code that uses named tensors with code that uses ordinary tensors.)

We also get the renaming operation for free:  $A_{i \rightarrow j} = [A_i]_{\rightarrow j}$  renames index  $i$  to  $j$ .

Finally, this notation alleviates the duality problem, as can be seen in the definition of a RNN:

$$\begin{aligned} x^{(t)} &\in \mathbb{R}^{\text{emb}[d]} \\ h^{(t)} &\in \mathbb{R}^{\text{state}[d]} \\ A &\in \mathbb{R}^{\text{state}[d] \times \text{state}'[d]} \\ B &\in \mathbb{R}^{\text{state}[d] \times \text{emb}[d]} \\ c &\in \mathbb{R}^{\text{state}[d]} \\ h_{\text{state}}^{(t+1)} &= \tanh \left( A_{\text{state}, \text{state}'} h_{\text{state}}^{(t)} + B_{\text{state}, \text{emb}} x_{\text{emb}}^{(t)} + c_{\text{state}} \right) \end{aligned}$$

or equivalently,

$$h^{(t+1)} = \tanh \left( A_{\text{state}'} \cdot h_{\text{state}}^{(t)} + B_{\text{emb}} \cdot x_{\text{emb}}^{(t)} + c \right)$$

Attention:

$$\begin{aligned} \text{Att}: \mathbb{R}^{\text{seq}'[n'] \times \text{key}[d_k]} \times \mathbb{R}^{\text{seq}[n] \times \text{key}[d_k]} \times \mathbb{R}^{\text{seq}[n] \times \text{val}[d_v]} &\rightarrow \mathbb{R}^{\text{seq}'[n'] \times \text{val}[d_v]} \\ \text{Att}(Q, K, V) &= \text{softmax} \left[ \frac{Q_{\text{key}} \cdot K_{\text{key}}}{\sqrt{d_k}} \right]_{\text{seq}} \cdot V_{\text{seq}} \end{aligned}$$

Multivariate normal distribution:

$$\begin{aligned} X &\in \mathbb{R}^{\text{batch}[b] \times \text{d}[k]} \\ \mu &\in \mathbb{R}^{\text{d}[k]} \\ \Sigma &\in \mathbb{R}^{\text{d}[k] \times \text{d}'[k]} \\ \mathcal{N}(X; \mu, \Sigma) &= \frac{\exp \left( -\frac{1}{2} [X - \mu]_{\text{d}}^{\top} \Sigma_{\text{d}, \text{d}'}^{-1} [X - \mu]_{\text{d}} \right)}{\sqrt{(2\pi)^k \det \Sigma_{\text{d}, \text{d}'}}} \end{aligned}$$

Because this notation can be a little more verbose (often requiring you to write index names twice), we'd keep around the notation  $A \cdot B$  as a shorthand for  $A_i \cdot B_i$ . We'd also keep named reductions, or at least  $\text{softmax}_i$ .

### 6.3 Indexing with a tensor of indices

Contributors: Tongfei Chen and Chu-Cheng Lin

NumPy defines two kinds of *advanced* (also known as *fancy*) indexing: by integer arrays and by Boolean arrays. Here, we generalize indexing by integer arrays to named tensors. That is, if  $A$  is a named tensor with  $D$  indices and  $I^1, \dots, I^D$  are named tensors, called “indexers,” what is  $A_{I^1, \dots, I^D}$ ?

We first consider the case where all the indexers have the same shape  $\mathcal{S}$ :

$$\begin{aligned} A &\in F^{i_1[X_1], \dots, i_D[X_D]} \\ I^d &\in X_D^{\mathcal{S}} \quad d = 1, \dots, D. \end{aligned}$$

Then  $A_{I^1, \dots, I^D}$  is the named tensor with shape  $\mathcal{S}$  such that for any  $s \in \text{ind } \mathcal{S}$ ,

$$[A_{I^1, \dots, I^D}]_s = A_{I_s^1, \dots, I_s^D}.$$

More generally, suppose the indexers have different but compatible shapes:

$$\begin{aligned} A &\in F^{i_1[X_1], \dots, i_D[X_D]} \\ I^d &\in X_D^{\mathcal{S}_d} \quad d = 1, \dots, D, \end{aligned}$$

where the  $\mathcal{S}_d$  are pairwise compatible. Then  $A_{I^1, \dots, I^D}$  is the named tensor with shape  $\mathcal{S} = \bigcup_d \mathcal{S}_d$  such that for any  $s \in \text{ind } \mathcal{S}$ ,

$$[A_{I^1, \dots, I^D}]_s = A_{I_{s|_{\mathcal{S}_1}}^1, \dots, I_{s|_{\mathcal{S}_D}}^D}.$$

Let’s consider a concrete example in natural language processing. Consider a batch of sentences encoded as a sequence of word vectors, that is, a tensor  $X \in \mathbb{R}^{\text{batch}[B], \text{sent}[N], \text{emb}[E]}$ . For each sentence, we would like to take out the encodings of a particular span for each sentence  $b \in [B]$  in the batch, resulting in a tensor  $Y \in \mathbb{R}^{\text{batch}[B], \text{span}[M], \text{emb}[E]}$ .

We create an indexer for the `sent` axis:  $I_{\text{sent}} \in [N]^{\text{batch}:B \times \text{span}:M}$  that selects the desired tokens. Then we can write

$$Y = X_{\text{batch}[I], \text{sent}[I], \text{emb}[I]}$$

where the other two indexers

$$\begin{aligned} I_{\text{batch}} &\in [B]^{\text{batch}[B]} \\ [I_{\text{batch}}]_{\text{batch}[b]} &= b \\ I_{\text{emb}} &\in [E]^{\text{emb}[E]} \\ [I_{\text{sent}}]_{\text{sent}[n]} &= n \end{aligned}$$

select all values of their respective indices.

## Acknowledgements

Thanks to Ekin Akyürek, Colin McDonald, Chung-chieh Shan, and Nishant Sinha for their input to this document (or the ideas in it).

## References

- Tongfei Chen. 2017. Typesafe abstractions for tensor operations. In *Proceedings of the 8th ACM SIGPLAN International Symposium on Scala*, SCALA 2017, pages 45–50.
- Charles R. Harris, K. Jarrod Millman, St’efan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fern’andez del R’io, Mark Wiebe, Pearu Peterson, Pierre G’erard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy. *Nature*, 585(7825):357–362.
- Dougal Maclaurin, Alexey Radul, Matthew J. Johnson, and Dimitrios Vytiniotis. 2019. Dex: array programming with typed indices. In *NeurIPS Workshop on Program Transformations for ML*.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Alexander Rush. 2019. Named tensors. Open-source software.
- Nishant Sinha. 2018. Tensor shape (annotation) library. Open-source software.
- Torch Contributors. 2019. Named tensors. PyTorch documentation.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, pages 5998–6008. Curran Associates, Inc.