

# Neural Networks

Kieran Harvie

Copyright ©July 27, 2023. All Rights Reserved.

This document is still **under construction**.  
Changes in content, structure, and readability are expected.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Derivatives</b>	<b>4</b>
2.1	A Tale of Two Derivative . . . . .	4
2.2	Dual Numbers . . . . .	4
2.2.1	Historical Note . . . . .	5
2.3	Automatic Differentiation . . . . .	6
2.3.1	Forward Accumulation . . . . .	6
2.3.2	Reverse Accumulation . . . . .	7
<b>3</b>	<b>Optimization</b>	<b>8</b>
3.1	Gradient Descent . . . . .	8
3.2	Newton's Method . . . . .	8
<b>4</b>	<b>Nodes</b>	<b>9</b>
4.1	Layered Structure . . . . .	9
<b>A</b>	<b>Appendix</b>	<b>10</b>
A.1	Logistic Function . . . . .	10
A.2	Table of Elementary Functions of dual numbers . . . . .	10

# Introduction

# Derivatives

General

## 2.1 A Tale of Two Derivative

Partial and total. Chain rule.  $\nabla$  operator maybe an appendix for vectors?

$$\frac{df}{dt} = \sum_i \frac{\partial f}{\partial x_i} \frac{dx_i}{dt} = \nabla f \cdot \left( \frac{dx_0}{dt}, \frac{dx_1}{dt}, \frac{dx_2}{dt}, \dots \right)$$

## 2.2 Dual Numbers

Like how we get Complex numbers from Real numbers by adjoining a new number  $i$  such that  $i^2 = -1$  we get the Dual numbers by adjoining a new number  $\epsilon$ , called the infinitesimal, such that  $\epsilon^2 = 0$ .

Some properties aren't that special:

$$(x + x'\epsilon) + (y + y'\epsilon) = (x + y) + (x' + y')\epsilon$$

But consider what happens for multiplication and division:

$$\begin{aligned}(x + x'\epsilon)(y + y'\epsilon) &= xy + (x'y + xy')\epsilon + x'y'\epsilon^2 \\ &= xy + (x'y + xy')\epsilon \\ \frac{x + x'\epsilon}{y + y'\epsilon} &= \frac{(x + x'\epsilon)(y - y'\epsilon)}{(y + y'\epsilon)(y - y'\epsilon)} \\ &= \frac{xy + (x'y - xy')\epsilon - x'y'\epsilon^2}{y^2 - y'^2\epsilon^2} \\ &= \frac{x}{y} + \frac{x'y - xy'}{y^2}\epsilon\end{aligned}$$

Define:  $f(\langle u, u' \rangle, \langle v, v' \rangle) = \langle f(u, v), \nabla f(u, v) \cdot (u', v') \rangle$

Because:

$$\begin{aligned} f\left(\left\langle u, \frac{du}{dt} \right\rangle, \left\langle v, \frac{dv}{dt} \right\rangle\right) &= \left\langle f(u, v), \nabla f(u, v) \cdot \left(\frac{du}{dt}, \frac{dv}{dt}\right) \right\rangle \\ &= \left\langle f(u, v), \frac{d}{dt} f(u, v) \right\rangle \end{aligned}$$

(Because of the similar form this argument also works with partials).

### 2.2.1 Historical Note

As observed with the multiplication and divisions examples, representing infinitesimals with  $\epsilon^2 = 0$  means discarding any powers of  $\epsilon$  larger than 1.

This actually mirrors how calculus was developed. In 1710 Leibniz codified the "Transcendental Law of Homogeneity" which states that when equating sums involving to only include the lowest order infinitesimal terms.

For example, if we have infinitesimals  $du$  and  $dv$  the Transcendental Law of Homogeneity would mean that:

$$dv^2 + dudv + 2dv = 2dv$$

This can be used to calculate derivative of a function by subtracting the version with finite variables from infinitesimals ones:

$$(v + dv)(u + du) - uv = vdu + udv + dudv = vdu + udv$$

Which clearly mimics the multiplication example.

A historical note on a historical note, this specific application of the rule is similar to "Adequation" which can be traced back to Pierre de Fermat in a 1636 treatise. Here we find the extrema of a function  $f$  at a value  $x$  "adequating"  $f(x)$  to  $f(x + e)$ , then dividing by  $e$  and discarding any remaining terms involving  $e$ .

For a worked example consider  $f(x) = x^2 + x + 1$  and represent "adequal-

ity" with  $\sim$ :

$$\begin{aligned}
f(x) &\sim f(x+e) \\
x^2 + x + 1 &\sim (x+e)^2 + (x+e)1 \\
&\sim x^2 + 2xe + e^2 + x + e + 1 \\
0 &\sim 2xe + e^2 + e \quad (\text{Cancellation}) \\
0 &\sim 2x + e + 1 \quad (\text{Division by } e) \\
0 &\sim 2x + 1 \quad (\text{Discarding } e)
\end{aligned}$$

Observe the similarity with the previous arguments and boils down to  $\frac{df}{dx} = 0$ .

The purpose of this historical tangent is to establish manipulating infinitesimals like this is not some fringe idea related solely to this application.

## 2.3 Automatic Differentiation

Automatic Differentiation are a collection of techniques for algorithmically calculating the partial derivative of a function in some general way. The core observation is that if the functions can be expressed in terms of elementary functions with known partial derivatives we should be able to use the chain rule to calculate the partial derivative of the original function.

### 2.3.1 Forward Accumulation

Forward Accumulation is the most direct form of automatic differentiation.

Consider the following relation:

$$f(\langle x, 1 \rangle, \langle y, 0 \rangle) = \left\langle f(x, y), \frac{\partial}{\partial x} f(x, y) \right\rangle$$

Since dual numbers are easy for computers to represent and manipulate we can calculate  $\frac{\partial f}{\partial x}$  by substituting in  $\langle x, 1 \rangle$  and  $\langle y, 0 \rangle$  and reading off the final infinitesimal.

Consider  $f(x, y) = \cos(xy) + y$ ,  $\cos$  is an elementary function for which we can easily verify:

$$\cos(\langle x, x' \rangle) = \langle \cos(x), -x' \sin(x) \rangle$$

Hence we can calculate the  $\frac{\partial f}{\partial x}$  at  $x = 0$  and  $y = 1$  as:

$$\begin{aligned} f(\langle 0, 1 \rangle, \langle 1, 0 \rangle) &= \cos(\langle 0, 1 \rangle \times \langle 1, 0 \rangle) + \langle 1, 0 \rangle \\ &= \cos(\langle 0, 1 \rangle) + \langle 1, 0 \rangle \\ &= \langle 1, 0 \rangle + \langle 1, 0 \rangle \\ &= \langle 2, 0 \rangle \end{aligned}$$

Like wise for  $\frac{\partial f}{\partial y}$ :

$$\begin{aligned} f(\langle 0, 0 \rangle, \langle 1, 1 \rangle) &= \cos(\langle 0, 0 \rangle \times \langle 1, 1 \rangle) + \langle 1, 1 \rangle \\ &= \cos(\langle 0, 0 \rangle) + \langle 1, 1 \rangle \\ &= \langle 1, 0 \rangle + \langle 1, 1 \rangle \\ &= \langle 2, 1 \rangle \end{aligned}$$

This method is call *forward* accumulation because we exclusively move from inputs,  $x$  and  $y$ , to intermediate results,  $\cos(x, y)$ , to outputs,  $f(x, y)$ . This has its advantages of letting us reuse how computers call functions to store intermediate results and avoid doing any overhead work. But its main flaw is only calculating one partial at a time.

We can calculate multiple partials at time with the next method.

### 2.3.2 Reverse Accumulation

In reverse accumulation we do a forward pass to calculate and store  $f(x, y)$ , and intermediates, but then do a second pass in the reverse direction to calculate the partials by fixing the final infinitesimal as 1. To see the intuition as to why this would work consider:

$$f\left(\left\langle x, \frac{dx}{df} \right\rangle, \left\langle y, \frac{dy}{df} \right\rangle\right) = \left\langle f(x, y), \frac{df}{df} \right\rangle = \langle f(x, y), 1 \rangle$$

Although this method requires multiple passes and the overhead of storing intermediate results it calculated all partials at the same time. Adjoint:

$$\bar{x} = \frac{\partial f}{\partial x}$$

Diamond with a forward value calculation and a reverse adjoint accumulation.



# Optimization

## 3.1 Gradient Descent

## 3.2 Newton's Method

Like a dynamic step size since we have this info? Only when root finding. (Use dual numbers).

Using error  $l_2$  all the functions are twice differentiable so you can do the Lagrange form of the remainder thing. Also only one root so how does that work with basins

$$x_{n+1} = x_n - \frac{F(x_n)}{|\nabla F(x_n)|^2} \nabla F(x_n)$$

# Nodes

Transfer (SUM) Activator (Function (logistic)) Loss/Cost (Error)

## 4.1 Layered Structure

Back propagate but only keep one row in memory at a time.

# Appendix

## A.1 Logistic Function

## A.2 Table of Elementary Functions of dual numbers