

Named Tensor Notation

David Chiang and Sasha Rush

December 7, 2020

Contents

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+:2017:

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

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 `numpy`, `pytorch` 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` chen2017, `typesafe`, `tsalib` tsalib, `NamedTensor` namedtensor, named tensors in PyTorch `named-tensors`, and `Dex` maclaurin+: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

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

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 §??.

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

A *named tuple* (also known as a *record*) looks like this:

$$\{\text{foo} : 2, \text{bar} : 3\}.$$

The order of the elements doesn’t matter:

$$\{\text{foo} : 2, \text{bar} : 3\} = \{\text{bar} : 3, \text{foo} : 2\}.$$

We use **sans-serif** font for names.

Then, a real *named tensor* is a function from named tuples to reals. Each of its indices has a name, and the ordering of the indices doesn’t matter. For example, here is a tensor with an index named **foo** ranging from 1 to 2 and an index named **bar** ranging from 1 to 3. More succinctly, we say that the *shape* of A is $\{\text{foo} : 2, \text{bar} : 3\}$.

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

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 $\mathbb{R}^{\text{foo}:2,\text{bar}:3}$ for the set of all tensors with shape $\{\text{foo} : 2, \text{bar} : 3\}$.

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

$$A_{\text{foo}:1} = \begin{matrix} & \text{bar} \\ & \begin{bmatrix} 3 & 1 & 4 \end{bmatrix} \end{matrix}$$

$$A_{\text{bar}:3} = \begin{matrix} & \text{foo} \\ \begin{bmatrix} 4 & 9 \end{bmatrix} & \end{matrix}.$$

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 §?? 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:

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

More elementwise unary operations:

kA	scalar multiplication by k
$-A$	negation
$\exp A$	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

$$B = \text{foo} \begin{matrix} & \text{bar} \\ \begin{bmatrix} 2 & 7 & 1 \\ 8 & 2 & 8 \end{bmatrix} \end{matrix}$$

then

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

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 .

$$A + 1 = \text{foo} \begin{array}{c} \text{bar} \\ \begin{bmatrix} 3+1 & 1+1 & 4+1 \\ 1+1 & 5+1 & 9+1 \end{bmatrix} \end{array}$$

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

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

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
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 §??, 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 $\{i : X, \dots\}$, then

$$\sum_{\text{foo}} A = \sum_{x \in X} A_{\text{foo}:x} = \begin{matrix} & \text{bar} \\ [4 & 6 & 13] \end{matrix}$$

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

$$\min_{\text{foo}} A = \min_{x \in X} A_{\text{foo}:x} = \begin{matrix} & \text{bar} \\ [1 & 1 & 4] \end{matrix}$$

$$\max_{\text{foo}} A = \max_{x \in X} A_{\text{foo}:x} = \begin{matrix} & \text{bar} \\ [3 & 5 & 9] \end{matrix}$$

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

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

(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

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

$$B \in \mathbb{R}^{\text{foo}:m, \text{bar}:n}$$

$$C \in \mathbb{R}^{\text{foo}:m, \text{bar}:n}$$

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:

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

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

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

Concatenation combines two vectors into one:

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

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

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 $\{\text{foo} : 2, \text{bar} : 2, \text{baz} : 2\}$:

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

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

3.1 Attention

Let d_k and d_v be positive integers, and let n and n' be the input and output sequence length. Define a function

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

In self-attention, Q , K , and V are all computed from the same sequence. Let d_{model} be a positive integer. The parameters are:

$$\begin{aligned}
 W^Q &\in \mathbb{R}^{\text{emb}:d_{\text{model}}, \text{key}:d_k} \\
 W^K &\in \mathbb{R}^{\text{emb}:d_{\text{model}}, \text{key}:d_k} \\
 W^V &\in \mathbb{R}^{\text{emb}:d_{\text{model}}, \text{val}:d_v} \\
 W^O &\in \mathbb{R}^{\text{val}:d_v, \text{emb}:d_{\text{model}}}
 \end{aligned}$$

Then define

$$\text{SelfAtt}: \mathbb{R}^{\text{seq}:n, \text{emb}:d_{\text{model}}} \rightarrow \mathbb{R}^{\text{seq}:n, \text{emb}:d_{\text{model}}}$$

$$\text{SelfAtt}(X; W^Q, W^K, W^V, W^O) = W^O \underset{\text{val}}{\cdot} [\text{Att}(Q, K, V)]_{\text{seq}' \rightarrow \text{seq}}$$

where

$$Q = W^Q \underset{\text{emb}}{\cdot} [X]_{\text{seq} \rightarrow \text{seq}'}$$

$$K = W^K \underset{\text{emb}}{\cdot} X$$

$$V = W^V \underset{\text{emb}}{\cdot} X.$$

To change this to multi-head self-attention with h attention heads, simply re-define

$$W^Q \in \mathbb{R}^{\text{head}:h, \text{emb}:d_{\text{model}}, \text{key}:d_k}$$

$$W^K \in \mathbb{R}^{\text{head}:h, \text{emb}:d_{\text{model}}, \text{key}:d_k}$$

$$W^V \in \mathbb{R}^{\text{head}:h, \text{emb}:d_{\text{model}}, \text{val}:d_v}$$

$$W^O \in \mathbb{R}^{\text{head}:h, \text{val}:d_v, \text{emb}:d_{\text{model}}}$$

and define

$$\text{MultiSelfAtt}: \mathbb{R}^{\text{seq}:n, \text{emb}:d_{\text{model}}} \rightarrow \mathbb{R}^{\text{seq}:n, \text{emb}:d_{\text{model}}}$$

$$\text{MultiSelfAtt}(X; W^Q, W^K, W^V, W^O) = \sum_{\text{head}} \text{SelfAtt}(X; W^Q, W^K, W^V, W^O).$$

3.2 RNN

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

$$\begin{aligned} x^{(t)} &\in \mathbb{R}^{\text{emb}:d} & t &= 1, \dots, n \\ h^{(t)} &\in \mathbb{R}^{\text{state}:d} & t &= 0, \dots, n \\ A &\in \mathbb{R}^{\text{state}:d, \text{state}':d} \\ B &\in \mathbb{R}^{\text{emb}:d, \text{state}':d} \\ c &\in \mathbb{R}^{\text{state}':d} \end{aligned}$$

$$h^{(t+1)} = \left[\tanh \left(A \underset{\text{state}}{\cdot} h^{(t)} + B \underset{\text{emb}}{\cdot} x^{(t)} + c \right) \right]_{\text{state}' \rightarrow \text{state}}$$

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 §??.

3.3 Fully-Connected Layers

Fully-connected layers are bit more verbose, but make more explicit which parameters connect which layers.

$$\begin{aligned}
 V &\in \mathbb{R}^{\text{output}:o, \text{hidden}:h} & c &\in \mathbb{R}^{\text{output}:o} \\
 W &\in \mathbb{R}^{\text{hidden}:h, \text{in}:i} & b &\in \mathbb{R}^{\text{hidden}:h} \\
 X &\in \mathbb{R}^{\text{batch}:b, \text{in}:i} \\
 Y &= \sigma \left(W_{\text{in}} \cdot X + b \right) \\
 Z &= \sigma \left(V_{\text{hidden}} \cdot Y + c \right)
 \end{aligned}$$

3.4 Deep Learning Norms

These three functions are often informally described using the same equation, but they each correspond to very different functions. They differ by which axes are normalized.

Batch Norm

$$\begin{aligned}
 X &\in \mathbb{R}^{\text{batch}:b, \text{channels}:c, \text{hidden}:h} \\
 \gamma, \beta &\in \mathbb{R}^{\text{batch}:b} \\
 \text{batchnorm}(X; \gamma, \beta) &= \frac{X - \text{mean}_{\text{batch}}(X)}{\sqrt{\text{var}_{\text{batch}}(X) + \epsilon}} \odot \gamma + \beta
 \end{aligned}$$

Instance Norm

$$\begin{aligned}
 X &\in \mathbb{R}^{\text{batch}:b, \text{channels}:c, \text{hidden}:h} \\
 \gamma, \beta &\in \mathbb{R}^{\text{hidden}:h} \\
 \text{instancenorm}(X; \gamma, \beta) &= \frac{X - \text{mean}_{\text{hidden}}(X)}{\sqrt{\text{var}_{\text{hidden}}(X) + \epsilon}} \odot \gamma + \beta
 \end{aligned}$$

Layer Norm

$$\begin{aligned}
 X &\in \mathbb{R}^{\text{batch}:b, \text{channels}:c, \text{hidden}:h} \\
 \gamma, \beta &\in \mathbb{R}^{\text{channels}:c, \text{hidden}:h} \\
 \text{layernorm}(X; \gamma, \beta) &= \frac{X - \text{mean}_{\text{hidden, channels}}(X)}{\sqrt{\text{var}_{\text{hidden, channels}}(X) + \epsilon}} \odot \gamma + \beta
 \end{aligned}$$

3.5 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.

$$\begin{aligned}
 X &\in \{0, 1\}^{\text{seq}:n, \text{vocab}:v} & \sum_{\text{vocab}} X &= 1 \\
 E &\in \mathbb{R}^{\text{vocab}:v, \text{hidden}:h} \\
 W &\in \mathbb{R}^{\text{class}:c, \text{hidden}:h} \\
 \text{cbow}(X; E, W) &= \text{softmax}(W \underset{\text{class}}{\cdot} E \underset{\text{vocab}}{\cdot} X)
 \end{aligned}$$

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

3.6 Bayes' Rule

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:

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

Then Bayes' rule is just:

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

3.7 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.

$$\begin{aligned}
 X &\in \{0, 1\}^{\text{row}:9, \text{col}:9, \text{assign}:9} \\
 \text{check}(X) &= \left(\sum_{\text{assign}} Y = \sum_{\text{Height}, \text{height}} Y = \sum_{\text{Width}, \text{width}} Y = \sum_{\text{height}, \text{width}} Y = 1 \right) \\
 Y &\in \{0, 1\}^{\text{Height}:3, \text{Width}:3, \text{height}:3, \text{width}:3, \text{assign}:9}
 \end{aligned}$$

$$Y_{\text{Height}:r, \text{height}:r', \text{Width}:c, \text{width}:c'} = X_{\text{height}:r \times 3 + r' - 1, \text{width}:c \times 3 + c' - 1},$$

3.8 Max Pooling

Max pooling used in image recognition takes a similar form as the Sudoku example.

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

3.9 1D Convolution

1D Convolution can be easily written by unrolling a tensor and then applying a standard dot product.

$$\begin{aligned}
X &\in \mathbb{R}^{\text{channels:}c, \text{seq:}n} \\
W &\in \mathbb{R}^{\text{out_channels:}c', \text{channels:}c, \text{kw:}k} \\
\text{conv1d}(X, W) &= W \underset{\text{channels, kw}}{\cdot} U \\
U &\in \mathbb{R}^{\text{channels:}c, \text{seq:}n-k+1, \text{kw:}k} \\
U_{\text{seq:}i, \text{kw:}j} &= X_{\text{seq:}i+j-1}
\end{aligned}$$

3.10 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 ,

$$\begin{aligned}
X &\in \mathbb{R}^{\text{batch:}b, \text{space:}k} \\
C &\in \mathbb{R}^{\text{clusters:}c, \text{space:}k}
\end{aligned}$$

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:

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

3.11 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.

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

3.12 Multivariate Normal

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).

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

3.13 Attention with Causal Masking

When the Transformer is used for generation, it is necessary to have an additional mask to ensure the model does not look at future words. This can be included in the attention definition with clear names.

$$\begin{aligned}
Q &\in \mathbb{R}^{\text{key}:d_v, \text{seq}':n} \\
K &\in \mathbb{R}^{\text{head}:h, \text{key}:d_k, \text{seq}:n} \\
V &\in \mathbb{R}^{\text{head}:h, \text{val}:d_v, \text{seq}:n} \\
\text{attention}(Q, K, V) &= \text{softmax}_{\text{seq}} \left(\frac{Q \cdot_{\text{key}} K}{\sqrt{d_k}} + M \right) \cdot_{\text{seq}} V \\
M &\in \mathbb{R}^{\text{seq}:n, \text{seq}':n} \\
M_{\text{seq}:i, \text{seq}':j} &= \begin{cases} 0 & i \leq j \\ -\infty & \text{otherwise} \end{cases}
\end{aligned}$$

3.14 Full Examples: Transformer and LeNet

As further proof of concept, we have written the full models for Transformer (<https://namedtensor.github.io/transformer.html>) and LeNet (<https://namedtensor.github.io/convnet.html>).

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:

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

5 Formal Definitions

5.1 Named tuples

A *named tuple* is a set of pairs, written as $\{i_1 : x_1, \dots, i_r : x_r\}$, where i_1, \dots, i_r are pairwise distinct *names*. We write both names and variables ranging over names using sans-serif font.

If t is a named tuple, we write $\text{dom } t$ for the set $\{i \mid (i : x) \in t \text{ for some } x\}$. If $i \in \text{dom } t$, we write $t.i$ for the unique x such that $(i : x) \in t$. We write the empty named tuple as \emptyset .

We define a partial ordering \sqsubseteq on named tuples: $t_1 \sqsubseteq t_2$ iff for all i, x , $(i : x) \in t_1$ implies $(i : x) \in t_2$. Then $t_1 \sqcap t_2$ is the greatest lower bound of t_1 and t_2 , and $t_1 \sqcup t_2$ is their least upper bound.

A *shape* is a set of pairs written as $\{i_1 : X_1, \dots, i_r : X_r\}$ where X_1, \dots, X_r are sets. We use shapes to define sets of named tuples:

$$\text{ind}\{i_1 : X_1, \dots, i_r : X_r\} = \{(i : 1 = x_1, \dots, i : r = x_r) \mid x_1 \in X_1, \dots, x_r \in X_r\}.$$

We define \sqsubseteq , \sqcup , and \sqcap on shapes just as with named tuples.

If $t \in \text{ind } \mathcal{T}$ and $\mathcal{S} \sqsubseteq \mathcal{T}$, then we write $t|_{\mathcal{S}}$ for the named tuple $\{(i : x) \in t \mid i \in \text{dom } \mathcal{S}\}$.

5.2 Named tensors

Let $[n] = \{1, \dots, n\}$. We deal with shapes of the form $\{i : 1 : [n_1], \dots, i : r : [n_r]\}$ so frequently that we define the shorthand $\{i : 1 : n_1, \dots, i : r : n_r\}$.

Let F be a field and let \mathcal{T} be a shape. Then a *named tensor over F with shape \mathcal{T}* is a mapping from $\text{ind } \mathcal{T}$ to F . We write the set of all named tensors with shape \mathcal{T} as $F^{\mathcal{T}}$. To avoid clutter, in place of $F^{\{i_1 : X_1, \dots, i_r : X_r\}}$, we usually write $F^{i_1 : X_1, \dots, i_r : X_r}$.

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 named tuple $t \in \text{ind } \mathcal{T}$; 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. Let \mathcal{U} be a shape such that $\mathcal{U} = \mathcal{S} \sqcup \mathcal{T}$ and $\mathcal{S} \sqcap \mathcal{T} = \emptyset$. If A is a tensor with shape \mathcal{S} and $s \in \text{ind } \mathcal{S}$, then we define A_s to be the named tensor with shape \mathcal{T} such that, for any $t \in \text{ind } \mathcal{T}$,

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

(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 §??, 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} such that $\mathcal{S} \sqcap \mathcal{U} = \emptyset$ and $\mathcal{T} \sqcap \mathcal{U} = \emptyset$, we can extend f to:

$$f: F^{\mathcal{S} \sqcup \mathcal{U}} \rightarrow G^{\mathcal{T} \sqcup \mathcal{U}}$$

$$[f(A)]_u = f(A_u) \quad \text{for all } u \in \text{ind } \mathcal{U}.$$

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 shape \mathcal{S}' , \mathcal{T}' such that $\mathcal{U}' = \mathcal{S}' \sqcup \mathcal{T}'$ exists, and

$$\begin{aligned} \mathcal{S} \sqcap \mathcal{S}' &= \emptyset \\ \mathcal{T} \sqcap \mathcal{T}' &= \emptyset \\ \mathcal{U} \sqcap \mathcal{U}' &= \emptyset \end{aligned}$$

we can extend f to:

$$f: F^{\mathcal{S} \sqcup \mathcal{S}'} \times G^{\mathcal{T} \sqcup \mathcal{T}'} \rightarrow H^{\mathcal{U} \sqcup \mathcal{U}'}$$

$$[f(A, B)]_u = f(A_{u|_{\mathcal{S}'}}, B_{u|_{\mathcal{T}'}}) \quad \text{for all } u \in \text{ind } \mathcal{U}'.$$

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

$$\cdot_i: F^{i:n} \times F^{i:n} \rightarrow F$$

$$A \cdot_i B = \sum_{i=1}^n A_{i:i} B_{i:i}.$$

6 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 §?? 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.

6.1 Contracting two names

We define a version of the contraction operator that can contract two indices with different names. If $i \in \text{dom } \mathcal{A}$ and $j \in \text{dom } \mathcal{B}$ and $\mathcal{A}.i = \mathcal{B}.j = X$, then we define

$$A \cdot_{i|j} B = \sum_{x \in X} A_{i:x} B_{j:x}$$

For example, the RNN would look like this.

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

6.2 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.

$$\begin{aligned} x^{(t)} &\in \mathbb{R}^{\text{emb}:d} \\ h^{(t)} &\in \mathbb{R}^{\text{state}:d} \\ A &\in \mathbb{R}^{\text{state}*:d, \text{state}:d} \\ B &\in \mathbb{R}^{\text{emb}*:d, \text{state}:d} \\ c &\in \mathbb{R}^{\text{state}:d} \\ h^{(t+1)} &= \tanh \left(A \cdot_{\text{state}} h^{(t)} + B \cdot_{\text{emb}} x^{(t)} + c \right) \end{aligned}$$

In general, if $i* \in \text{dom } \mathcal{A}$ and $i \in \text{dom } \mathcal{B}$ and $\mathcal{A}.i* = \mathcal{B}.i = X$, then we define

$$A \cdot_i B = \sum_{x \in X} A_{i*:x} B_{i:x}$$

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. \cdot_i contracts i with i , and we need another notation like $\cdot_{i(*)}$ or \times_i for contracting i^* with i .
3. \cdot_i always contracts i^* with i ; there's no way to contract i with i .

6.3 Named and numbered indices

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

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

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 $C = A \cdot_i B$.

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:

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

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:

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

Multivariate normal distribution:

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

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

Acknowledgements

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