# Backprop as Functor

Brendan Fong, with David Spivak, Rémy Tuyéras, Mike Johnson

2nd Workshop on Open Games
Oxford
5 July 2018
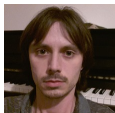
Consider the function:

$$\text{Cat?: Pictures} = \mathbb{R}^{100 \times 100 \times 3} \longrightarrow \langle \texttt{cat}, \texttt{not\_cat} \rangle = \mathbb{R}^2$$

 $\longmapsto 1.00|\texttt{cat}\rangle + 0.00|\texttt{not\_cat}\rangle$

 $\longmapsto 0.12|\texttt{cat}\rangle + 0.95|\texttt{not\_cat}\rangle$

 $\longmapsto 1.00|\texttt{cat}\rangle + 1.00|\texttt{not\_cat}\rangle$

How do we program it?

**Outline**

# I. Supervised Learning, Compositionally

**Goal: learn a function from examples**

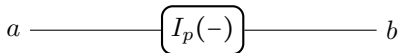Fix sets $A$, $B$. For all $f: A \to B$, use pairs $(a, f(a))$ to *approximate* $f$.

**Method: use the following data**

Hypothesis set: $P$

Implementation function: $I: P \times A \to B$

Update function $U: P \times A \times B \to P$

Request function $r: P \times A \times B \to A$

$$a \quad\overline{\qquad\qquad}\boxed{I_p(-)}\overline{\qquad\qquad}\quad b$$

A **learner** $A \to B$ is a tuple* $(P, I, U, r)$.

---

*actually an equivalence class.

**Goal: learn a function from examples**

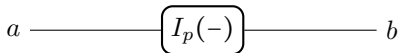Fix sets $A$, $B$. For all $f: A \to B$, use pairs $(a, f(a))$ to *approximate* $f$.

**Method: use the following data**

Hypothesis set: $P \leftarrow$ Strategies

Implementation function: $I: P \times A \to B \leftarrow$ Play

Update function $U: P \times A \times B \to P \leftarrow$ Equilibrium

Request function $r: P \times A \times B \to A \leftarrow$ Coutility

$$a \;\overline{\qquad\boxed{I_p(-)}\qquad}\; b$$

A **learner** $A \to B$ is a tuple* $(P, I, U, r)$.

---

*actually an equivalence class.

**Goal: learn a function from examples**

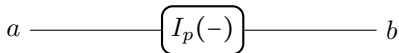Fix sets $A$, $B$. For all $f: A \to B$, use pairs $(a, f(a))$ to *approximate* $f$.

**Method: use the following data**

    Hypothesis set: $P$

    Implementation function: $I: P \times A \to B$

    Update function $U: P \times A \times B \to P$

    Request function $r: P \times A \times B \to A$

$$a \quad\rule{3cm}{0.4pt}\quad \boxed{I_p(-)} \quad\rule{3cm}{0.4pt}\quad b$$

A **learner** $A \to B$ is a tuple[*] $(P, I, U, r)$.

---

[*]actually an equivalence class.

The symmetric monoidal category Learn has
>     **objects**: sets
>     **morphisms**: learners $(P, I, U, r)$.

*How does **composition** work? Suppose we have a pair of learners:*

$$A \xrightarrow{(P,I,U,r)} B \xrightarrow{(Q,J,V,s)} C.$$

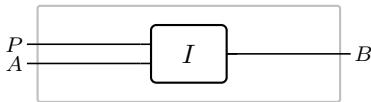*How does **composition** work? Suppose we have a pair of learners:*

$$A \xrightarrow{(P,I,U,r)} B \xrightarrow{(Q,J,V,s)} C.$$

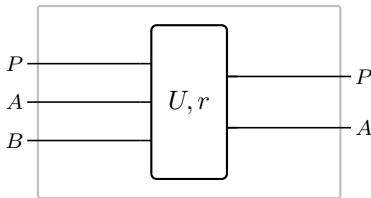The new parameter space is just the product $Q \times P$.

*How does **composition** work? Suppose we have a pair of learners:*

$$A \xrightarrow{(P,I,U,r)} B \xrightarrow{(Q,J,V,s)} C.$$

Let's represent our learners with string diagrams:
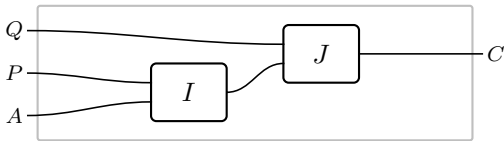


$$I \colon P \times A \longrightarrow B$$



$$(U,r) \colon P \times A \times B \longrightarrow P \times A$$

*How does **composition** work? Suppose we have a pair of learners:*

$$A \xrightarrow{(P,I,U,r)} B \xrightarrow{(Q,J,V,s)} C.$$

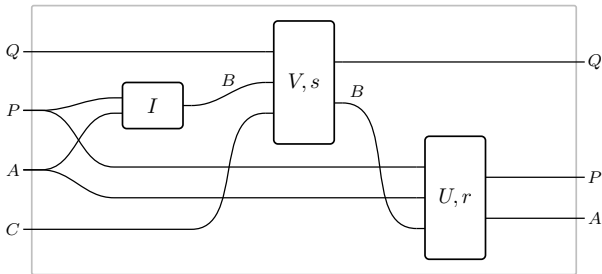Composing implementation functions is straightforward:



$$(q, p, a) \longmapsto J(q, I(p, a))$$

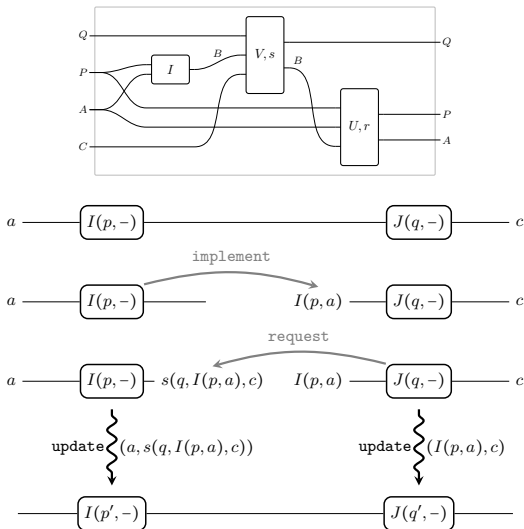*How does **composition** work? Suppose we have a pair of learners:*

$$A \xrightarrow{(P,I,U,r)} B \xrightarrow{(Q,J,V,s)} C.$$

Composing update/request functions is more complicated:



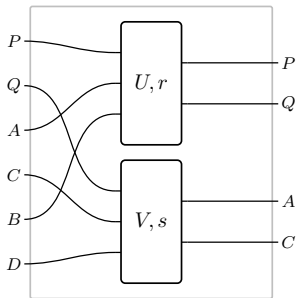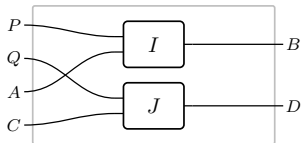$$(q,p,a,c) \longmapsto \Big( V\big(q, I(p,a), c\big), U\big(p, a, s(q, I(p,a), c)\big), r\big(p, a, s(q, I(p,a), c)\big)\Big).$$

*Key idea: composition creates local training data.*

The **monoidal product** of $(P, I, U, r): A \to B$ and $(Q, J, V, s): C \to D$ is given by

A compositional framework for supervised learning:

    **Learning**: parameter updates.

    **Supervised**: training is by (input, output) pairs.

    **Compositional**: we can build new learners from old.

A compositional framework for supervised learning:

**Learning**: parameter updates.

**Supervised**: training is by (input, output) pairs.

**Compositional**: we can build new learners from old.

But how can we explicitly construct a learner?

# II. Specifying Parametrised Functions
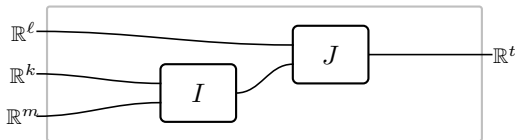
The prop Para has
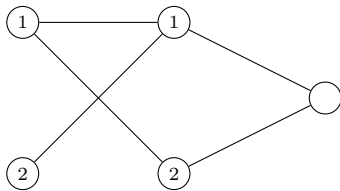
> **objects**: natural numbers
>
> **morphisms** $m \to n$: *differentiable* functions
>
> $$I : \mathbb{R}^k \times \mathbb{R}^m \to \mathbb{R}^n.$$

**Composition** is as for implementation functions in Learn:

*Neural networks (sequences of bipartite graphs) are a compositional, combinatorial language for specifying differentiable parametrised functions.*



$I\colon (\mathbb{R}^5 \times \mathbb{R}^3) \times \mathbb{R}^2 \longrightarrow \mathbb{R};$

$$(p, q, a) \longmapsto \sigma\big(q_1\sigma(p_{11}a_1 + p_{12}a_2 + p_{1b}) + q_2\sigma(p_{21}a_1 + p_{2b}) + q_b\big).$$

where $\sigma\colon \mathbb{R} \to \mathbb{R}$ is a differentiable function known as the activation.

The prop NNet has

> **objects**: natural numbers.
>
> **morphisms** $m \to n$: neural networks with $m$ inputs and $n$ outputs.
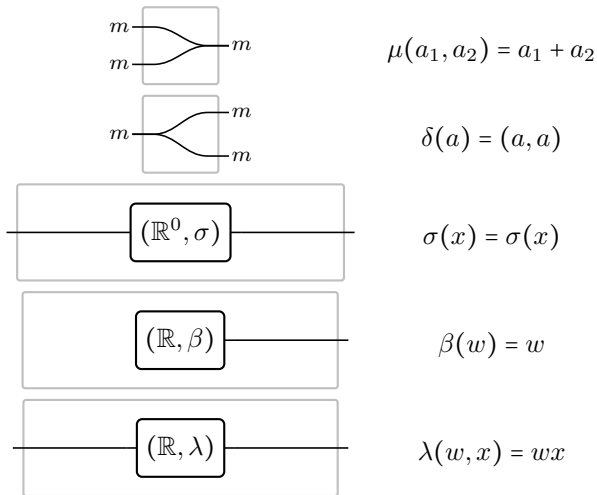>
> **composition**: concatenation of neural networks.

**Theorem**
A differentiable function $\sigma: \mathbb{R} \to \mathbb{R}$ defines a prop functor

$$I_\sigma: \mathsf{NNet} \longrightarrow \mathsf{Para}.$$

*Differentiable parametrised functions can also be constructed using string diagrams in* Para.

The image of NNet under $I_\sigma$ is contained in the composite of:

 $\mu(a_1, a_2) = a_1 + a_2$

 $\delta(a) = (a, a)$

 $\sigma(x) = \sigma(x)$

 $\beta(w) = w$
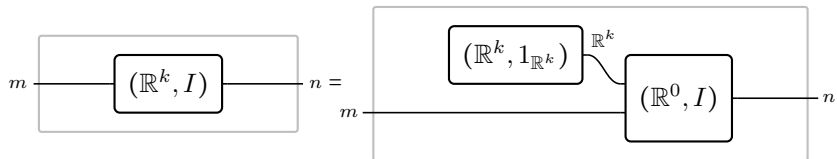
 $\lambda(w, x) = wx$

*Differentiable parametrised functions can also be constructed using string diagrams in* Para.
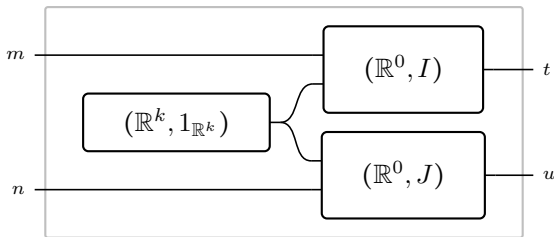
*Weight-tying is a technique that identifies parameters that describe the same structure.*

We factorise.



Then copy.

# III. Backprop: Updates and Requests via Gradient Descent

**Theorem**
Fix $\epsilon > 0$, $e \colon \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ such that $\frac{\partial e}{\partial x}(x_0, -) \colon \mathbb{R} \to \mathbb{R}$ has inverse $h_{x_0}$ for each $x_0$.

There is a faithful, injective-on-objects, strong symmetric monoidal functor

$$L_{\epsilon,e} \colon \mathsf{Para} \longrightarrow \mathsf{Learn}$$

sending each object $m$ to $\mathbb{R}^m$, and each morphism $(\mathbb{R}^k, I) \colon m \to n$ to the learner $(\mathbb{R}^k, I, U_I, r_I) \colon \mathbb{R}^m \to \mathbb{R}^n$ defined by

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

$$r_I(p, a, b) = h_a \Big( \nabla_a E_I(p, a, b) \Big),$$

Here $E_I(p, a, b) = \sum_i e(I(p, a)_i, b_i)$ and $h_a$ denotes component-wise application of $h_{a_i}$.

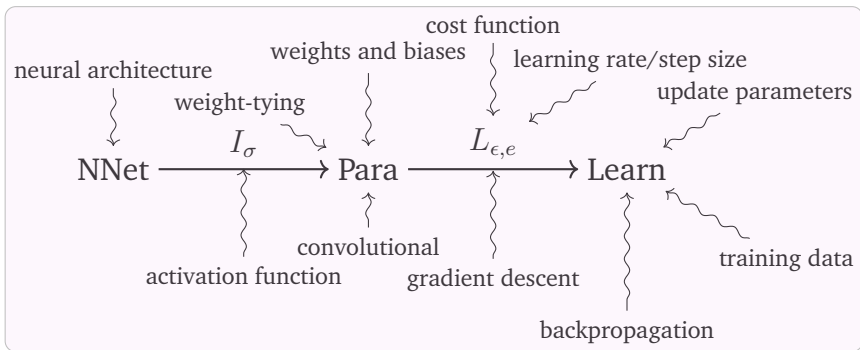Let $e$ be the *quadratic error* $\mathrm{quad}(x,y) = \frac{1}{2}(x-y)^2$.

**Corollary**
For every $\epsilon > 0$, there is a strong symmetric monoidal functor

$$L_{\epsilon,\mathsf{quad}}\colon \mathsf{Para} \longrightarrow \mathsf{Learn}$$

sending $(\mathbb{R}^k, I)\colon m \to n$ to the learner $(\mathbb{R}^k, I, U_I, r_I)\colon \mathbb{R}^m \to \mathbb{R}^n$ defined by

$$U_I(p,a,b)_k = p_k - \epsilon \sum_j (I_j(p,a) - b_j)\frac{\partial I_j}{\partial p_k}$$

$$r_I(p,a,b)_i = a_i - \sum_j (I_j(p,a) - b_j)\frac{\partial I_j}{\partial a_i}.$$

$$\text{NNet} \xrightarrow{I_\sigma} \text{Para} \xrightarrow{L_{\epsilon,e}} \text{Learn}$$

neural architecture

weight-tying

weights and biases

cost function

learning rate/step size

update parameters

activation function

convolutional

gradient descent

backpropagation

training data

# IV. Learners, Lenses, and Open Games

An **asymmetric lens** $(p, g): A \to B$ is a learner with trivial state space.

| **Learner** $A \to B$ | **Asymmetric lens** $A \to B$ |
|---|---|
| Hypotheses $P$ | — |
| Implementation $I: P \times A \to B$ | Put $p: A \to B$ |
| Update $U: P \times A \times B \to P$ | — |
| Request $r: P \times A \times B \to A$ | Get $g: A \times B \to A$ |

**Theorem**
There is a faithful, identity-on-objects symmetric monoidal functor from Learn to the category of *spans of asymmetric lenses* mapping

$$(P, I, U, r): A \to B$$

to

$$A \xleftarrow{(\pi_2, (\pi_1, \pi_3))} P \times A \xrightarrow{(I, (U, r))} B.$$

| **Learner** $A \to B$ | **Open game** $(X, S) \to (Y, R)$ |
|:---:|:---:|
| Hypotheses $P$ | Strategy profiles $\Sigma$ |
| Implementation $I \colon P \times A \to B$ | Play $P \colon X \times \Sigma \to Y$ |
| Update $U \colon P \times A \times B \to P$ | Equilibrium $E \colon X \times (Y \to R) \to \mathcal{P}\Sigma$ |
| Request $r \colon P \times A \times B \to A$ | Coutility $C \colon X \times \Sigma \times R \to S$ |

**Summary**

**For more:**

https://arxiv.org/abs/1711.10455

http://www.brendanfong.com/