# CATEGORICAL FOUNDATIONS OF GRADIENT-BASED LEARNING

(CRUTTWELL, GAVRANOVIĆ, GHANI, WILSON, ZANASI)
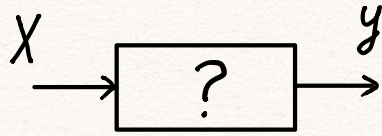
GOAL:
PROVIDE A CATEGORICAL FRAMEWORK
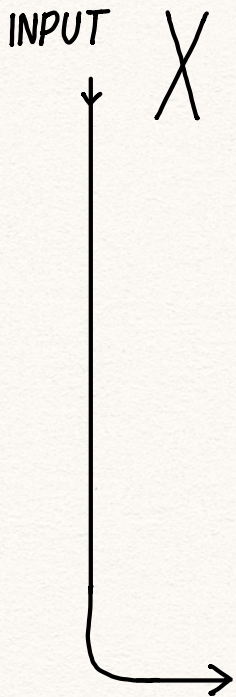
FOR DEEP LEARNING
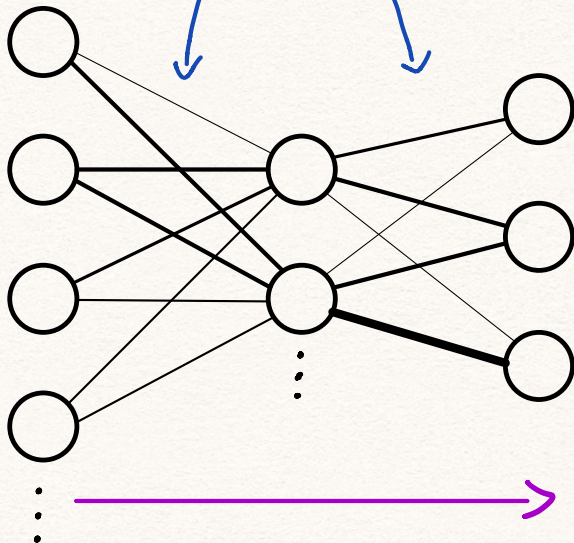
# SUPERVISED LEARNING WITH NEURAL NETWORKS
## IN ONE SLIDE:

$$X \rightarrow \boxed{?} \rightarrow y$$

DATASET : List $X \times Y$



INPUT $X$

NEURAL NETWORK
WEIGHTS

1 CAT
+
0 DOG
+
0 HORSE

LABEL $y$

PREDICTION $y$

0.4 CAT
+
0.5 DOG
+
0.1 HORSE

LOSS $\mathbb{R}$

GRADIENT DESCENT
~
„OPTIMIZER"

- NN IS COMPUTATION PARAMETERIZED BY WEIGHTS

- BACKPROPAGATION OF CHANGES

- PARAMETER UPDATE — „OPTIMIZERS"

THIS SIMPLE STORY PERMEATES DEEP LEARNING!

# PLAN FOR TODAY?

TAKE A BIRD'S EYE VIEW OF NEURAL NETWORKS

- TRACE OUT THE INFORMATION FLOW ABOVE

- PRECISELY WRITE DOWN ALL THE HIGH-LEVEL NOTIONS IN ISOLATION:
  - DIFFERENTIATION        - REVERSE DERIVATIVE CATS.
  - BIDIRECTIONALITY        - OPTICS/LENSES
  - PARAMETERIZATION       - PARA

AND STUDY THEIR INTERACTION.

## PARAMETERIZED OPTICS
AS A COMMON STRUCTURE BEHIND
- NEURAL NETWORKS
- LOSS FUNCTIONS
- OPTIMIZERS

- PAUL: CONCRETE EXAMPLES OF NEURAL NETWORKS

# DIFFERENTIATION

- CARTESIAN (FORWARD) DIFFERENTIAL CATEGORIES
(Blute et. al.)
- CARTESIAN REVERSE DIFFERENTIAL CATEGORIES (CRDC)
(Cockett et. al.)

## DEFINITION.

A CRDC $\mathcal{C}$ is a Cartesian left-additive category which for every map

$$f: A \longrightarrow B$$

has a REVERSE DIFFERENTIAL COMBINATOR

$$R[f]: A \times B \longrightarrow A$$

$\left(\text{compare } D[f]: A \times A \longrightarrow B\right)$

subject to 7 axioms.

## EXAMPLE. Smooth is a CRDC. $Poly_{\mathbb{Z}_2}$ IS A RDC

## EXAMPLE. Let $\mathbb{R}^2 \xrightarrow{f} \mathbb{R}$ in Smooth.
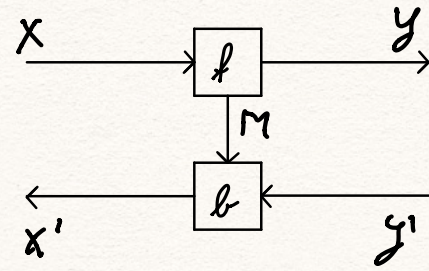$$(x,y) \longmapsto x^2 + 3yx$$

Then $R[f]: \mathbb{R}^2 \times \mathbb{R} \longrightarrow \mathbb{R}^2$
$$(x,y), w \longmapsto (2xw, 3xw)$$

PLAN: STUDY CRDC'S THROUGH OPTICS/LENSES

# OPTICS/LENSES

**DEFINITION.** Let $\mathcal{C}$ be a SMC. Category $\text{Optic}(\mathcal{C})$:

- Objects — pairs of objects $\binom{X}{X'}$ in $\mathcal{C}$

- $\text{Optic}(\mathcal{C})\left(\binom{X}{X'}, \binom{Y}{Y'}\right) = \int^{M:\mathcal{C}} \mathcal{C}(X, Y \otimes M) \times \mathcal{C}(Y' \otimes M, X')$  ← COEND

$$(M, f, b) \qquad f: X \longrightarrow Y \otimes M$$
$$b: Y' \otimes M \longrightarrow X'$$

**PROP.** If $\mathcal{C}$ is Cartesian,

$$\int^{M:\mathcal{C}} \mathcal{C}(X, M \times Y) \times \mathcal{C}(M \times Y', X')$$
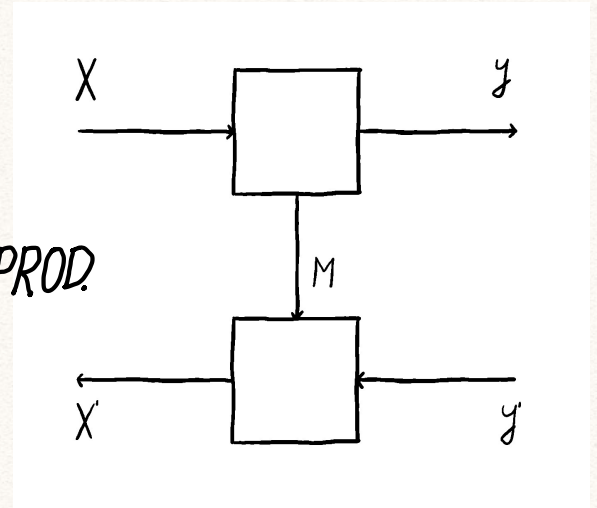$$\cong \text{ UNIV. PROPERTY OF PROD.}$$
$$\int^{M:\mathcal{C}} \mathcal{C}(X, Y) \times \mathcal{C}(X, M) \times \mathcal{C}(M \times Y', X')$$
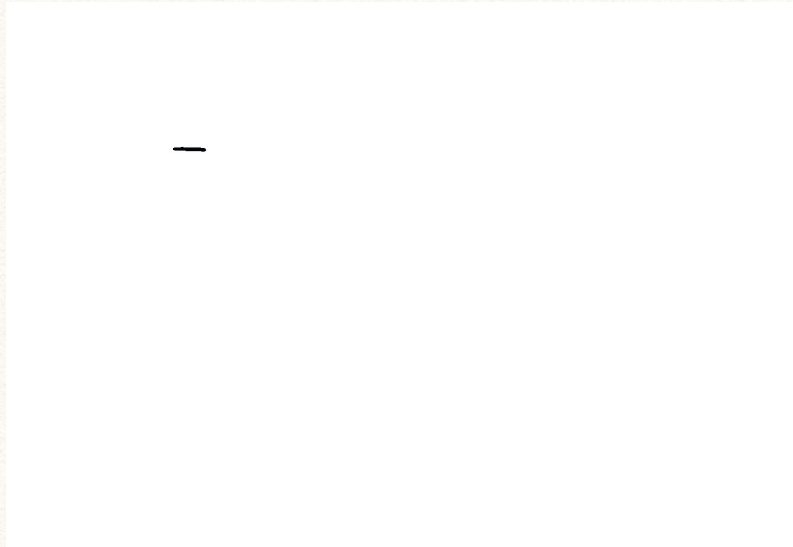$$\cong \text{ YONEDA REDUCTION}$$
$$\int \mathcal{C}(X, Y) \times \mathcal{C}(X \times Y', X')$$

$\underset{\text{get}}{\phantom{\int}} \qquad \underset{\text{put}}{\phantom{\int}}$
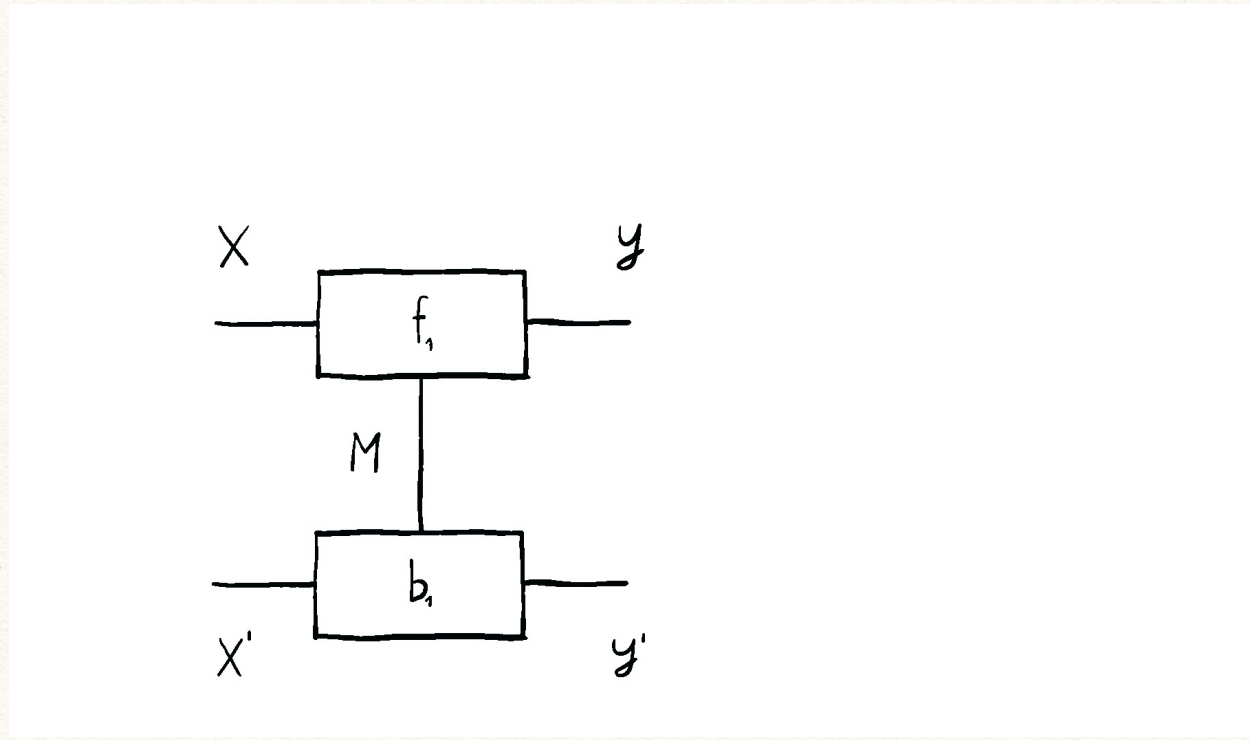
then $\text{Optic}(\mathcal{C}) \cong \text{Lens}(\mathcal{C})$
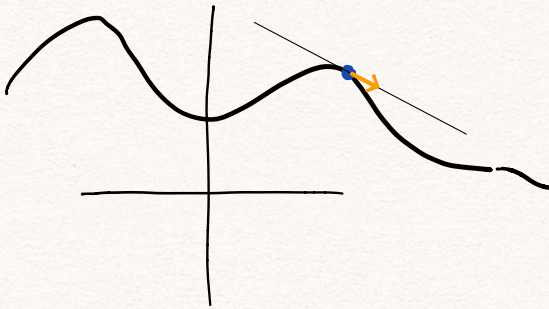
BIDIRECTIONAL INFORMATION FLOW

—

# OPTICS CAN BE COMPOSED



**PROPOSITION.** Optic($\mathcal{C}$) is symmetric monoidal.

# EXAMPLE.

## GRADIENT DESCENT



$$P \times P' \xrightarrow{\;u\;} P$$
$$(p, \nabla p) \longmapsto p - \alpha \nabla p$$

is a lens, for $\mathcal{C} := \text{Smooth}$

$$\binom{P}{P} \xrightarrow{(id_p, u)} \binom{P}{P'}$$



-It's a surprise tool that will help us later

# EXAMPLE.  STATEFUL OPTIMIZERS

$$\binom{S \times P}{S \times P} \longrightarrow \binom{P}{P'}$$

- ## MOMENTUM,

$$\text{get}: P \times P \longrightarrow P$$
$$(\nu, p) \longmapsto p$$

$$\text{put}: P \times P \times P \longrightarrow P \times P$$
$$(\nu, p, \nabla p) \longmapsto (\nu', p - \nu')$$
$$\text{where } \nu' = \gamma \nu + \varepsilon p'$$

- ## NESTEROV  MOMENTUM

$$\text{get}: P \times P \longrightarrow P$$
$$(\nu, p) \longmapsto p - \gamma \nu$$

put - same as above

- ## ADAGRAD

- ## ADAM
  . . .

# BACK TO CRDC's:

$$f: A \longrightarrow B$$
$$R[f]: A \times B \longrightarrow A$$

$\sim$ '.get' MAP OF A LENS
$\sim$ '.put' MAP OF A LENS

## <span style="color:blue">PROPOSITION.</span>

For each CRDC $\mathcal{C}$ there is a symmetric monoidal functor

$$\mathcal{C} \xrightarrow{\quad F \quad} Lens(\mathcal{C}) \cong Optic(\mathcal{C})$$

$$
\begin{array}{ccc}
A & \longmapsto & (A, A) \\
\downarrow f & & \downarrow (f, R[f]) \\
B & \longmapsto & (B, B)
\end{array}
$$

· THIS IS OUR FRAMEWORK FOR BACKPROPAGATION

# PARAMETERIZATION

Fix a SMC $(\mathcal{C}, \otimes, I)$.

**DEF.** Bicategory $\text{Para}(\mathcal{C})$

Objects – objects of $\mathcal{C}$

CATEGORY
OF ELEMENTS

$$\text{Para}(\mathcal{C})(A,B) = \int^{P:\mathcal{C}} \mathcal{C}(P \otimes A, B)$$

$$A \xrightarrow{(P:\mathcal{C},\; f:P \otimes A \to B)} B$$

2-cells are reparameterizations: a 2-cell

is a map $Q \xrightarrow{r} P$ such that

$$A \; \underset{(Q,g)}{\overset{(P,f)}{\rightrightarrows}} \; B \quad \Downarrow r$$

$$Q \otimes A \xrightarrow{r \otimes A} P \otimes A$$
$$\searrow{g} \qquad \swarrow{f}$$
$$B$$

$(\text{Set}, \times, 1)$

$\text{Para}(\text{Set})$

SETS AND
PARAMETERIZED FUNCTIONS

$(\text{Smooth}, \times, 1)$

$\text{Para}(\text{Smooth})$

EUCLIDEAN SPACES AND
PARAMETERIZED SMOOTH
FUNCTIONS

$(\text{Optic}(\mathcal{C}), \otimes, 1)$

$\text{Para}(\text{Optic}(\mathcal{C}))$

PAIRS OF OBJECTS AND
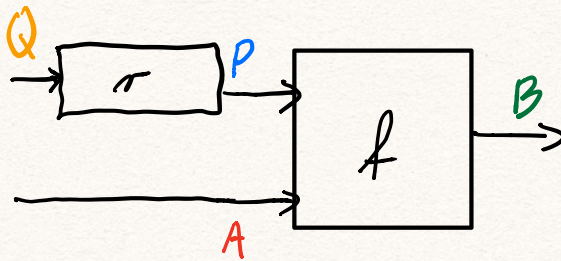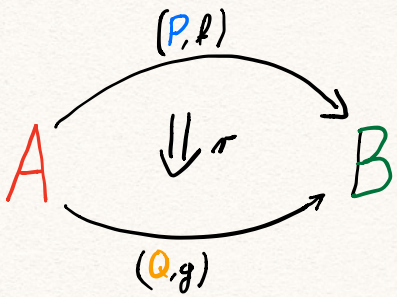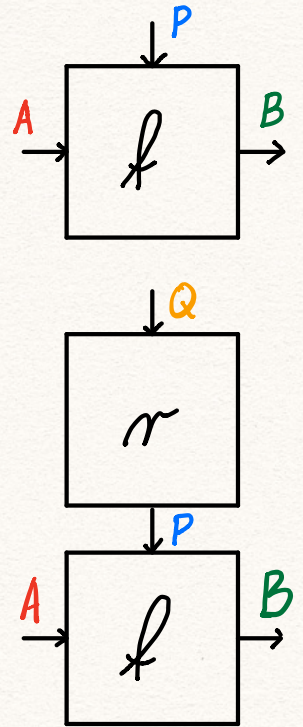PARAMETERIZED OPTICS

. . .

# GRAPHICAL LANGUAGE

| TEXTUAL NOTATION | STANDARD STRING DIAGRAM | 2D STRING DIAGRAM |
|---|---|---|

$$f: P \otimes A \longrightarrow B$$



$(P, f)$

$A \overset{(P,f)}{\underset{(Q,g)}{\Longrightarrow}} B \quad r$

# HOW DOES COMPOSITION WORK?

# RECAP



$w_1$     $w_2$

$w_1$     $w_2$

$\mathbb{R}^{4\times3}$     $\mathbb{R}^{2\times3}$

$\iff$

$\mathbb{R}^4 \rightarrow \boxed{f} \xrightarrow{\mathbb{R}^2} \boxed{g} \xrightarrow{\mathbb{R}^3}$

Para IS NATURAL WITH RESPECT TO BASE CHANGE.

# DEFINITION.

Let $G: \mathcal{C} \longrightarrow \mathcal{D}$ be a symm. monoidal functor. We define

$$\text{Para}(G): \text{Para}(\mathcal{C}) \longrightarrow \text{Para}(\mathcal{D})$$



$$
\begin{array}{ccc}
A & \longmapsto & GA \\
\downarrow {\scriptstyle (P,f)} & & \downarrow {\scriptstyle (GP, f')} \\
B & \longmapsto & GB
\end{array}
$$

where $f'$ is the composite

$$G(P) \otimes G(A) \qquad\qquad\qquad G(B)$$

+ MORE.
  Para IS RICH IN CATEGORICAL STRUCTURE.
  - Cokleisli category of a graded comonad
  - Double category
  - Actegorical Para
    ...

# PARAMETERIZED OPTICS

$$\mathcal{C} \longmapsto Optic(\mathcal{C}) \longmapsto Para(Optic(\mathcal{C}))$$

· Objects – objects of $Optic(\mathcal{C})$ – pairs $\begin{pmatrix} X \\ X' \end{pmatrix}$ in $\mathcal{C}$

· Morphisms $\begin{pmatrix} X \\ X' \end{pmatrix} \xrightarrow{\left( \binom{P}{P'}, f \right)} \begin{pmatrix} Y \\ Y' \end{pmatrix}$ where $f: \begin{pmatrix} P \otimes X \\ P' \otimes X' \end{pmatrix} \longrightarrow \begin{pmatrix} Y \\ Y' \end{pmatrix}$

$(M, f, b)$



· WE CAN COMPOSE PARAMETERIZED OPTICS

- We automatically get two parameter ports

- A 2-cell $\begin{pmatrix} X \\ S \end{pmatrix} \overset{\begin{pmatrix} P \\ Q, f \end{pmatrix}}{\underset{\begin{pmatrix} Z \\ W, g \end{pmatrix}}{\Longrightarrow r}} \begin{pmatrix} Y \\ R \end{pmatrix}$ is an <u>optic</u>

$$\begin{pmatrix} Z \\ W \end{pmatrix} \overset{r}{\longrightarrow} \begin{pmatrix} P \\ Q \end{pmatrix}$$

<u>THEOREM.</u>

GRADIENT DESCENT IS A 2-cell IN $\mathrm{Para}(\mathrm{Optic}(\mathcal{C}))$.

(Since it is a lens)

# THEOREM.

APPLYING $Para$ TO THE CRDC FUNCTOR

$$\mathcal{C} \xrightarrow{\ F\ } Optic(\mathcal{C})$$

RESULTS IN A FUNCTOR

$$Para(\mathcal{C}) \xrightarrow{\ Para(F)\ } Para(Optic(\mathcal{C}))$$



- FUNCTORIALITY IS IMPORTANT!

# EXAMPLE. A NEURAL NETWORK + A LOSS FUNCTION

WE CAN PUT THE PIECES TOGETHER.

# SUPERVISED LEARNING

# Categorical Foundations of Gradient-Based Learning

# Recap

So far[1]:

▶ Para and Lens
▶ Optimizers, loss functions, models all (parametrised) lenses
▶ Putting them together, we get this:



Now: make these boxes more transparent...

[1]Cruttwell et al., "Categorical Foundations of Gradient-Based Learning."

# Next up

▶ Theme: **How to Build a Neural Network out of Lenses**
▶ Choosing the model is a creative process
▶ For an example problem, we'll look at the structure of one choice of model

Two goals:

▶ Show how to build a simple neural network out of lenses
▶ How to replace "classical" picture of neural networks with string diagrams

# To String Diagrams

# Three Levels of Detail 1: Learning

▶ The most "zoomed-out" view is the learner
▶ We look at the model as a kind of black box

# Three Levels of Detail 2: Model Architecture

▶ The high-level structure of the model as a composition of "layers"

▶ Think of layers[2] as subroutines

▶ DL literature already starting to look string-diagrammatic[3]



---

[2] Ambiguous terminology warning: "Layer" conflates objects and morphisms
[3] Kaiser et al., "One Model to Learn Them All."

# Three Levels of Detail 3: Layer

Finally: what are the pair of maps in our base category that make up a layer?



What's in here?

# This Section of the Talk

- ▶ Supervised Learning & Reverse derivatives[4]
- ▶ End-To-End Example of a Neural Network
- ▶ Other Layer Examples
  - ▶ Weight Tying
  - ▶ Convolutional Layers
- ▶ Other settings (Circuits and POLY$_{\mathbf{Z}_2}$)

---

[4]Cockett et al., "Reverse Derivative Categories."

## Supervised Learning

In supervised learning, we want to learn a map

$$f : A \to B$$

from a dataset of examples

$$(a, b) \in A \times B$$

Now, based on our beliefs about the structure of $A$ and $B$, we design a *parametrised* map:

$$\text{model} : P \times A \to B$$

and we search for some $\theta \in P$ such that $\text{model}(\theta, -)$ best represents the data.

## Gradient-Based Learning

We want to use a datapoint $(a, b) \in A \times B$ to improve $\theta$, so we need a map

$$??? : P \times A \times B \to P$$

The reverse derivative is almost what we want. For a map $f : A \to B$,

$$R[f] : A \times B' \to A'$$

(while in an RDC $A' = A$ and $B' = B$, it's useful think of the "primed" objects as representing **changes**)

So the reverse derivative of our model morphism has the following type:

$$R[\text{model}] : P \times A \times B' \to P' \times A'$$

# Updates, "Displacement" and Reverse Derivatives

This is not quite enough: we have two problems:

1. We have a "true" value $b \in B$ and a "predicted" value $\mathrm{model}(\theta, a) \in B$ but we need a $B'$
2. The reverse derivative gives us a $P'$ and we want a $P$

This is exactly what the update and loss lenses are for:

$$\mathrm{update}_{\mathrm{put}} : P \times P' \to P$$

$$\mathrm{loss}_{\mathrm{put}} : B \times B \to B' \times B'$$

$$R[\mathrm{model}] : P \times A \times B' \to P' \times A'$$

# Reverse Derivatives, Graphically

CARTESIAN STRUCTURE

LEFT-ADDITIVE STRUCTURE

Copy

discard

add

zero

$x \mapsto \langle x, x \rangle$

$x \mapsto \langle \rangle$

$x_1, x_2 \mapsto \langle x_1 + x_2 \rangle$

$\langle \rangle \mapsto \langle 0 \rangle$

Addition & zero maps

$f + g : A \to B$

$0 : A \to B$

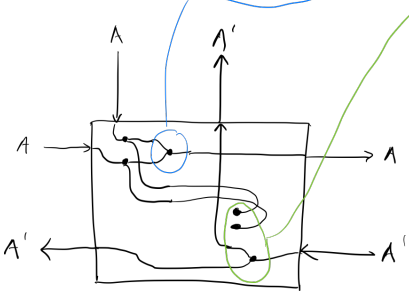# Reverse Derivatives, Graphically

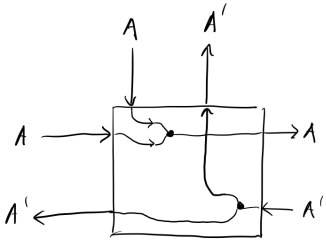# Reverse Derivatives, Graphically

# Neural Networks 1: Dense Layers

Now let's unpack a dense layer...

# Neural Networks 2: Bias "Layer"

# Neural Networks 3: 'Linear' Layer

▶ Parameters $P = \mathbb{R}^{b \cdot a}$ are the coefficients of a matrix
▶ Input $A = \mathbb{R}^a$ is an $a$-dimensional vector
▶ Forward pass multiplies the matrix by the vector:
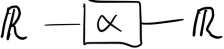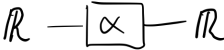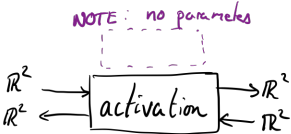
$$\mathsf{get}(M, x) \mapsto Mx$$

▶ Reverse pass does the "obvious" thing that typechecks: if we think of the get map as having the type

$$\mathsf{get} : \mathsf{Mat}(A, B) \times \mathsf{Vec}(A) \to \mathsf{Vec}(B)$$
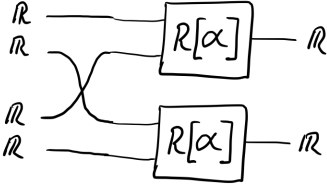
Then the codomain of the put map should be $\mathsf{Mat}(A, B) \times \mathsf{Vec}(A)$:

$$\mathsf{put}(M, x, y) \mapsto \langle y \otimes x, M^T y \rangle$$
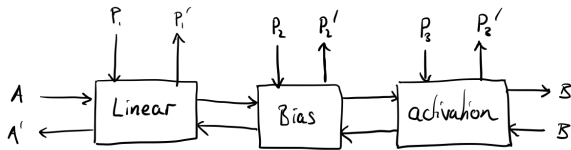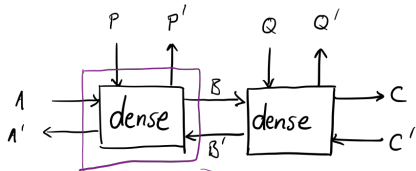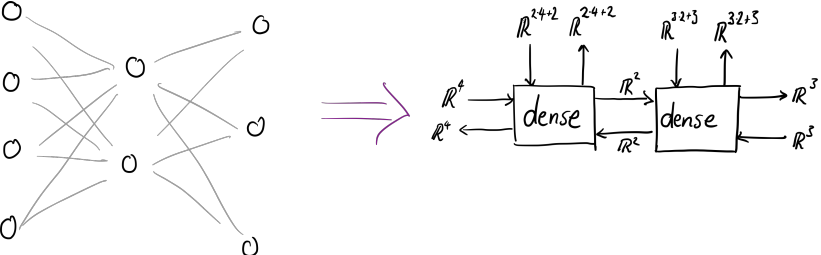
# Neural Networks 4: Activation Layer
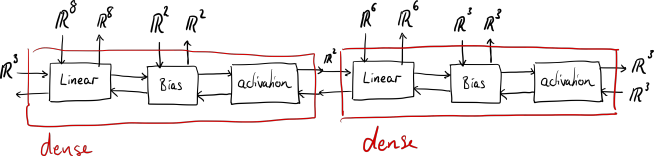
# Neural Networks 5: Dense Layers (again)

# Neural Networks 6: Hidden Layer Neural Network
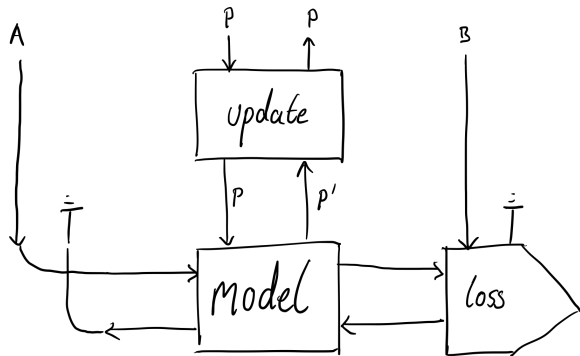
Returning to the "standard" picture of a neural network:
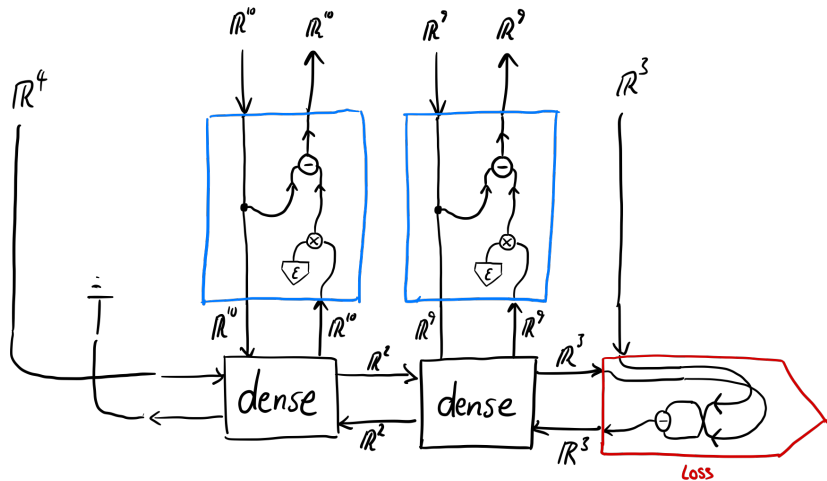


Expanding out "dense":

# Update & Loss

Now let's substitute all parts into the full picture

# Full Picture (Again)

# Code

- ► Code implementing these ideas can be found here:
  https://github.com/statusfailed/numeric-optics-python/
- ► Includes this hidden layer neural network model
- ► Also includes a convolutional model for the MNIST dataset
  (more on this shortly…)

# More

- ▶ Other Layer Examples
    - ▶ Weight Tying
    - ▶ Convolutional Layers
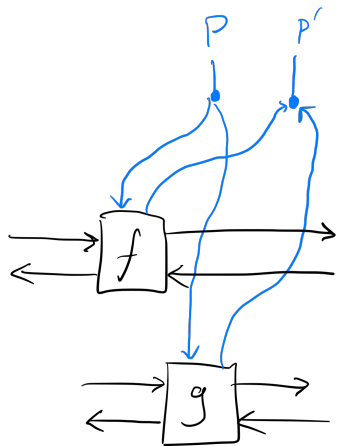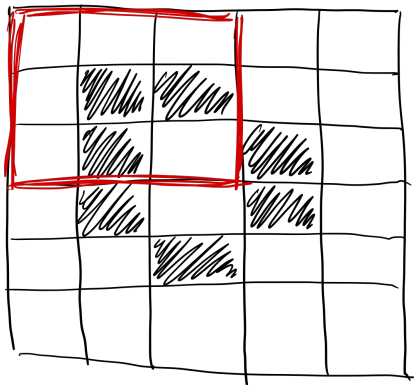- ▶ Other settings (Circuits and $\text{POLY}_{\mathbf{Z}_2}$)

"Weight Tying"

# Image Processing

▶ Example problem: image processing, e.g. digit recognition
▶ Convolution layer: features with spatial locality

# Convolutional Layers

# Other Settings: POLY$_{\mathbf{Z}_2}$

▶ POLY$_{\mathbf{Z}_2}$ is an RDC
▶ We can still think of morphisms as functions
▶ Gradient-based learning still works[5]
▶ Strange possibilities for layers: the LUT

---

[5]Wilson and Zanasi, "Reverse Derivative Ascent."

# References

Cockett, Robin, Geoffrey Cruttwell, Jonathan Gallagher, Jean-Simon Pacaud Lemay, Benjamin MacAdam, Gordon Plotkin, and Dorette Pronk. "Reverse Derivative Categories," 2019. http://arxiv.org/abs/1910.07065.

Cruttwell, G. S. H., Bruno Gavranović, Neil Ghani, Paul Wilson, and Fabio Zanasi. "Categorical Foundations of Gradient-Based Learning," 2021. http://arxiv.org/abs/2103.01931.

Kaiser, Lukasz, Aidan N. Gomez, Noam Shazeer, Ashish Vaswani, Niki Parmar, Llion Jones, and Jakob Uszkoreit. "One Model to Learn Them All," 2017. http://arxiv.org/abs/1706.05137.

Wilson, Paul, and Fabio Zanasi. "Reverse Derivative Ascent: A Categorical Approach to Learning Boolean Circuits." *Electronic Proceedings in Theoretical Computer Science* 333 (February 2021): 247–60. https://doi.org/10.4204/eptcs.333.17.