

Backprop as Functor: A compositional perspective on supervised learning

David I. Spivak (joint with Brendan Fong and Rémy Tuyéras)

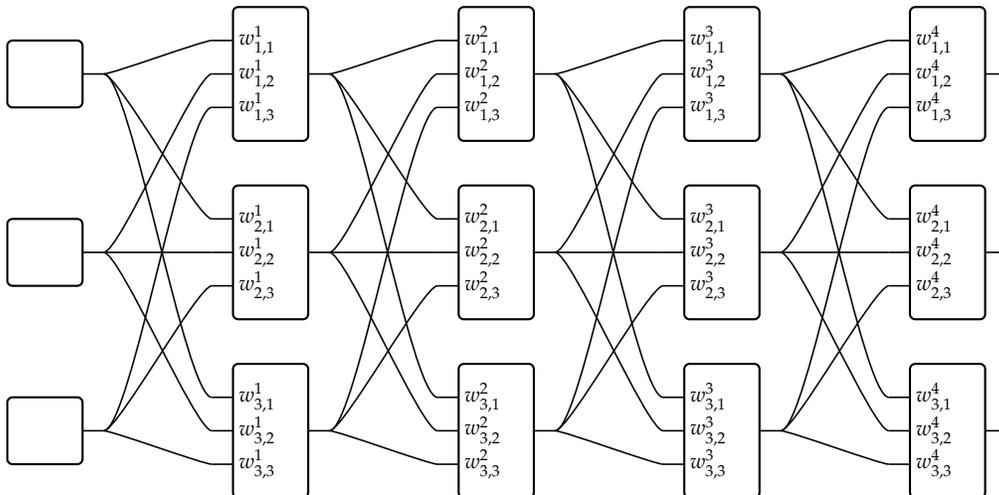
2018/05/06 at Max Planck Institute

Abstract

Abstract: Neural networks can be trained to perform functions, such as classifying images. The usual description of this process involves keywords like neural architecture, activation function, cost function, back propagation, training data, weights and biases, and weight-tying.

In this talk we will define a symmetric monoidal category \mathbf{Learn} , in which objects are sets and morphisms are roughly “functions that adapt to training data”. The back propagation algorithm can then be viewed as a strong monoidal functor from a category of parameterized functions between Euclidean spaces to our category \mathbf{Learn} .

This presentation is algebraic, not algorithmic; in particular it does not give immediate insight into improving the speed or accuracy of neural networks. The point of the talk is simply to articulate the various structures that one observes in this subject—including all the keywords mentioned above—and thereby get a categorical foothold for further study. For example, by articulating the structure in this way, we find functorial connections to the well-established category of lenses in database theory and the much more recently-developed category of compositional economic games.



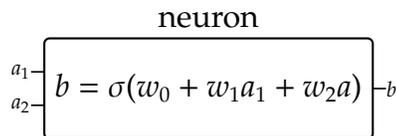
Backprop as Functor: A compositional perspective on supervised learning

Joint with Brendan, Remy

I. Introduction

A. What is supervised learning?

1. The usual story:
 - a. neural architecture, weights+biases, activation fn.,



- b. cost fn., step size
 - c. training data
 - d. gradient descent, backpropagation
 2. Extra bells and whistles: Convolutional&recurrent neural nets
 - a. Convolutional neural nets: image classification
 - b. Recurrent neural nets: time series
 - c. Idea: *Shift-invariance*: look for same patterns in different places
 - d. Uses something called “weight tying”
- B. Our goal: articulate the structures categorically
 1. What we have
 - a. Monoidal category setting **Learn** for general supervised learning
 - b. Special case: functors $\text{NNet} \rightarrow \text{Para} \rightarrow \text{Learn}$
 - c. Weight tying, time permitting
 2. What remains to be done
 - a. Address accuracy, speed, etc. (nothing yet)
 - b. Flesh out connection to lenses and combinatorial games

II. Compositional supervised learning

A. Basic structure isn't compositional

1. Basic structure of morphism $A \rightarrow B$:
 - a. Parameter space P ,
 - b. Parameterized function $I: P \times A \rightarrow B$, and
 - c. Update function: $U: P \times A \times B \rightarrow P$.
2. Problem with composition: need local training data

B. The fix: request function

1. Add to the list: $r: P \times A \times B \rightarrow A$.

2. Mainly used for compositionality: generalized backprop
3. Also for *dreaming*: “make input more cat-like”

C. The symmetric monoidal category **Learn**

1. $\text{Ob}(\mathbf{Learn}) = \text{Ob}(\mathbf{Set})$
2. $\mathbf{Learn}(A, B) = (P, I, U, r)$ up to equivalence
 - a. Any surjective $f: P \twoheadrightarrow P'$ commuting with I, I', U, U', r, r' is equivalence.
3. Monoidal product: \times across the board; unit: **1**
4. Neural network diagrams are string diagrams.

III. The strong monoidal functor $\mathbf{Para} \rightarrow \mathbf{Learn}$

A. The prop **Para**

1. $\text{Ob}(\mathbf{Para}) = \mathbb{N}$; think \mathbb{R}^n ; I’ll often denote them as \mathbb{R}^n
2. $\mathbf{Para}(m, n) = \{(k, I) \mid k \in \mathbb{N}, I: \mathbb{R}^k \times \mathbb{R}^m \rightarrow \mathbb{R}^n \text{ smooth}\}$
3. Neural networks are special morphisms in **Para**

B. The functors $L_{\varepsilon, e}: \mathbf{Para} \rightarrow \mathbf{Learn}$

1. Choices
 - a. Choose step size $\varepsilon > 0$
 - b. Choose cost function $e: \mathbb{R}^{\{g, b\}} \rightarrow \mathbb{R}$ satisfying technical condition:
 - c. $\frac{\partial e}{\partial g}(g_0, -): \mathbb{R}^{\{b\}} \xrightarrow{\cong} \mathbb{R}$ an isomorphism for each “guess” g_0
2. Obtain a strong monoidal functor $L_{\varepsilon, e}: \mathbf{Para} \rightarrow \mathbf{Learn}$
 - a. For simplicity, assume quadratic error: $e(g, b) = \frac{1}{2}(g - b)^2$.
 - b. On objects: $n \mapsto \mathbb{R}^n$.
 - c. On morphisms: $(P, I) \mapsto (P, I, U, r)$ where

$$U_I(p, a, b) := p - \varepsilon \nabla_p E_I(p, a, b)$$

$$r_I(p, a, b) := a - \nabla_a E_I(p, a, b)$$

$$\text{where } E_I(p, a, b) = \frac{1}{2} \|I(p, a) - b\|^2.$$

C. What does strong monoidal functoriality mean?

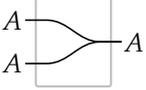
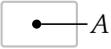
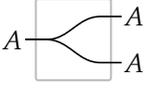
1. We can compose string diagrams
2. Gradient descent can be done locally, and requests are back-propagated.
3. (Cheaper).

IV. Other structures

A. A bimonoid on each object $\mathbb{R}^n \in \mathbf{Learn}$

1. Let **FVect** denote category of f.d. \mathbb{R} -vector spaces and linear maps.
2. Have $\mathbf{FVect} \hookrightarrow \mathbf{Para} \hookrightarrow \mathbf{Learn}$ (for any ε, e).

- a. image objects of $\mathbf{FVect} \hookrightarrow \mathbf{Para}$ have parameter space $\{*\}$.
 3. Bimonoid structure on each object of the image. [see next page]
- B. Weight-tying
1. Each morphism in \mathbf{Para} can be factored.
 2. $(P, I): A \rightarrow B$ vs. $(P, \text{id}): \mathbf{1} \rightarrow P$ and $(\mathbf{1}, I): P \times A \rightarrow B$.
 3. Weight tying [see next page]

	Implementation	Request
Multiplication, μ $(1, I_\mu, !, r_\mu)$ 	$I_\mu: A \times A \longrightarrow A;$ $(a_1, a_2) \mapsto a_1 + a_2$	$r_\mu: (A \times A) \times A \longrightarrow A \times A;$ $((a_1, a_2), a_3) \mapsto (a_3 - a_2, a_3 - a_1)$
Unit, η $(1, I_\eta, !, r_\eta)$ 	$I_\eta: \mathbb{R}^0 \longrightarrow A;$ $0 \mapsto 0$	$r_\eta: A \longrightarrow \mathbb{R}^0;$ $a \mapsto 0$
Comultiplication, δ $(1, I_\delta, !, r_\delta)$ 	$I_\delta: A \longrightarrow A \times A;$ $a \mapsto (a, a)$	$r_\delta: A \times (A \times A) \longrightarrow A;$ $(a_1, (a_2, a_3)) \mapsto a_2 + a_3 - a_1$
Counit, ϵ $(1, I_\epsilon, !, r_\epsilon)$ 	$I_\epsilon: A \longrightarrow \mathbb{R}^0;$ $a \mapsto 0$	$r_\epsilon: A \longrightarrow A;$ $a \mapsto 0$

