

# Named Tensor Notation

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

Version 0.2

## Abstract

We propose a notation for tensors with named axes. Instead of writing  $A \in \mathbb{R}^{b \times w \times h \times c}$  for an order-4 tensor corresponding to a batch of images, we write  $A \in \mathbb{R}^{\text{batch} \times \text{width} \times \text{height} \times \text{channel}}$ . This tensor can be indexed using these names, as in  $A_{\text{batch}(25), \text{width}(12), \text{height}(30), \text{channel}(1)}$ , which is the same as  $A_{\text{channel}(1), \text{width}(12), \text{height}(30), \text{batch}(25)}$ . Named tensor notation relieves the author, reader, and future implementers from needing to keep track of the order of axes. It also makes it easier to partially index tensors, as in  $A_{\text{batch}(17)}$ , and to seamlessly extend operations on low-order tensors to higher order ones (e.g., extend an operation on images to batches of images, or extend the attention mechanism to multiple attention heads).

After a brief overview of our notation, we illustrate it through several examples from modern machine learning, from building blocks like attention and convolution to full model like Transformer and LeNet. Our proposals build on ideas from many previous papers and software libraries. We hope that this document will encourage more authors to use named tensors, resulting in clearer papers and less bug-prone implementations.

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.

## Contents

1	Introduction	2
2	Informal Overview	3
3	Examples	5
4	$\text{\LaTeX}$ Macros	15
5	Formal Definitions	16
6	Extensions	18

# 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 features? We would need to know the sizes of  $Q$ ,  $K$ , and  $V$  to know that it's taken over the features. Is the softmax taken over the query sequence or the key/value sequence? The usual notation doesn't even offer a way to answer this question. With multiple attention heads or multiple sentences in a minibatch, the notation becomes more complicated still.

In this document, we propose mathematical notation for tensors with *named axes*. The notation has a formal underpinning, but is hopefully intuitive enough that researchers in machine learning should be able to understand it without much effort.

As a preview, in our notation, the above equation becomes

$$\begin{aligned} \text{Att}: \mathbb{R}^{\text{seq}' \times \text{key}} \times \mathbb{R}^{\text{seq} \times \text{key}} \times \mathbb{R}^{\text{seq} \times \text{val}} &\rightarrow \mathbb{R}^{\text{seq}' \times \text{val}} \\ \text{Att}(Q, K, V) = \text{softmax}_{\text{seq}} \left( \frac{Q \underset{\text{key}}{\odot} K}{\sqrt{|\text{key}|}} \right) \underset{\text{seq}}{\odot} V. \end{aligned}$$

Here,  $K$  does not have rows or columns, so the reader does not need to remember which corresponds to the keys' features and which corresponds to elements of the sequence; instead, the dot product  $Q \underset{\text{key}}{\odot} K$  is explicitly over the **key** axis.

The resulting tensor has a **seq** axis for the key/value sequence and a **seq'** axis for the query sequence, and the softmax is explicitly over **seq**, as is the dot product with  $V$ . This formula works as written if we add a **heads** axis for multiple attention heads, or a **batch** axis for multiple sequences in a minibatch.

Our notation is inspired by libraries for programming with multidimensional arrays (Harris et al., 2020; Paszke et al., 2019) and extensions that use named axes, like Nexus (Chen, 2017), tsalib (Sinha, 2018), NamedTensor (Rush, 2019), named tensors in PyTorch (Torch Contributors, 2019), and Dex (Maclaurin et al., 2019). However, our focus is on mathematical notation rather than code.

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

In standard notation, a vector, matrix, or tensor is indexed by an integer or sequence of integers. If  $A \in \mathbb{R}^{3 \times 3}$ , then the order of the two axes matters:  $A_{1,3}$  and  $A_{3,1}$  are not the same element. It's up to the reader to remember what each axis of each tensor is for. We think this is a problem and propose a solution.

### 2.1 Named tensors

In a *named tensor*, we give each axis a name. For example, if  $A$  represents an image, we can make it a named tensor like so (writing it two equivalent ways to show that the order of axes does not matter):

$$A \in \mathbb{R}^{\text{height}[3] \times \text{width}[3]} = \mathbb{R}^{\text{width}[3] \times \text{height}[3]}$$

$$A = \text{height} \begin{array}{c} \text{width} \\ \begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 & 9 \\ 2 & 6 & 5 \end{bmatrix} \end{array} = \text{width} \begin{array}{c} \text{height} \\ \begin{bmatrix} 3 & 1 & 2 \\ 1 & 5 & 6 \\ 4 & 9 & 5 \end{bmatrix} \end{array}.$$

We access elements of  $A$  using named indices, whose order again does not matter:  $A_{\text{height}(1), \text{width}(3)} = A_{\text{width}(3), \text{height}(1)} = 4$ . We also allow partial indexing:

$$A_{\text{height}(1)} = \begin{array}{c} \text{width} \\ \begin{bmatrix} 3 & 1 & 4 \end{bmatrix} \end{array} \qquad A_{\text{width}(3)} = \begin{array}{c} \text{height} \\ \begin{bmatrix} 4 & 9 & 5 \end{bmatrix} \end{array}.$$

In many contexts, an axis name is used with only one size. If so, we can simply write `height` for the unique axis with name `height`, as in  $\mathbb{R}^{\text{height} \times \text{width}}$ . We can leave the size of an axis unspecified at first, and specify its size later (like in a section on experimental details): for example, `|height| = |width| = 28` to specify its exact size or just `|height| = |width|` to specify that it's a square image.

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, if we wanted  $A$  to have red, green, and blue channels, we'd name the axis `channels`, and if we wanted to represent a minibatch of images, we'd name the axis `batch`. Please see §3 for more examples.

### 2.2 Named tensor operations

Operations on named tensors are defined by taking a function on low-order tensors and extending it to higher-order tensors.

#### 2.2.1 Elementwise operations and broadcasting

Any function from a scalar to a scalar can be applied elementwise to a named tensor, and any function from two scalars to a scalar can be applied to two

named tensors with the same shape. For example:

$$\frac{1}{1 + \exp(-A)} = \text{height} \begin{array}{c} \text{width} \\ \begin{bmatrix} \frac{1}{1+\exp(-3)} & \frac{1}{1+\exp(-1)} & \frac{1}{1+\exp(-4)} \\ \frac{1}{1+\exp(-1)} & \frac{1}{1+\exp(-5)} & \frac{1}{1+\exp(-9)} \\ \frac{1}{1+\exp(-2)} & \frac{1}{1+\exp(-6)} & \frac{1}{1+\exp(-5)} \end{bmatrix} \end{array}.$$

But if we apply a binary function/operator to tensors with different shapes, they are *broadcast* against each other (similarly to NumPy and derivatives). Let

$$B \in \mathbb{R}^{\text{height}[3]} \qquad C \in \mathbb{R}^{\text{width}[3]}$$

$$B = \text{height} \begin{bmatrix} 2 \\ 7 \\ 1 \end{bmatrix} \qquad C = \begin{array}{c} \text{width} \\ [1 \quad 4 \quad 1] \end{array}.$$

(We write  $B$  as a column just to make the broadcasting easier to visualize.) Then, to evaluate  $A + B$ , we effectively replace  $B$  with a new tensor  $B'$  that contains a copy of  $B$  for every index of axis **width**. Likewise for  $A + C$ :

$$A + B = \text{height} \begin{array}{c} \text{width} \\ \begin{bmatrix} 3+2 & 1+2 & 4+2 \\ 1+7 & 5+7 & 9+7 \\ 2+1 & 6+1 & 5+1 \end{bmatrix} \end{array} \quad A + C = \text{height} \begin{array}{c} \text{width} \\ \begin{bmatrix} 3+1 & 1+4 & 4+1 \\ 1+1 & 5+4 & 9+1 \\ 2+1 & 6+4 & 5+1 \end{bmatrix} \end{array}.$$

### 2.2.2 Reductions

The same broadcasting rules apply to functions from vectors to scalars, called *reductions*. Unlike with functions on scalars, we always have to specify which axis reductions apply to, using a subscript. (This is equivalent to the **axis** argument in NumPy and **dim** in PyTorch.)

For example, we can sum over the **height** axis or the **width** axis of  $A$ :

$$\sum_{\text{height}} A = \sum_i A_{\text{height}(i)} = \begin{array}{c} \text{width} \\ [3+1+2 \quad 1+5+6 \quad 4+9+5] \end{array}$$

$$\sum_{\text{width}} A = \sum_j A_{\text{width}(j)} = \begin{array}{c} \text{height} \\ [3+1+4 \quad 1+5+9 \quad 2+6+5] \end{array}.$$

We can also write multiple names to perform the reduction over multiple axes at once. For example,

$$\sum_{\text{height, width}} A = \sum_i \sum_j A_{\text{height}(i), \text{width}(j)} = 3+1+4+1+5+9+2+6+5.$$

The vector dot-product is a function from *two* vectors to a scalar, which generalizes to named tensors to give the ubiquitous *contraction* operator. You can think of it as elementwise multiplication, then summation over one axis:

$$A \underset{\text{width}}{\odot} C = \sum_j A_{\text{width}(j)} B_{\text{width}(j)} = \text{height} \begin{bmatrix} 3 \cdot 1 + 1 \cdot 4 + 4 \cdot 1 \\ 1 \cdot 1 + 5 \cdot 4 + 9 \cdot 1 \\ 2 \cdot 1 + 6 \cdot 4 + 5 \cdot 1 \end{bmatrix}.$$

Again, we can write multiple names to contract multiple axes at once. An operator  $\odot$  with no axis name under it contracts zero axes and is equivalent to elementwise multiplication, so we use  $\odot$  for elementwise multiplication as well.

### 2.2.3 Renaming and reshaping

It's often useful to rename an axis (analogous to a transpose operation in standard notation):

$$[A]_{\text{height} \rightarrow \text{height}'} = \text{height}' \overset{\text{width}}{\begin{bmatrix} 3 & 1 & 4 \\ 1 & 5 & 9 \\ 2 & 6 & 5 \end{bmatrix}}.$$

We can also reshape two or more axes into one axis:

$$[A]_{(\text{height}, \text{width}) \rightarrow \text{layer}} = \overset{\text{layer}}{[3 \quad 1 \quad 4 \quad 1 \quad 5 \quad 9 \quad 2 \quad 6 \quad 5]}$$

The order of elements in the new axis is undefined. If you need a particular order, you can write a more specific definition.

## 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 Some statistics

$$\begin{aligned}\min_{\text{ax}} A &= \min\{A_{\text{ax}(i)} \mid 1 \leq i \leq n\} \\ \max_{\text{ax}} A &= \max\{A_{\text{ax}(i)} \mid 1 \leq i \leq n\} \\ \text{norm}_{\text{ax}} A &= \sqrt{\sum_{\text{ax}} A^2} \\ \text{mean}_{\text{ax}} A &= \frac{1}{n} \sum_{\text{ax}} A \\ \text{var}_{\text{ax}} A &= \frac{1}{n} \sum_{\text{ax}} (A - \text{mean}_{\text{ax}} A)^2.\end{aligned}$$

The min and max operators are overloaded, as is the summation operator defined above (§2.2.2). If the operator is applied to a tensor and has an axis under it, then it's a reduction performed over the axis. But if it is applied to a set of tensors and has no axis under it, then it's an elementwise operation performed over the set.

#### 3.1.2 Softmax and argmax

Most activation functions are elementwise operations (sigmoid, tanh, ReLU), so they are straightforward to use in our notation; the softmax, however, is interesting because it's defined as a function from vectors to vectors:

$$\text{softmax}_{\text{ax}} A = \frac{\exp A}{\sum_{\text{ax}} \exp A}.$$

As with reductions, we write an axis below the softmax operator, but this axis is retained in the output.

Closely related are argmax and argmin, which we define to compute one-hot vectors with a one at the position containing the maximum or minimum value.

$$\begin{aligned}\text{argmax}_{\text{ax}} A &= \lim_{\alpha \rightarrow \infty} \text{softmax}_{\text{ax}} \alpha A \\ \text{argmin}_{\text{ax}} A &= \lim_{\alpha \rightarrow -\infty} \text{softmax}_{\text{ax}} \alpha A.\end{aligned}$$

### 3.1.3 Fully-connected layers

A feedforward neural network looks like this:

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

The layer sizes can be set by writing  $|\text{input}| = 100$ , etc. Alternatively, we could have called the axes `layer[n0]`, `layer[n1]`, 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:

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

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  $\text{FullConn}^l$  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:

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

### 3.1.4 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.

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

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.

### 3.1.5 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:

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

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

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

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

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

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

### 3.1.6 Convolution

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

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

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

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:

$$\begin{aligned} \text{conv2d}: \mathbb{R}^{\text{channels} \times \text{height}[h] \times \text{width}[w]} &\rightarrow \mathbb{R}^{\text{height}[h'] \times \text{width}[w']} \\ \text{conv2d}(X; W, b) &= W \underset{\text{channels, kh, kw}}{\odot} U + b \end{aligned}$$

where

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

### 3.1.7 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

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

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

where

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

### 3.1.8 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

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

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:

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

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

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

where LayerNorm, SelfAtt and FFN are defined below.

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

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

We defined attention in §3.1.5; 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

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

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

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

### 3.3 LeNet

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

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

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

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

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

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

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

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

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

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

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

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

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

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  $p(\mathbf{B} \mid \mathbf{A})$  and  $p(\mathbf{A})$  as tensors:

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

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

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

### 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.

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

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.

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

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:

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

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  $\text{ax}$  and  $A \odot (\underset{\text{ax}}{\operatorname{argmaxk}} \underset{\text{ax}, k}{\operatorname{maxk}} A) = \underset{\text{ax}, k}{\operatorname{max}} A$ .

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

### 3.4.6 Multivariate normal distribution

To define a multivariate normal distribution, we need some matrix operations. These have two axis names written under them, for rows and columns, respectively. Determinant and inverse have the following signatures:

$$\begin{aligned} \det_{\text{ax1}, \text{ax2}} &: F^{\text{ax1}[n] \times \text{ax2}[n]} \rightarrow F \\ \text{inv}_{\text{ax1}, \text{ax2}} &: F^{\text{ax1}[n] \times \text{ax2}[n]} \rightarrow F^{\text{ax1}[n] \times \text{ax2}[n]}. \end{aligned}$$

(We write  $\text{inv}$  instead of  $\cdot^{-1}$  because there's no way to write axis names under the latter.)

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).

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

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

where

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

## 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:

- Basics
  - Use `\name{foo}` to write an axis name: `foo`.
  - Use `\mathbb{R}^{\text{nset}{foo}{2}}` to write a set of tensors:  $\mathbb{R}^{\text{foo}[2]}$ .
  - Use `A_{\text{nidx}{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:  $A \underset{\text{foo}}{\odot} B$ .
  - 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  $\underset{\text{foo}}{\text{qux}} A$ .

## 5 Formal Definitions

### 5.1 Records and shapes

A *named index* is a pair, written  $\mathbf{ax}(i)$ , where  $\mathbf{ax}$  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  $\{\mathbf{ax}_1(i_1), \dots, \mathbf{ax}_r(i_r)\}$ , where  $\mathbf{ax}_1, \dots, \mathbf{ax}_r$  are pairwise distinct names.

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

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

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

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

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

We say two shapes  $\mathcal{S}$  and  $\mathcal{T}$  are *compatible* if whenever  $\mathbf{ax}(I) \in \mathcal{S}$  and  $\mathbf{ax}(J) \in \mathcal{T}$ , then  $I = J$ . We say that  $\mathcal{S}$  and  $\mathcal{T}$  are *orthogonal* if there is no  $\mathbf{ax}$  such that  $\mathbf{ax}(I) \in \mathcal{S}$  and  $\mathbf{ax}(J) \in \mathcal{T}$  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  $A_{\{\mathbf{ax}_1(x_1), \dots, \mathbf{ax}_r(x_r)\}}$ , we usually write  $A_{\mathbf{ax}_1(x_1), \dots, \mathbf{ax}_r(x_r)}$ . When a named tensor is an expression like  $(A + B)$ , we surround it with square brackets like this:  $[A + B]_{\mathbf{ax}_1(x_1), \dots, \mathbf{ax}_r(x_r)}$ .

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 Named tensor operations

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 can be defined as:

$$\odot_{\mathbf{ax}}: F^{\mathbf{ax}[n]} \times F^{\mathbf{ax}[n]} \rightarrow F$$

$$A \odot_{\mathbf{ax}} B = \sum_{i=1}^n A_{\mathbf{ax}(i)} B_{\mathbf{ax}(i)}.$$

## 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 ( $V = \{\text{cat}, \text{dog}, \dots\}$ ), 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  $\mathbf{u}$  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, instead

of defining an  $\text{inv}$  operator that takes two subscripts, we could write

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

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

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

## 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 §3.1.4 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):

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

For example, the RNN would look like this.

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

### 6.2.2 Starred axis names

If  $\mathbf{ax}$  is a name, we also allow a tensor to have an axis  $\mathbf{ax}^*$  (alternatively: superscript  $\mathbf{ax}$ ). Multiplication contracts starred axes in the left operand with non-starred axes in the right operand.

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

The contraction operator can be defined as:

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

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

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

### 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:

$A_{\mathbf{ax}}$	rename axis $\mathbf{ax}$ to 1
$A_{\mathbf{ax}1, \mathbf{ax}2}$	rename axis $\mathbf{ax}1$ to 1 and $\mathbf{ax}2$ to 2
$A_{\rightarrow \mathbf{ax}}$	rename axis 1 to $\mathbf{ax}$
$A_{\rightarrow \mathbf{ax}1, \mathbf{ax}2}$	rename axis 1 to $\mathbf{ax}1$ and 2 to $\mathbf{ax}2$

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  $C = A \underset{\text{ax}}{\odot} B$ .

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  $\max_{\text{ax}}$ ), 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  $\text{ax1}$  to  $\text{ax2}$ .

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

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

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 axis names twice), we'd keep around the notation  $A \odot_{\text{ax}} B$  as a shorthand for  $A_{\text{ax}} \cdot B_{\text{ax}}$ . We'd also keep named reductions, or at least  $\text{softmax}_{\text{ax}}$ .

### 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  $\iota^1, \dots, \iota^D$  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, Adam Poliak, Matt Post, Chungchieh Shan, Nishant Sinha, and Yee Whye Teh 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.