

Named Tensor Notation

David Chiang Sasha Rush Boaz Barak
University of Notre Dame Cornell University Harvard University

Version 0.2

Contents

1	Introduction	2
2	Informal Overview	3
2.1	Named tensors	3
2.2	Named tensor operations	4
3	Examples	8
3.1	Building blocks	8
3.2	Transformer	12
3.3	LeNet	14
3.4	Other examples	14
4	LTeX Macros	17
5	Formal Definitions	18
5.1	Records and shapes	18
5.2	Named tensors	18
5.3	Extending functions to named tensors	19
6	Extensions	20
6.1	Index types	20
6.2	Duality	21
6.3	Indexing with a tensor of indices	25

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 axes by *position*, making code bug-prone. Quite a few libraries have been developed to identify axes 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).

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 axes by name instead of by position. This document describes our proposal to do both.

As a preview, the above equation becomes

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

making it unambiguous which axis 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 axes.

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 *shape*, which looks like this:

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

which stands for records where `foo`'s index ranges from 1 to 2, and `bar`'s index ranges from 1 to 3. The pairs `foo[2]` and `bar[3]` are called *axes*, and again their order 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 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 $\mathbb{R}^{m \times n}$, we write $\mathbb{R}^{\text{foo}[2] \times \text{bar}[3]}$ for the set of all tensors with shape $\text{foo}[2] \times \text{bar}[3]$.

In many contexts, an axis name is used with only one size. In this case, you can simply write `foo` for the unique axis with name `foo`, as in $\mathbb{R}^{\text{foo} \times \text{bar}}$. It's common to leave the size of an axis unspecified at first, and specify its size later (like in a section or appendix on experimental details). To do this, write $|\text{foo}| = 2$, $|\text{bar}| = 3$.

What are good choices for axis 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 axes, 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

$$A + B = \text{foo} \begin{array}{c} \text{bar} \\ \left[\begin{array}{ccc} 3 + 2 & 1 + 7 & 4 + 1 \\ 1 + 8 & 5 + 2 & 9 + 8 \end{array} \right] \end{array}.$$

But things get more complicated when A and B don't have the same shape. If A and B each have an axis with the same name (and size), the two axes are *aligned*, as above. But if A has an axis named `foo` 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 axis `foo`.

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

Similarly, if B has an axis named `foo` and A doesn't, then we effectively replace A with a new tensor A' that contains a copy of A for every value of axis `foo`. 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 axis 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,

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

More reductions: If A has an axis $\text{foo}[n]$, then

$$\begin{aligned}\sum_{\text{foo}} A &= \sum_{i=1}^n A_{\text{foo}(i)} = \begin{bmatrix} 4 & 6 & 13 \end{bmatrix}^{\text{bar}} \\ \text{norm}_{\text{foo}} A &= \sqrt{\sum_{\text{foo}} A^2} = \begin{bmatrix} \sqrt{10} & \sqrt{26} & \sqrt{97} \end{bmatrix}^{\text{bar}} \\ \min_{\text{foo}} A &= \min_i A_{\text{foo}(i)} = \begin{bmatrix} 1 & 1 & 4 \end{bmatrix}^{\text{bar}} \\ \max_{\text{foo}} A &= \max_i A_{\text{foo}(i)} = \begin{bmatrix} 3 & 5 & 9 \end{bmatrix}^{\text{bar}} \\ \text{mean}_{\text{foo}} A &= \frac{1}{n} A = \begin{bmatrix} 2 & 3 & 6.5 \end{bmatrix}^{\text{bar}} \\ \text{var}_{\text{foo}} A &= \frac{1}{n} \sum_{\text{foo}} (A - \text{mean}_{\text{foo}} A)^2 = \begin{bmatrix} 1 & 4 & 6.25 \end{bmatrix}^{\text{bar}}\end{aligned}$$

(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 axes 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 axis. It can be used, for example, for matrix multiplication:

$$\begin{aligned}C &= \text{bar} \begin{bmatrix} 1 & -1 \\ 2 & -2 \\ 3 & -3 \end{bmatrix}^{\text{baz}} \\ A \cdot_{\text{bar}} C &= \text{foo} \begin{bmatrix} 17 & -17 \\ 53 & -53 \end{bmatrix}^{\text{baz}}\end{aligned}$$

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

$$\begin{aligned}A &\in \mathbb{R}^{\text{foo}[m]} \\ B &\in \mathbb{R}^{\text{foo}[m] \times \text{bar}[n]} \\ C &\in \mathbb{R}^{\text{foo}[m] \times \text{bar}[n]}\end{aligned}$$

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:

$$\underset{\text{foo}}{\text{softmax}} A = \frac{\exp A}{\sum_{\text{foo}} \exp A} \approx \underset{\text{foo}}{\text{foo}} \overset{\text{bar}}{\begin{bmatrix} 0.881 & 0.018 & 0.007 \\ 0.119 & 0.982 & 0.993 \end{bmatrix}}.$$

Also its “hard” version, which maps a vector to a one-hot vector where the maximum element becomes 1 and all other elements become 0 (breaking ties arbitrarily):

$$\underset{\text{foo}}{\text{argmax}} A = \underset{\text{foo}}{\text{foo}} \overset{\text{bar}}{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix}}.$$

Concatenation combines two vectors into one:

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

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

2.2.5 Renaming and reshaping

It's also very handy to have a function that renames an axis:

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

We can reshape two or more axes into one axis:

$$[A]_{(\text{foo}, \text{bar}) \rightarrow \text{baz}} = \overset{\text{baz}}{\begin{bmatrix} 3 & 1 & 4 & 1 & 5 & 9 \end{bmatrix}}$$

Similarly, we can reshape one axis into two or more axes, or even multiple axes into multiple axes. The order of elements in the new axis or axes is undefined. If you need a particular order, you can write a more specific definition.

2.2.6 Matrices

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

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

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

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

A feedforward neural network looks like this:

$$\begin{aligned}
 X^0 &\in \mathbb{R}^{\text{input}} \\
 X^1 &= \sigma(W^1_{\text{input}} \cdot X^0 + b^1) & W^1 &\in \mathbb{R}^{\text{hidden1} \times \text{input}} & b^1 &\in \mathbb{R}^{\text{hidden1}} \\
 X^2 &= \sigma(W^2_{\text{hidden1}} \cdot X^1 + b^2) & W^2 &\in \mathbb{R}^{\text{hidden2} \times \text{hidden1}} & b^2 &\in \mathbb{R}^{\text{hidden2}} \\
 X^3 &= \sigma(W^3_{\text{hidden2}} \cdot X^2 + b^3) & W^3 &\in \mathbb{R}^{\text{output} \times \text{hidden2}} & b^3 &\in \mathbb{R}^{\text{output}}
 \end{aligned}$$

The layer sizes can be set by writing $|\text{input}| = 100$, etc. Alternatively, we could have called the axes $\text{layer}[n_0]$, $\text{layer}[n_1]$, etc. and set the n_l .

If you don't like repeating the equations for fully-connected layers, you can put them inside a function:

$$\text{FullConn}^l(x) = \left[\sigma \left(W_{\text{layer}}^l \cdot x + b^l \right) \right]_{\text{layer}' \rightarrow \text{layer}}$$

where

$$\begin{aligned} W^l &\in \mathbb{R}^{\text{layer}[n_l] \times \text{layer}[n_{l-1}]} \\ b^l &\in \mathbb{R}^{\text{layer}[n_l]}. \end{aligned}$$

Now FullConn^l encapsulates both the equation for layer l as well as its parameters (analogous to what TensorFlow and PyTorch call *modules*).

Then the network can be defined like this:

$$\begin{aligned} X^0 &\in \mathbb{R}^{\text{layer}[n_0]} \\ X^1 &= \text{FullConn}^1(X^0) \\ X^2 &= \text{FullConn}^2(X^1) \\ X^3 &= \text{FullConn}^3(X^2). \end{aligned}$$

3.1.2 Recurrent neural networks

As a second example, let's define a simple (Elman) RNN. This is similar to the feedforward network, except that the number of timesteps is variable and they all share parameters.

$$\begin{aligned} x^t &\in \mathbb{R}^{\text{input}} & t &= 1, \dots, n \\ W^h &\in \mathbb{R}^{\text{hidden} \times \text{hidden}'} & |\text{hidden}| &= |\text{hidden}'| \\ W^i &\in \mathbb{R}^{\text{input} \times \text{hidden}'} \\ b &\in \mathbb{R}^{\text{hidden}'} \\ h^0 &\in \mathbb{R}^{\text{hidden}} \\ h^t &= \left[\sigma \left(W_{\text{hidden}}^h \cdot h^{t-1} + W_{\text{input}}^i \cdot x^t + b \right) \right]_{\text{hidden}' \rightarrow \text{hidden}} & t &= 1, \dots, n \end{aligned}$$

Here the axis 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 (W^h) to one axis and put the result in the same axis. For possible solutions, see §6.2.

3.1.3 Attention

In the introduction (§1), we mentioned some difficulties in interpreting the equation for attention as it's usually written. In our notation, it looks like this:

$$\text{Att}: \mathbb{R}^{\text{key}} \times \mathbb{R}^{\text{seq} \times \text{key}} \times \mathbb{R}^{\text{seq} \times \text{val}} \rightarrow \mathbb{R}^{\text{val}}$$

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

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

$$\text{Att}: \mathbb{R}^{\text{key}} \times \mathbb{R}^{\text{seq} \times \text{key}} \times \mathbb{R}^{\text{seq} \times \text{val}} \times \mathbb{R}^{\text{seq}} \rightarrow \mathbb{R}^{\text{val}}$$

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

Models often use attention to compute a sequence of values, not just a single value. If Q has (say) a seq' axis, then the above definition computes a sequence of values along the seq' axis. If Q , K , and V have a **head** axis for multiple attention heads, then it will compute multi-head attention.

3.1.4 Convolution

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

$$\text{conv1d}: \mathbb{R}^{\text{channels} \times \text{seq}[n]} \rightarrow \mathbb{R}^{\text{seq}[n']}$$

$$\text{conv1d}(X; W, b) = W \underset{\text{channels, kernel}}{\cdot} U + b$$

where

$$n' = n - |\text{kernel}| + 1$$

$$W \in \mathbb{R}^{\text{channels} \times \text{kernel}}$$

$$U \in \mathbb{R}^{\text{channels} \times \text{seq}[n'] \times \text{kernel}}$$

$$U_{\text{seq}(i), \text{kernel}(j)} = X_{\text{seq}(i+j-1)}$$

$$b \in \mathbb{R}.$$

A 2-dimensional convolution:

$$\text{conv2d}: \mathbb{R}^{\text{channels} \times \text{height}[h] \times \text{width}[w]} \rightarrow \mathbb{R}^{\text{height}[h'] \times \text{width}[w']}$$

$$\text{conv2d}(X; W, b) = W \underset{\text{channels, kh, kw}}{\cdot} U + b$$

where

$$\begin{aligned}
h' &= h - |\mathbf{k}h| + 1 \\
w' &= w - |\mathbf{k}w| + 1 \\
W &\in \mathbb{R}^{\text{channels} \times \mathbf{k}h \times \mathbf{k}w} \\
U &\in \mathbb{R}^{\text{channels} \times \text{height}[h'] \times \text{width}[w'] \times \mathbf{k}h \times \mathbf{k}w} \\
U_{\text{height}(i), \text{width}(j), \mathbf{k}h(ki), \mathbf{k}w(kj)} &= X_{\text{height}(i+\mathbf{k}i-1), \text{width}(j+\mathbf{k}j-1)} \\
b &\in \mathbb{R}.
\end{aligned}$$

3.1.5 Max pooling

$$\begin{aligned}
\text{maxpool1d}_k &: \mathbb{R}^{\text{seq}[n]} \rightarrow \mathbb{R}^{\text{seq}[n/k]} \\
\text{maxpool1d}_k(X) &= \max_k U
\end{aligned}$$

where

$$\begin{aligned}
U &\in \mathbb{R}^{\text{seq}[n/k] \times \mathbf{k}[k]} \\
U_{\text{seq}(i), \mathbf{k}(di)} &= X_{\text{seq}(i \times k + di - 1)}.
\end{aligned}$$

$$\begin{aligned}
\text{maxpool2d}_{kh, kw} &: \mathbb{R}^{\text{height}[h] \times \text{width}[w]} \rightarrow \mathbb{R}^{\text{height}[h/kh] \times \text{width}[w/kw]} \\
\text{maxpool2d}_{kh, kw}(X) &= \max_{\mathbf{k}h, \mathbf{k}w} U
\end{aligned}$$

where

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

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:

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

where

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

Now, suppose that the input has three axes:

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

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

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

3.2 Transformer

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

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

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

$$\begin{aligned}X^0 &= W + P \\ T^1 &= \text{LayerNorm}^1(\text{SelfAtt}^1(X^0)) + X^0 \\ X^1 &= \text{LayerNorm}^{1'}(\text{FFN}^1(T^1)) + T^1 \\ &\vdots \\ T^L &= \text{LayerNorm}^L(\text{SelfAtt}^L(X^{L-1})) + X^{L-1} \\ X^L &= \text{LayerNorm}^{L'}(\text{FFN}^L(T^L)) + T^L \\ O &= \text{softmax}(E \underset{\text{vocab}}{\cdot} \underset{\text{layer}}{X^L})\end{aligned}$$

where LayerNorm, SelfAtt and FFN are defined below.

Layer normalization ($l = 1, 1', \dots, L, L'$):

$$\begin{aligned} \text{LayerNorm}^l: \mathbb{R}^{\text{layer}} &\rightarrow \mathbb{R}^{\text{layer}} \\ \text{LayerNorm}^l(X) &= \text{xnorm}_{\text{layer}}(X; \beta^l, \gamma^l). \end{aligned}$$

We defined attention in §3.1.3; the Transformer uses multi-head self-attention, in which queries, keys, and values are all computed from the same sequence.

$$\begin{aligned} \text{SelfAtt}^l: \mathbb{R}^{\text{seq} \times \text{layer}} &\rightarrow \mathbb{R}^{\text{seq} \times \text{layer}} \\ \text{SelfAtt}^l(X) &= Y \end{aligned}$$

where

$$\begin{aligned} |\text{seq}| &= |\text{seq}'| \\ |\text{key}| = |\text{val}| &= |\text{layer}| / |\text{heads}| \\ Q &= \left[W^{l,Q} \cdot_{\text{layer}} X \right]_{\text{seq} \rightarrow \text{seq}'} & W^{l,Q} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\ K &= W^{l,K} \cdot_{\text{layer}} X & W^{l,K} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\ V &= W^{l,V} \cdot_{\text{layer}} X & W^{l,V} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{val}} \\ M &\in \mathbb{R}^{\text{seq} \times \text{seq}'} \\ M_{\text{seq}(i), \text{seq}'(j)} &= \begin{cases} 0 & i \leq j \\ -\infty & \text{otherwise} \end{cases} \\ Y &= \sum_{\text{heads}} W^{l,O} \cdot_{\text{val}} [\text{Att}(Q, K, V, M)]_{\text{seq}' \rightarrow \text{seq}} & W^{l,O} &\in \mathbb{R}^{\text{heads} \times \text{val} \times \text{layer}} \end{aligned}$$

Feedforward neural networks:

$$\begin{aligned} \text{FFN}^l: \mathbb{R}^{\text{layer}} &\rightarrow \mathbb{R}^{\text{layer}} \\ \text{FFN}^l(X) &= X^2 \end{aligned}$$

where

$$\begin{aligned} X^1 &= \text{relu}(W^{l,1} \cdot_{\text{layer}} X + b^{l,1}) & W^{l,1} &\in \mathbb{R}^{\text{hidden} \times \text{layer}} & b^{l,1} &\in \mathbb{R}^{\text{hidden}} \\ X^2 &= \text{relu}(W^{l,2} \cdot_{\text{hidden}} X^1 + b^{l,2}) & W^{l,2} &\in \mathbb{R}^{\text{layer} \times \text{hidden}} & b^{l,2} &\in \mathbb{R}^{\text{hidden}}. \end{aligned}$$

3.3 LeNet

$$X^0 \in \mathbb{R}^{\text{batch} \times \text{channels}[c_0] \times \text{height} \times \text{width}}$$

$$T^1 = \text{relu}(\text{conv}^1(X^0))$$

$$X^1 = \text{maxpool}^1(T^1)$$

$$T^2 = \text{relu}(\text{conv}^2(X^1))$$

$$X^2 = [\text{maxpool}^2(T^2)]_{(\text{height}, \text{width}, \text{channels}) \rightarrow \text{layer}}$$

$$X^3 = \text{relu}(W^3 \underset{\text{layer}}{\cdot} X^2 + b^3) \quad W^3 \in \mathbb{R}^{\text{hidden} \times \text{layer}} \quad b^3 \in \mathbb{R}^{\text{hidden}}$$

$$O = \text{softmax}_{\text{classes}}(W^4 \underset{\text{hidden}}{\cdot} X^3 + b^4) \quad W^4 \in \mathbb{R}^{\text{classes} \times \text{hidden}} \quad b^4 \in \mathbb{R}^{\text{classes}}$$

The flattening operation in the equation for X^2 is defined in §2.2.5. Alternatively, we could have written

$$X^2 = \text{maxpool}^2(T^2)$$

$$X^3 = \text{relu}(W^3 \underset{\text{height}, \text{width}, \text{channels}}{\cdot} X^2 + b^3) \quad W^3 \in \mathbb{R}^{\text{hidden} \times \text{height} \times \text{width} \times \text{channels}}.$$

The convolution and pooling operations are defined as follows:

$$\text{conv}^l(X) = [\text{conv2d}(X; W^l, b^l)]_{\text{channels}' \rightarrow \text{channels}}$$

where

$$|\text{channels}| = |\text{channels}'|$$

$$W^l \in \mathbb{R}^{\text{channels}'[c_l] \times \text{channels}[c_{l-1}] \times \text{kh}[kh_l] \times \text{kw}[kw_l]}$$

$$b^l \in \mathbb{R}^{\text{channels}'[c_l]}$$

and

$$\text{maxpool}^l(X) = \text{maxpool2d}_{ph^l, ph^l}(X).$$

3.4 Other examples

3.4.1 Discrete random variables

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

$$\begin{aligned} p(B \mid A) &\in [0, 1]^{\mathbf{A} \times \mathbf{B}} & \sum_{\mathbf{B}} p(B \mid A) &= 1 \\ p(A) &\in [0, 1]^{\mathbf{A}} & \sum_{\mathbf{A}} p(A) &= 1 \end{aligned}$$

Then many common operations on probability distributions can be expressed in terms of tensor operations:

$$\begin{aligned}
p(A, B) &= p(B \mid A) \odot p(A) && \text{chain rule} \\
p(B) &= \sum_A p(A, B) = p(B \mid A) \cdot_A p(A) && \text{marginalization} \\
p(A \mid B) &= \frac{p(A, B)}{p(B)} = \frac{p(B \mid A) \odot p(A)}{p(B \mid A) \cdot_A p(A)}. && \text{Bayes' rule}
\end{aligned}$$

3.4.2 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.

$$\begin{aligned}
\text{cbow}: \{0, 1\}^{\text{seq} \times \text{vocab}} &\rightarrow \mathbb{R}^{\text{seq} \times \text{classes}} \\
\text{cbow}(X; E, W) &= \text{softmax}_{\text{class}}(W_{\text{class}} \cdot_{\text{hidden}} E_{\text{hidden}} \cdot_{\text{vocab}} X)
\end{aligned}$$

where

$$\begin{aligned}
\sum_{\text{vocab}} X &= 1 \\
E &\in \mathbb{R}^{\text{vocab} \times \text{hidden}} \\
W &\in \mathbb{R}^{\text{classes} \times \text{hidden}}.
\end{aligned}$$

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

3.4.3 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.

$$\begin{aligned}
\text{check}: \{0, 1\}^{\text{height}[9] \times \text{width}[9] \times \text{assign}[9]} &\rightarrow \{0, 1\} \\
\text{check}(X) &= \mathbb{I} \left[\begin{array}{l} \sum_{\text{assign}} X = 1 \wedge \sum_{\text{height}, \text{width}} Y = 1 \wedge \\ \sum_{\text{height}} X = 1 \wedge \sum_{\text{width}} X = 1 \end{array} \right]
\end{aligned}$$

where

$$\begin{aligned}
Y &\in \{0, 1\}^{\text{height}'[3] \times \text{width}'[3] \times \text{height}[3] \times \text{width}[3] \times \text{assign}[9]} \\
Y_{\text{height}'(h'), \text{height}(h), \text{width}'(w'), \text{width}(w)} &= X_{\text{height}(3h' + h - 1), \text{width}(3w' + w - 1)}.
\end{aligned}$$

3.4.4 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 ,

$$\begin{aligned} X &\in \mathbb{R}^{\text{batch} \times \text{space}} \\ C &\in \mathbb{R}^{\text{clusters} \times \text{space}} \end{aligned}$$

we repeat the following update: Compute cluster assignments

$$Q = \underset{\text{clusters}}{\operatorname{argmin}} \underset{\text{space}}{\operatorname{norm}}(C - X)$$

then recompute the cluster centers:

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

3.4.5 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.

$$\begin{aligned} H &\in \mathbb{R}^{\text{beam}} \\ S &\in \{0, 1\}^{\text{beam} \times \text{state}} \\ f &: \{0, 1\}^{\text{state}} \rightarrow \mathbb{R}^{\text{state}} \end{aligned} \quad \sum_{\text{state}} S = 1$$

Then we repeat the following update:

$$\begin{aligned} H' &= \underset{\text{beam}}{\operatorname{max}}(H \odot f(S)) \\ H &\leftarrow \underset{\text{state, beam}}{\operatorname{maxk}} H' \\ S &\leftarrow \underset{\text{state, beam}}{\operatorname{argmaxk}} H' \end{aligned}$$

where

$$\begin{aligned} \underset{\text{ax}, k}{\operatorname{maxk}}: \mathbb{R}^{\text{ax}} &\rightarrow \mathbb{R}^k \\ \underset{\text{ax}, k}{\operatorname{argmaxk}}: \mathbb{R}^{\text{ax}} &\rightarrow \{0, 1\}^{\text{ax}, k} \end{aligned}$$

are defined such that $[\underset{\text{ax}, k}{\operatorname{maxk}} A]_{k(i)}$ is the i -th largest value along axis ax and $A \cdot (\underset{\text{ax}}{\operatorname{argmaxk}} A) = \underset{\text{ax}, k}{\operatorname{max}} A$.

We can add a batch axis to H and S and the above equations will work unchanged.

3.4.6 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).

$$\mathcal{N}: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\mathcal{N}(X; \mu, \Sigma) = \frac{\exp\left(-\frac{1}{2} \left(\underset{d1, d2}{\text{inv}}(\Sigma) \underset{d1}{\cdot} [X - \mu]_{d \rightarrow d1} \right) \underset{d2}{\cdot} [X - \mu]_{d \rightarrow d2} \right)}{\sqrt{(2\pi)^{|d|} \underset{d1, d2}{\det}(\Sigma)}}$$

where

$$\begin{aligned} |d| &= |d1| = |d2| \\ \mu &\in \mathbb{R}^d \\ \Sigma &\in \mathbb{R}^{d1 \times d2}. \end{aligned}$$

4 L^AT_EX Macros

Many of the L^AT_EX 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^AT_EX source file (after `\documentclass{article}` but before `\begin{document}`).

The style file contains a small number of macros:

- Basics
 - Use `\name{foo}` to write an axis name: `foo`.
 - Use `\mathbb{R}^{\set{foo}{2}}` to write a set of tensors: $\mathbb{R}^{\text{foo}[2]}$.
 - Use `A_{\nid{foo}{1}}` to index a tensor: $A_{\text{foo}(1)}$.
 - Use `\nmov{foo}{bar}{A}` for renaming: $[A]_{\text{foo} \rightarrow \text{bar}}$.
- Binary operators
 - Use `A \ndot{foo} B` for contraction: $A \underset{\text{foo}}{\cdot} B$.
 - Use `A \ncat{foo} B` for concatenation: $A \underset{\text{foo}}{\oplus} B$.
 - In general, you can use `\nbin` to make a new binary operator with a name under it: `A \nbin{foo}{\star} B` gives you $A \underset{\text{foo}}{\star} B$.

- Functions

- Use `\nsum{foo} A` for summation: $\sum_{\text{foo}} A$.
- In general, you can use `\nfun` to make a function with a name under it: `\nfun{foo}{qux} A` gives you $\text{qux}_{\text{foo}} A$.

5 Formal Definitions

5.1 Records and shapes

A *named index* is a pair, written $\text{ax}(i)$, where ax is a *name* and i is usually a natural number. We write both names and variables ranging over names using sans-serif font.

A *record* is a set of named indices $\{\text{ax}_1(i_1), \dots, \text{ax}_r(i_r)\}$, where $\text{ax}_1, \dots, \text{ax}_r$ are pairwise distinct names.

An *axis* is a pair, written $\text{ax}[I]$, where ax is a name and I is a set of *indices*.

We deal with axes of the form $\text{ax}[[n]]$ (that is, $\text{ax}[\{1, \dots, n\}]$) so frequently that we abbreviate this as $\text{ax}[n]$.

In many contexts, there is only one axis with name ax , and so we refer to the axis simply as ax . The context always makes it clear whether ax is a name or an axis. If ax is an axis, we write $\text{ind}(\text{ax})$ for its index set, and we write $|\text{ax}|$ as shorthand for $|\text{ind}(\text{ax})|$.

A *shape* is a set of axes, written $\text{ax}_1[I_1] \times \dots \times \text{ax}_r[I_r]$, where $\text{ax}_1, \dots, \text{ax}_r$ are pairwise distinct names. A shape defines a set of records:

$$\text{rec}(\text{ax}_1[I_1] \times \dots \times \text{ax}_r[I_r]) = \{\{\text{ax}_1(i_1), \dots, \text{ax}_r(i_r)\} \mid i_1 \in I_1, \dots, i_r \in I_r\}.$$

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

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

5.2 Named tensors

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

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 $A \in F^{\mathcal{S}}$, then we access an element of A by applying it to a record $s \in \text{rec } \mathcal{S}$; but we write this using the usual subscript notation: A_s rather than $A(s)$. To avoid clutter, in place of $A_{\{\text{ax}_1(x_1), \dots, \text{ax}_r(x_r)\}}$, we usually write $A_{\text{ax}_1(x_1), \dots, \text{ax}_r(x_r)}$. When a named tensor is an expression like $(A + B)$, we surround it with square brackets like this: $[A + B]_{\text{ax}_1(x_1), \dots, \text{ax}_r(x_r)}$.

We also allow partial indexing. If A is a tensor with shape \mathcal{T} and $s \in \text{rec } \mathcal{S}$ 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 $t \in \text{rec}(\mathcal{T} \setminus \mathcal{S})$,

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

(For the edge case $\mathcal{T} = \emptyset$, 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:

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

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:

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

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”:

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

6 Extensions

6.1 Index types

We have defined an axis as a pair $\mathbf{ax}[I]$, where \mathbf{ax} is a name and I is a set, usually $[n]$ for some n . In this section, we consider some other possibilities for I .

6.1.1 Non-integral types

The sets I don't have to contain integers. For example, if V is the vocabulary of a natural language ($V = \{\mathbf{cat}, \mathbf{dog}, \dots\}$), 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 \mathbf{u} is a symbol and $n > 0$, define $[n]\mathbf{u} = \{1\mathbf{u}, 2\mathbf{u}, \dots, n\mathbf{u}\}$. You could think of \mathbf{u} as analogous to a physical unit, like kilograms. The elements of $[n]\mathbf{u}$ can be added and subtracted like integers ($a\mathbf{u} + b\mathbf{u} = (a + b)\mathbf{u}$) or multiplied by unitless integers ($c \cdot a\mathbf{u} = (c \cdot a)\mathbf{u}$), but numbers with different units are different ($a\mathbf{u} \neq a\mathbf{v}$).

Then the set $[n]\mathbf{u}$ could be used as an index set, which would prevent the axis from being aligned with another axis 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, the RNN would look like this:

$$\begin{aligned} x^t &\in \mathbb{R}^{\text{input}} & t = 1, \dots, n \\ W^h &\in \mathbb{R}^{\text{hidden}[d \times d]} \\ W^i &\in \mathbb{R}^{\text{input} \times \text{hidden}[d]} \\ b &\in \mathbb{R}^{\text{hidden}[d]} \\ h^0 &\in \mathbb{R}^{\text{hidden}[d]} \\ h^t &= \sigma \left(W^h_{\text{hidden}} \circ h^{t-1} + W^i_{\text{input}} \cdot x^t + b \right) & t = 1, \dots, n \end{aligned}$$

where \circ does matrix-matrix and matrix-vector multiplication.

This wouldn't really help with self-attention. The multivariate normal distribution would need a transpose.

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 axes, because a linear transformation must map one axis to another axis; if we want to map an axis 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 axes with different names (and the same size):

$$\begin{aligned} \cdot_{\text{ax1}|\text{ax2}} : F^{\text{ax1}[n]} \times F^{\text{ax2}[n]} &\rightarrow F \\ A \cdot_{\text{ax1}|\text{ax2}} B &= \sum_{i=1}^n A_{\text{ax1}(i)} B_{\text{ax2}(i)}. \end{aligned}$$

For example, the RNN would look like this.

$$\begin{aligned} x^t &\in \mathbb{R}^{\text{input}} & t &= 1, \dots, n \\ W^h &\in \mathbb{R}^{\text{hidden} \times \text{hidden}'} & |\text{hidden}| &= |\text{hidden}'| \\ W^i &\in \mathbb{R}^{\text{input} \times \text{hidden}} \\ b &\in \mathbb{R}^{\text{hidden}} \\ h^0 &\in \mathbb{R}^{\text{hidden}} \\ h^t &= \sigma \left(W^h_{\text{hidden}'|\text{hidden}} \cdot h^{t-1} + W^i_{\text{input}} \cdot x^t + b \right) & t &= 1, \dots, n \end{aligned}$$

In self-attention, it seems unavoidable to have one renaming (corresponding to

the transpose in the original formula):

$$\text{Att}: \mathbb{R}^{\text{key}} \times \mathbb{R}^{\text{seq} \times \text{key}} \times \mathbb{R}^{\text{seq} \times \text{val}} \rightarrow \mathbb{R}^{\text{val}}$$

$$\begin{aligned} \text{Att}(Q, K, V) &= \underset{\text{seq}'}{\text{softmax}} \left(\frac{Q \underset{\text{key}}{\cdot} [K]_{\text{seq} \rightarrow \text{seq}'}}{\sqrt{|\text{key}|}} \right) \underset{\text{seq}'|\text{seq}}{\cdot} V \\ |\text{seq}| &= |\text{seq}'| \\ Q &= W^{l,Q} \underset{\text{layer}}{\cdot} X & W^{l,Q} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\ K &= W^{l,K} \underset{\text{layer}}{\cdot} X & W^{l,K} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\ V &= W^{l,V} \underset{\text{layer}}{\cdot} X & W^{l,V} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{val}} \\ Y &= \sum_{\text{heads}} W^{l,O} \underset{\text{val}}{\cdot} \text{Att}(Q, K, V, M) & W^{l,O} &\in \mathbb{R}^{\text{heads} \times \text{val} \times \text{layer}} \end{aligned}$$

Multivariate normal:

$$\begin{aligned} \mathcal{N}: \mathbb{R}^{\text{d}} &\rightarrow \mathbb{R} \\ \mathcal{N}(X; \mu, \Sigma) &= \frac{\exp \left(-\frac{1}{2} \left(\underset{\text{d1}, \text{d2}}{\text{inv}}(\Sigma) \underset{\text{d1}|\text{d}}{\cdot} (X - \mu) \right) \underset{\text{d2}|\text{d}}{\cdot} (X - \mu) \right)}{\sqrt{(2\pi)^{|\text{d}|} \underset{\text{d1}, \text{d2}}{\det}(\Sigma)}} \end{aligned}$$

where

$$\begin{aligned} |\text{d}| &= |\text{d1}| = |\text{d2}| \\ \mu &\in \mathbb{R}^{\text{d}} \\ \Sigma &\in \mathbb{R}^{\text{d1} \times \text{d2}}. \end{aligned}$$

6.2.2 Starred axis names

If ax is a name, we also allow a tensor to have an axis ax^* with the same index set. (Alternatively, we could write one as a superscript and one as a subscript, as is done with contravariant and covariant indices.) Multiplication contracts starred axes in the left operand with non-starred axes in the right operand. There are a few variants of this idea that have been floated:

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

RNN:

$$\begin{aligned}
x^t &\in \mathbb{R}^{\text{input}} & t = 1, \dots, n \\
W^h &\in \mathbb{R}^{\text{hidden} \times \text{hidden} *} \\
W^i &\in \mathbb{R}^{\text{hidden} \times \text{input} *} \\
b &\in \mathbb{R}^{\text{hidden}} \\
h^0 &\in \mathbb{R}^{\text{hidden}} \\
h^t &= \sigma \left(W^h_{\text{hidden}} \cdot h^{t-1} + W^i_{\text{input}} \cdot x^t + b \right) & t = 1, \dots, n
\end{aligned}$$

Self-attention: Again one renaming seems unavoidable.

$$\text{Att}: \mathbb{R}^{\text{key}} \times \mathbb{R}^{\text{seq} \times \text{key}} \times \mathbb{R}^{\text{seq} \times \text{val}} \rightarrow \mathbb{R}^{\text{val}}$$

$$\begin{aligned}
\text{Att}(Q, K, V) &= \text{softmax}_{\text{seq} *} \left(\frac{Q \cdot [K]_{\text{seq} \rightarrow \text{seq} *}}{\sqrt{|\text{key}|}} \right) \cdot_{\text{seq}} V \\
Q &= W^{l, Q}_{\text{layer}} \cdot X & W^{l, Q} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\
K &= W^{l, K}_{\text{layer}} \cdot X & W^{l, K} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\
V &= W^{l, V}_{\text{layer}} \cdot X & W^{l, V} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{val}} \\
Y &= \sum_{\text{heads}} W^{l, O}_{\text{val}} \cdot \text{Att}(Q, K, V, M) & W^{l, O} &\in \mathbb{R}^{\text{heads} \times \text{val} \times \text{layer}}
\end{aligned}$$

Starred axes don't really help much with the multivariate normal distribution.

6.2.3 Named and numbered axes

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

$$\begin{aligned}
A_{\text{ax}} &\quad \text{rename axis ax to 1} \\
A_{\text{ax1}, \text{ax2}} &\quad \text{rename axis ax1 to 1 and ax2 to 2} \\
A_{\rightarrow \text{ax}} &\quad \text{rename axis 1 to ax} \\
A_{\rightarrow \text{ax1}, \text{ax2}} &\quad \text{rename axis 1 to ax1 and 2 to ax2}
\end{aligned}$$

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

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

$$C = A_{\text{ax}} \cdot B_{\text{ax}}$$

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

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_{ax}), but occasionally it looked quite different (e.g., matrix inversion). With numbered axes, 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: $A_{\text{ax1} \rightarrow \text{ax2}} = [A_{\text{ax1}}]_{\rightarrow \text{ax2}}$ renames axis ax1 to ax2 .

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

$$\begin{aligned}
x^t &\in \mathbb{R}^{\text{input}} & t &= 1, \dots, n \\
W^{\text{h}} &\in \mathbb{R}^{\text{hidden} \times \text{hidden}'} & |\text{hidden}| &= |\text{hidden}'| \\
W^{\text{i}} &\in \mathbb{R}^{\text{input} \times \text{hidden}} \\
b &\in \mathbb{R}^{\text{hidden}} \\
h^0 &\in \mathbb{R}^{\text{hidden}} \\
h^t &= \sigma \left(W_{\text{hidden}'}^{\text{h}} \cdot h_{\text{hidden}}^{t-1} + W_{\text{input}}^{\text{i}} \cdot x_{\text{input}}^t + b_{\text{hidden}} \right) & t &= 1, \dots, n
\end{aligned}$$

Attention requires either a renaming or a transpose:

$$\text{Att}: \mathbb{R}^{\text{key}} \times \mathbb{R}^{\text{seq} \times \text{key}} \times \mathbb{R}^{\text{seq} \times \text{val}} \rightarrow \mathbb{R}^{\text{val}}$$

$$\begin{aligned}
\text{Att}(Q, K, V) &= \text{softmax} \left[\frac{Q_{\text{key}} \cdot [K_{\text{key}}]_{\text{seq} \rightarrow \text{seq}'}}{\sqrt{|\text{key}|}} \right]_{\text{seq}'} \cdot V_{\text{seq}} \\
|\text{seq}| &= |\text{seq}'| \\
Q &= W_{\text{layer}}^{l, Q} \cdot X & W^{l, Q} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\
K &= W_{\text{layer}}^{l, K} \cdot X & W^{l, K} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{key}} \\
V &= W_{\text{layer}}^{l, V} \cdot X & W^{l, V} &\in \mathbb{R}^{\text{heads} \times \text{layer} \times \text{val}} \\
Y &= \sum_{\text{heads}} W_{\text{val}}^{l, O} \cdot \text{Att}(Q, K, V, M) & W^{l, O} &\in \mathbb{R}^{\text{heads} \times \text{val} \times \text{layer}}
\end{aligned}$$

The multivariate normal distribution also needs a renaming or transpose:

$$\begin{aligned}
X &\in \mathbb{R}^{\text{d}} \\
\mu &\in \mathbb{R}^{\text{d}} \\
\Sigma &\in \mathbb{R}^{\text{d1} \times \text{d2}} \\
\mathcal{N}(X; \mu, \Sigma) &= \frac{\exp \left(-\frac{1}{2} [X - \mu]_{\text{d}}^{\top} \Sigma_{\text{d1}, \text{d2}}^{-1} [X - \mu]_{\text{d}} \right)}{\sqrt{(2\pi)^k \det \Sigma_{\text{d1}, \text{d2}}}}
\end{aligned}$$

Because this notation can be a little more verbose (often requiring you to write axis names twice), we'd keep around the notation $A \cdot_{\text{ax}} B$ as a shorthand for $A_{\text{ax}} \cdot B_{\text{ax}}$.

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 ι^1, \dots, ι^D are named tensors, called “indexers,” what is $A_{\iota^1, \dots, \iota^D}$?

Advanced indexing could be derived by taking a function

$$\begin{aligned} \text{index}_{\text{ax}}: F^{\text{ax}[I]} \times I &\rightarrow F \\ \text{index}_{\text{ax}}(A, i) &= A_{\text{ax}(i)} \end{aligned}$$

and extending it to higher-order tensors in its second argument according to the rules in §5.3. But because that's somewhat abstract, we give a more concrete definition below.

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

$$\begin{aligned} A &\in F^{\text{ax}_1[I_1] \times \dots \times \text{ax}_D[I_D]} \\ \iota^d &\in I_d^{\mathcal{S}} \quad d = 1, \dots, D. \end{aligned}$$

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

$$[A_{\iota^1, \dots, \iota^D}]_s = A_{\iota_s^1, \dots, \iota_s^D}.$$

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

$$\begin{aligned} A &\in F^{\text{ax}_1[I_1] \times \dots \times \text{ax}_D[I_D]} \\ \iota^d &\in I_d^{\mathcal{S}_d} \quad d = 1, \dots, D, \end{aligned}$$

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

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

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 $X \in \mathbb{R}^{\text{batch}[B] \times \text{sent}[N] \times \text{emb}[E]}$. 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 $Y \in \mathbb{R}^{\text{batch}[B] \times \text{span}[M] \times \text{emb}[E]}$.

We create an indexer for the `sent` axis: $\iota \in [N]^{\text{batch}[B] \times \text{span}[M]}$ that selects the desired tokens. Also define the function

$$\begin{aligned} \text{arange}_{\text{ax}}(I) &\in I^{\text{ax}[I]} \\ \left[\text{arange}_{\text{ax}}(I) \right]_{\text{ax}(i)} &= i \end{aligned}$$

which generalizes the NumPy function of the same name.

Then we can write

$$Y = X_{\text{batch}(\iota), \underset{\text{sent}}{\text{sent}}(\text{arange}(n)), \underset{\text{emb}}{\text{emb}}(\text{arange}(E))}.$$

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, Stéfan 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 Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-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.