

Neural network knowledge distillation in tensor networks

Dereck Piché

March 10, 2023

Abstract

1 Introduction

void

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 called the "student".

2.1 Response-Based Knowledge Distillation

The first approach is to look exclusively at the outputs of the student and teachers. 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

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

Here, Q is the distribution we are aiming for and 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.

Consider a input vector \mathbf{x} .

3.1 Tensor Products and Transformations

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.

A particular dimension is of interest here, the transformation

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

With, this particular transformation, the Kronecker Product gives us a $\Phi(\mathbf{x})$ which is a matrix where each element is a basis of the space of multilinear functions on the elements of the vector \mathbf{x} .

3.2 Linear Combinations

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.

3.3 The *Matrix Product State* Tensor Network

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

Let $\mathbf{x} \in \mathbb{R}^d$. Let $f(\mathbf{x})\mathbb{R}^d \mapsto \mathbb{R}^{z^d}$ function that returns a vector \mathbf{v} . Let $f(\mathbf{x})$ be defined such that \mathbf{v} contains every \mathbf{x}_i -variate monome of degree $\leq z$ where the variables are the elements of \mathbf{x} .

Let

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

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

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

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

We can choose Θ such that

$$M(f(\mathbf{x}), \Theta) = \begin{bmatrix} x_1 \\ x_1 \\ (\dots) \\ x_d \end{bmatrix} = \lambda \quad (9)$$

We can rewrite λ as

$$\lambda = [\varepsilon_{1,1} \quad (\dots) \quad \varepsilon_{1,z} \quad \varepsilon_{2,1} \quad (\dots) \quad \varepsilon_{2,z} \quad (\dots) \quad \varepsilon_{d,z}]^T \quad (10)$$

Every monome of ζ is of the form

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

where $\sum_i k_i \leq z$. However, for every $K = k_1, k_2, k_d$ meeting this condition, we can rewrite as

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

However, we clearly see that this term is a multilinear monome of the variables in λ .