

Mathematics for Deep Learning

Zheguang Zhao
zheguang.zhao@gmail.com

Last updated: February 28, 2019

1 Introduction

The goal of this article is to establish some mathematical background of one of the fundamental tool for building deep learning: the computation graphs.

Computation graphs have been used extensively not only to communicate different deep learning architectures but also to program as the building block in software frameworks such as TensorFlow and Keras. One advantage of using computation graph is to be able to think about neural networks at a level of abstraction that is clear and effective, while keeping the mathematical details just one step away.

However the idea of computation graphs is not new. It has roots in calculus, long before computers were invented. So in the following I will first establish the correspondence between computation graphs and functions. Most importantly, I will show how deep learning networks are just composition of functions, and composition of functions is just connectivity of multiple computation subgraphs. One consequence of this is that we can change a neural network by simply modifying its graph structure.

2 Single-variable real-valued function

Let $z \in \mathbb{R}$, and $f(z)$ be a scalar map from scalars to scalars, i.e. $\mathbb{R} \rightarrow \mathbb{R}$:

$$\frac{df(z)}{dz} := \lim_{\delta_z \rightarrow 0} \frac{f(z + \delta_z) - f(z)}{\delta_z} \quad (1)$$

The computation graphs representation of $f(z)$ and its derivative $f'(z)$ are both one-to-one, which reflects that they are both scalar-to-scalar functions:

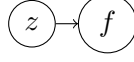


Figure 1: $f(z)$

Figure 2: Single-variable real-valued function.

3 Vector-valued functions

Let $t \in \mathbb{R}$, and $\mathbf{x} = \mathbf{f}(t) := (x_1(t), x_2(t), \dots, x_3(t))$ be a map of scalars to vectors, i.e. $\mathbb{R} \rightarrow \mathbb{R}^3$

The derivative of this vector-valued function w.r.t. its scalar variable is simply the the derivatives of its components w.r.t. this scalar variable:

$$\frac{d\mathbf{x}}{dt} := \left(\frac{dx_1}{dt}, \frac{dx_2}{dt}, \frac{dx_3}{dt} \right) \quad (2)$$

We need to use three nodes to denote the components of \mathbf{x} in the computation graph:

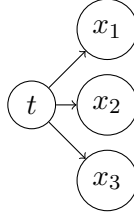


Figure 3: Vector-valued function

4 Multi-variable real-valued functions

Let $z = f(\mathbf{x}) = f(x_1, x_2, x_3)$ be a map from $\mathbb{R}^3 \rightarrow \mathbb{R}$.

The derivative of z w.r.t. \mathbf{x} is also called the gradient of z , denoted as $\nabla_{\mathbf{x}} z$:

$$\nabla_{\mathbf{x}} z = \frac{dz}{d\mathbf{x}} := \left(\frac{\partial z}{\partial x_1}, \frac{\partial z}{\partial x_2}, \frac{\partial z}{\partial x_3} \right) \quad (3)$$

A multi-variable real-valued function can be represented as many-to-one graph:

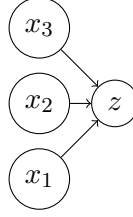


Figure 4: Multi-variable real-valued function

5 Function composition

Let $\mathbf{x} \in \mathbb{R}^3$, $t, z \in \mathbb{R}$, $\mathbf{x} = \mathbf{f}(t)$ be a vector-valued function, $z = g(\mathbf{x})$ be a multi-variable real-valued function. Then we can compose a map $h = \mathbf{f} \circ g$ from $\mathbb{R} \rightarrow \mathbb{R}$ as $z = h(t) = (\mathbf{f} \circ g)(t) = g(\mathbf{f}(t))$.

The computation graph of h is just a concatenation of the subgraphs of a vector-valued function \mathbf{f} (Figure 3) and g (Figure 4):

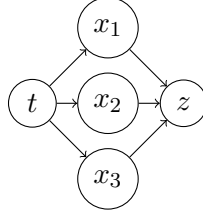


Figure 5: $z = h(t) = (\mathbf{f} \circ g)(t)$

What is the derivative of z w.r.t. t ? Because h is a function composition, we need to use chain rule.

$$\frac{dz}{dt} = \frac{dh(t)}{dt} \quad (4)$$

$$= \frac{d(\mathbf{f} \circ g)(t)}{dt} \quad (5)$$

$$= \frac{dg(\mathbf{x})}{d\mathbf{x}} \cdot \frac{d\mathbf{f}(t)}{dt} \quad (6)$$

$$= \frac{dz}{d\mathbf{x}} \cdot \frac{d\mathbf{x}}{dt} \quad (7)$$

$$= \left(\frac{\partial z}{\partial x_1}, \frac{\partial z}{\partial x_2}, \frac{\partial z}{\partial x_3} \right) \cdot \left(\frac{dx_1}{dt}, \frac{dx_2}{dt}, \frac{dx_3}{dt} \right) \quad (8)$$

$$= \frac{\partial z}{\partial x_1} \frac{dx_1}{dt} + \frac{\partial z}{\partial x_2} \frac{dx_2}{dt} + \frac{\partial z}{\partial x_3} \frac{dx_3}{dt} \quad (9)$$

This shows that the derivative is also just the dot product of the two functions of the two subgraphs. This establishes the correspondence between the graph composition and the function composition.

6 License

This work is licensed under a Creative Commons “Attribution-NonCommercial-ShareAlike 3.0 Unported” license.

