

Neural network knowledge distillation in tensor networks

Dereck Piché

March 21, 2023

Abstract

1 Introduction

TODO

2 Knowledge Distillation

Knowledge Distillation is a machine learning practice which involves taking a trained model and using it's parameters to train another one. The already trained model is referred to as the "teacher", and the model in which his "knowledge" is to be distilled is referred to as the "student".

2.1 Response-Based Knowledge Distillation

While this is a relatively novel technique, there are already several distinct approaches introduced by researchers. We used two of these approaches. Our inspiration for the distillation methodology was found in a 2021 survey which resumed the emerging practice Gou et al. [1]

The first approach we used to look exclusively at the outputs of the student and teacher. Now, we apply the logits of our functions element-wise to the outputs and the softmax.

$$\text{softmax}(v_i) = \frac{e^{v_i}}{\sum_j e^{v_j}} \quad (1)$$

The softmax function's objective is to transform the logits into probability distributions for the different classes. We will now apply a loss function L to

these two functions. Since we are trying to reduce the divergence between two distributions, we will use the Kullback-Leibler divergence loss

$$\text{KL}(\mathbf{P}, \mathbf{Q}) = \frac{1}{n} \sum_i \mathbf{Q} \frac{\log_e(\mathbf{Q})}{\log_e(\mathbf{P})} \quad (2)$$

Here, \mathbf{Q} is the distribution we are aiming for and \mathbf{P} is the one we have.

3 Tensor Networks

Tensor Networks come from the study of quantum phenomena. They started being used recently as machine learning models. Tensor Networks can be thought out as two things: a visual notation system and a set of methods for tensor manipulation.

Tensors Before explaining these two, we shall disambiguate the meaning of "tensor". In this report (and very often in the context of machine learning), we use the word "tensor" to refer to the mathematical arrays of arbitrary indices. Here, the number of indices is called the "order" of the tensor, meaning that vectors are simply tensors of order 1 and matrices tensors of order 2.

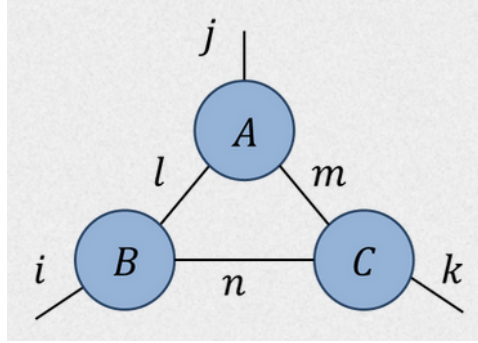
Contraction At the heart of tensor networks is the *contraction* operation. Tensor Networks are used to compute a larger network by "contracting" several smaller tensors over chosen indices. A "contraction" is simply an operation where we sum over indices. For example, the contraction of \mathbf{A}_{ijk} and \mathbf{B}_{ijk} on index j will produce the tensor $\mathbf{C}_{ik} = \sum_j \mathbf{A}_{ijk} \mathbf{B}_{ijk}$. Evidently, the two indices present in a contraction must be of the same size.

Notational system The principal motivation behind the creation of Tensor Networks to compute or approximate large tensors *contracting* several smaller tensors over chosen indices. This is all well and good until our conventional summation notation begins overclocking our primal brains. So, in order to make the manipulation of these network of tensors, mathematicians created a notational system. Tensors are represented as nodes where each vertex connected to the node represents one of the tensor's indices. If a vertex connects two nodes, it means that the indices shall be contracted in order to produce the post-contraction tensor. It should be evident that by the definitions above, no node (tensor) is completely isolated in a tensor

network, as it would be completely purposeless. The shape of the post-contraction tensor can be easily visually identified, since it is found in the unconnected vertices. A simple tensor network can be found in figure 1.

Tensor Network Methods The term *Tensor Network Methods* is used to refer to, you guessed it, the methods. There are several architectures of Tensor Networks that are frequently used, such as the *Matrix Product State (MPS)*, the *Tenso*

Figure 1: Simple tensor network illustrating the notation.



4 Kernel

Consider an input vector $\mathbf{x} \in \mathbb{R}^d$. Let $\phi(\mathbf{x}) : \mathbb{R} \mapsto \mathbb{R}^{d_\phi}$. Then, let us take the tensor product of apply $\phi(\mathbf{x})$ applied to every element of the vector \mathbf{x} .

$$\Phi(\mathbf{x}) = \phi(x_1) \otimes \phi(x_2) \otimes (\dots) \phi(x_d) \quad (3)$$

We obtain a tensor $\Phi(\mathbf{x})$, of which the sum of it's element contain the basis of a space of products of the elements of the transformed elements of $\phi(\mathbf{x})$.

In order to reduce the abstractness of this statement, we can take say that \otimes refers to the Kronecker Product.

4.1 Multilinear feature map

In this project, we attributed a particular importance to the local feature map

$$\phi^*(\mathbf{x}) = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \quad (4)$$

With, this particular transformation, the Kronecker Product gives us a global feature map $\Phi(\mathbf{x})$ which is a tensor where each element is a basis of the space of multilinear functions on the elements of the vector \mathbf{x} . Here, the order or shape of the tensor is superfluous and gives no theoretical advantage.

4.2 Linear Combinations

Now, the question is how do we use $\Phi(\mathbf{x})$, how does it become useful? Well, it becomes useful when we take a linear map of it's elements. Let our $\Phi(\mathbf{x})$ be an arbitrarily shaped tensor of d_Φ dimensions. The final form of our function $f(\mathbf{x})$ will be

$$f(\mathbf{x}) = \sum_i^{d_\Phi} \theta_i [\Phi(\mathbf{x})]_i \quad (5)$$

Now, we have a function that can, under some choice of transformation ϕ , become extremely expressive. For our particular local feature map ϕ^* , it can represent any multilinear function of the input vector \mathbf{x} .

5 Using The Matrix Product State Tensor Network

In this project, we used the *Matrix Product State* (MPS) tensor network to approximate the function $\mathbf{T}\Phi(\mathbf{x})$.

5.1

5.2 Expressivity of MPS combinations with transformation $[1, \mathbf{x}]^t$

Let $\mathbf{x} \in \mathbb{R}^d$. Let $f(\mathbf{x})$ be a function that returns a vector \mathbf{v} , which contains every element required to form a basis of the space of x_i -variate multilinear functions.

Let

$$\Theta = \{\theta_1, \theta_2, (\dots), \theta_{d_{\text{theta}}}\} \quad (6)$$

in

$$M(f(\mathbf{x}), \Theta) = \begin{bmatrix} m(f(\mathbf{x}), \theta_1) \\ m(f(\mathbf{x}), \theta_2) \\ (\dots) \\ m(f(\mathbf{x}), \theta_{d_{\text{theta}}}) \end{bmatrix} \quad (7)$$

Where $m(f(x), \theta_i) : \mathbb{R}^{2^d} \mapsto (\mathbb{R}^d \mapsto \mathbb{R})$ is of the form

$$m(f(x), \theta_i) = m(v, \theta_i) = \sum_j \theta_{i,j} \cdot v_j \quad (8)$$

We can trivially choose $\Theta^* \in \Theta$ such that

$$M(f(x), \Theta^*) = \begin{bmatrix} \text{repeated } z \text{ times } \begin{Bmatrix} x_1 \\ \dots \end{Bmatrix} \\ \text{repeated } z \text{ times } \begin{Bmatrix} x_2 \\ \dots \end{Bmatrix} \\ \dots \\ \text{repeated } z \text{ times } \begin{Bmatrix} x_n \\ \dots \end{Bmatrix} \end{bmatrix} = \lambda \quad (9)$$

We can rewrite λ as

$$\lambda = \begin{bmatrix} \varepsilon_{1,1} \\ \dots \\ \varepsilon_{1,z} \\ \varepsilon_{2,1} \\ \dots \\ \varepsilon_{2,z} \\ \dots \\ \varepsilon_{d,z} \end{bmatrix} = \lambda^* \quad (10)$$

Let Z be the space of x_i -variate polynomial functions of degree $\leq z$. Then every monomial of any x_i -variate polynomial function $\zeta \in Z$ is of the form.

$$x_1^{k_1} x_2^{k_2} (\dots) x_d^{k_d} \quad (11)$$

under the condition $\sum_i k_i \leq z$. However, for every $K = \{k_1, k_2, (\dots), k_d\}$ meeting this condition, we can rewrite the monomial as

$$\left(\prod_{i_1=1}^{k_1} \varepsilon_{1,i_1} \right) \left(\prod_{i_2=1}^{k_2} \varepsilon_{2,i_2} \right) (\dots) \left(\prod_{i_d=1}^{k_d} \varepsilon_{d,i_d} \right) \quad (12)$$

by using the elements of λ^* from (10). However, we clearly see that this term is a multilinear monomial of the variables in λ^* . This implies that $f(\lambda^*)$ returns a basis-vector of x_i -variate polynomial function of degree $\leq z$.

In other words,

$$\xi(x, \Theta) = M(f(M(f(x), \Theta^*)), \Theta) \quad (13)$$

can express any x_i -variate polynomial function of degree $\leq z$ under fixed Θ .

6 Further Exploration

TODO talk about capturing locality in the mappings
TODO talk about capturing locality in general with tensors

References

- [1] Jianping Gou et al. “Knowledge Distillation: A Survey”. In: *International Journal of Computer Vision* 129.6 (Mar. 2021), pp. 1789–1819. DOI: 10.1007/s11263-021-01453-z. URL: <https://doi.org/10.1007/s11263-021-01453-z>.