# Gradient Descent & Neural Networks

## A Categorical Approach

Following Fong, Spivak, Tuyéras, Capucci, Gavranović, Hedges, Rischel …

I: Gradient Descent
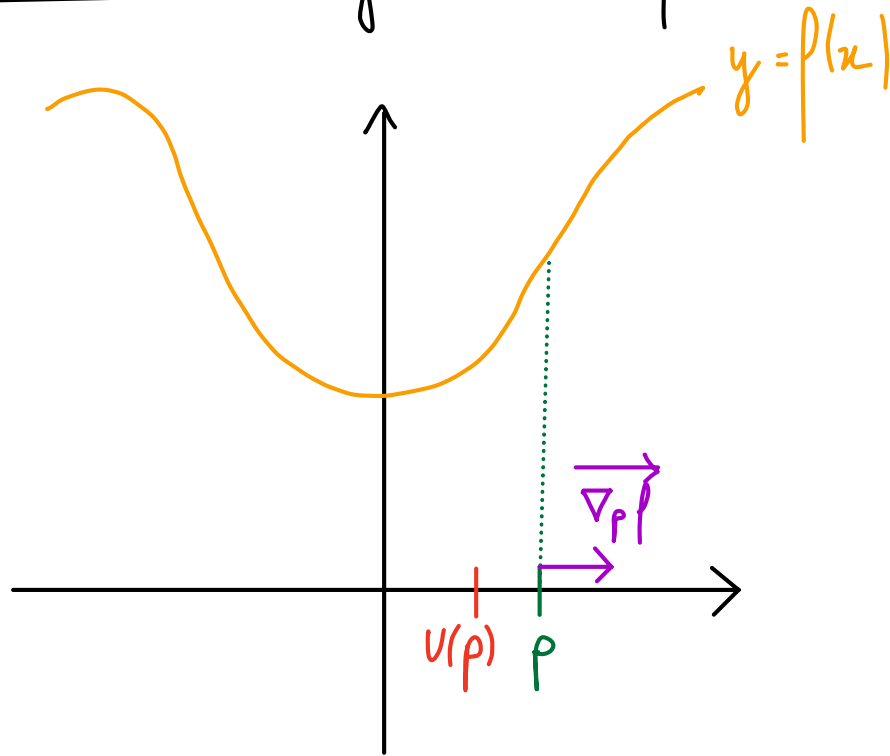
II: Categorical model: Para, Optic,...

III: Naïve Neural Networks

IV: Categorical Neural Networks

V: Training Your Network

# I Gradient descent: algorithm to find local minima

$$y = f(x)$$

$$\overrightarrow{\nabla_p f}$$

$$U(p) \quad p$$

$$U(p) = p - \varepsilon \overrightarrow{\nabla_p f}$$

$f$: total error/cost function

$P$: "First guess" parameter

$\overrightarrow{\nabla_p f}$: gradient at $p$.

$\varepsilon$: Step / learning rate

$U$: update function.

# GD for Function Approximation



$(a, b)$ : sample data

1) A shape of function: $2^{nd}$ degree polynomial

$$I(\beta_1, \beta_2, \beta_3, x) = \beta_1 + \beta_2 x + \beta_3 x^2$$

2) An error (distance) function:

$$e(x, y) = \frac{1}{2}(x - y)^2$$

Total error: $E(\beta, a, b) = \frac{1}{2}(I(\beta, a) - b)^2$

3) A step: $\varepsilon > 0$.

Get:

1) An update function:

$$U(p, a, b) = p - \varepsilon \nabla_p E(p, a, b)$$

2) A request function:

$$r(p, a, b) = \left(\frac{\partial e}{\partial x}\right)^{-1} \left(\nabla_a E(p, a, b)\right)$$

Instead of a map $\mathbb{R} \longrightarrow \mathbb{R}$, we have:

$$\mathbb{I}: \mathbb{R}^3 \times \mathbb{R} \longrightarrow \mathbb{R} \quad,$$

$$U: \mathbb{R}^3 \times \mathbb{R} \times \mathbb{R} \longrightarrow \mathbb{R}^3 \quad,$$

$$\Omega: \mathbb{R}^3 \times \mathbb{R} \times \mathbb{R} \longrightarrow \mathbb{R} \quad.$$

# II: Categorical Model

**Definition:** Let $\mathcal{C}$ be a category with products and terminal $T$.

The category $\text{Para}(\mathcal{C})$ has:

- Objects: same as $\mathcal{C}$

- Morphisms $A \xrightarrow{P,f} B$:
  pairs $(P \in ob(\mathcal{C}), \ P \times A \xrightarrow{f} B)$

identities: $A \xrightarrow{T, 1_A} B$

Composition:

$$A \xrightarrow{P,f} B \xrightarrow{Q,g} C$$
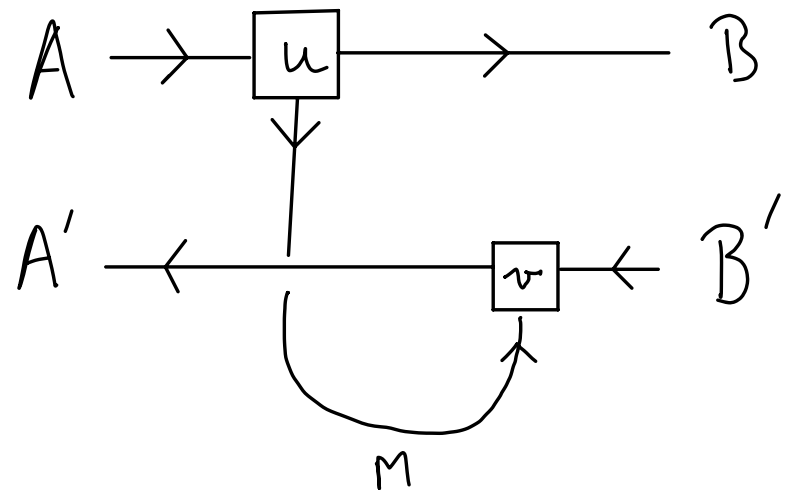$$A \xrightarrow{Q \times P, \ g \circ Q \times f} C$$

$$Q \times P \times A \xrightarrow{Q \times f} Q \times B \xrightarrow{g} C$$

Definition: Let $\mathcal{C}$ be a category with ...

The category $Optic(\mathcal{C})$ has:

- Objects: pairs $\begin{pmatrix} A \\ A' \end{pmatrix} \in ob(\mathcal{C})^2$

- Morphisms $\begin{pmatrix} A \\ A' \end{pmatrix} \begin{array}{c} \xrightarrow{f} \\ \xleftarrow{g} \end{array} \begin{pmatrix} B \\ B' \end{pmatrix}$ are classes of triples

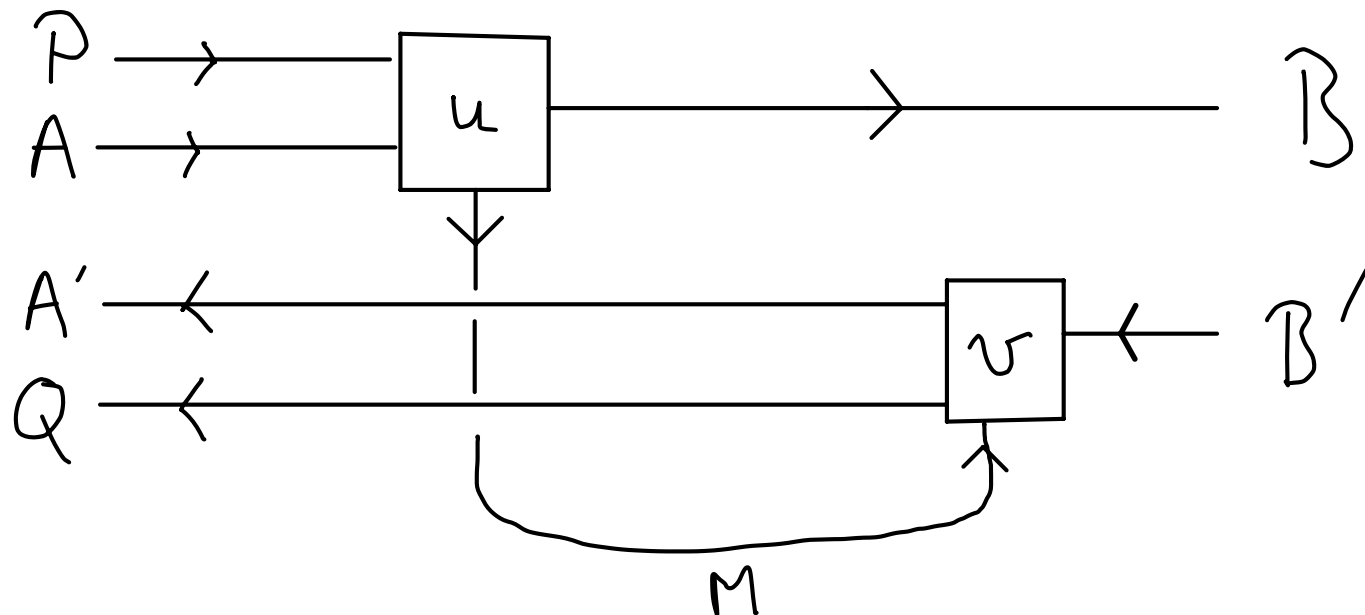$$\left( M \in ob(\mathcal{C}), \; f: A \longrightarrow B \times M, \; g: B' \times M \longrightarrow A' \right) \big/ \sim$$

$$Optic(\mathcal{C})\left( \begin{pmatrix} A \\ A' \end{pmatrix}, \begin{pmatrix} B \\ B' \end{pmatrix} \right)$$
$$:= \int^{M \in ob(\mathcal{C})} \mathcal{C}(A, M \times B) \times \mathcal{C}(M \times B', A')$$

# Parametrised Optics: $\text{Para}(\text{Optic}(\mathcal{C}))$ has

- objects: $\text{ob}(\mathcal{C})^{\times 2}$,

- morphism: $\begin{pmatrix} A \\ A' \end{pmatrix} \xrightleftharpoons[Q]{P} \begin{pmatrix} B \\ B' \end{pmatrix}$ :
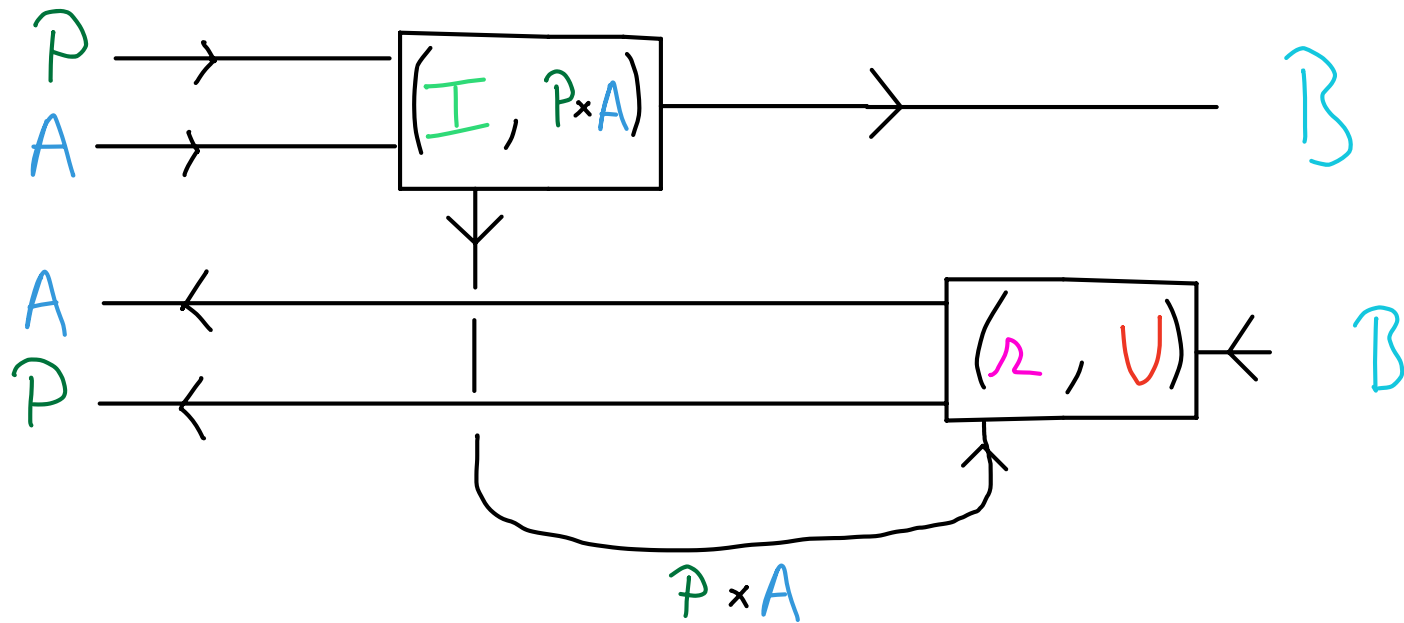
**Def:** **Smooth** has: $ob(\text{Smooth}) = \mathbb{N}$,

$$\text{Smooth}(n, m) = C^{\infty}(\mathbb{R}^n, \mathbb{R}^m)$$

$$n \times m := n + m$$

**Theorem** (Fong - Spivak - Tuyéras): Given a step $\varepsilon$ and an error function $e$, there is a Gradient Descent functor:

$$GD_{\varepsilon, e} : \text{Para}(\text{Smooth}) \longrightarrow \text{Para}(\text{Optic}(\text{Smooth}))$$

such that:

$$GD\ (A) = \begin{pmatrix} A \\ A \end{pmatrix}$$

$$GD\ \left( A \xrightarrow{P,I} B \right) = \begin{pmatrix} A \\ A \end{pmatrix} \underset{P}{\overset{P}{\rightleftarrows}} \begin{pmatrix} B \\ B \end{pmatrix} \quad \text{is represented by:}$$



$$U(p, a, b) = p - \varepsilon \nabla_p E(p, a, b) \quad \text{and} \quad r(p, a, b) = \left( \frac{\partial e}{\partial x} \right)^{-1} \left( \nabla_a E(p, a, b) \right)$$
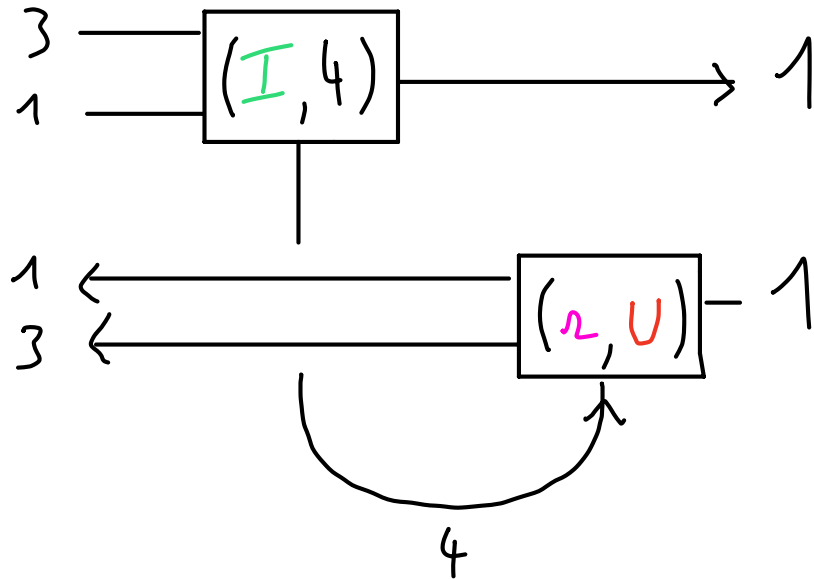
In our Example:

$$I(p_1, p_2, p_3, x) = p_1 + p_2 x + p_3 x^2$$

$$e(x, y) = \frac{1}{2}(x - y)^2$$

$$1 \xrightarrow{3, I} 1 \in Para(Smooth)$$

$$GD_{\varepsilon, e}\left(1 \xrightarrow{3, I} 1\right) =$$

# Reparametrisation

Para ($\mathcal{C}$) is (in fact) a 2-category with:

**2-cells**:

$$A \xrightarrow{P, \ell} B$$
$$\Downarrow r$$
$$A \xrightarrow{P', \ell'} B$$

$=$

$$P \times A \xrightarrow{\ell} B$$
$$\uparrow r \times A \quad \sigma \quad \nearrow$$
$$P' \times A \xrightarrow{\ell'}$$

In Para (Optic),

$$\binom{A}{B} \underset{Q}{\overset{P}{\rightleftarrows}} \binom{A'}{B'}$$

$$\Downarrow (r, s)$$

$$\binom{A}{B} \underset{Q'}{\overset{P'}{\rightleftarrows}} \binom{A'}{B'}$$

$=$

# Generalisation of GD using Differential Categories

Def: A **Cartesian Differential Category** is a cartesian

(left additive) Category $\mathcal{C}$ equipped with a **differential**

**combinator** :
$$\frac{A \xrightarrow{\ f\ } B}{A \times A \xrightarrow{\ Df\ } B} \ D$$
satisfying ...

Example: SMOOTH , with
$$\frac{n \xrightarrow{\ f\ } m}{n + n \xrightarrow{\ Df\ } m} \ D$$
$$(x, y) \longmapsto D_x f(y)$$

Remember: $r(p, a, b) = \left(\frac{\partial e}{\partial x}\right)^{-1} \left(\nabla_a E(p, a, b)\right)$

# Generalisation of GD using Differential Categories

Def: A **Reverse Differential Category** is a cartesian

(left additive) Category $\mathcal{C}$ equipped with a **differential**

**combinator**:

$$\frac{A \xrightarrow{f} B}{A \times B \xrightarrow{Rf} B} R$$

satisfying ...

Example: SMOOTH, with

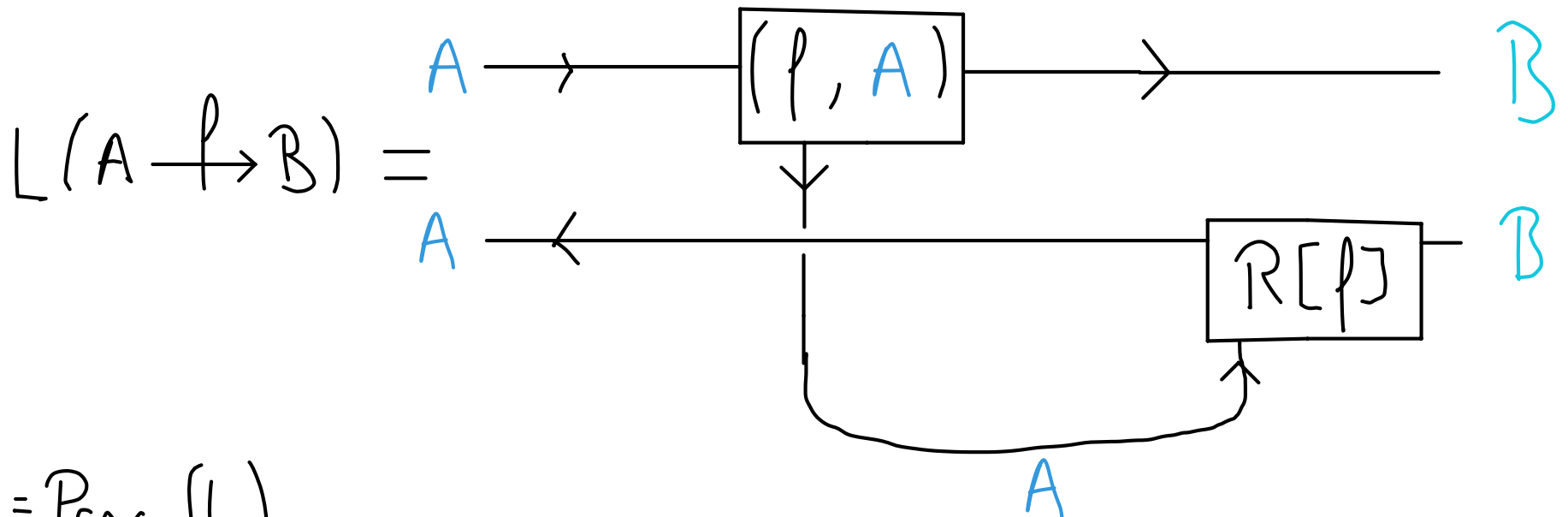$$\frac{n \xrightarrow{f} m}{n + m \xrightarrow{Rf} n} R$$

$$(x, y) \longmapsto (D_x f)^t(y)$$

# Theorem (Cockett-Cruttwell-Gallagher-Lemay-MacAdam-Plotkin-Pronk)

For any R.D.C $\mathcal{C}$, there is a functor:

$$L : \mathcal{C} \longrightarrow \text{Optic}(\mathcal{C}) \qquad \text{such that:}$$

$$L(A) = \begin{pmatrix} A \\ A \end{pmatrix}, \text{ and}$$
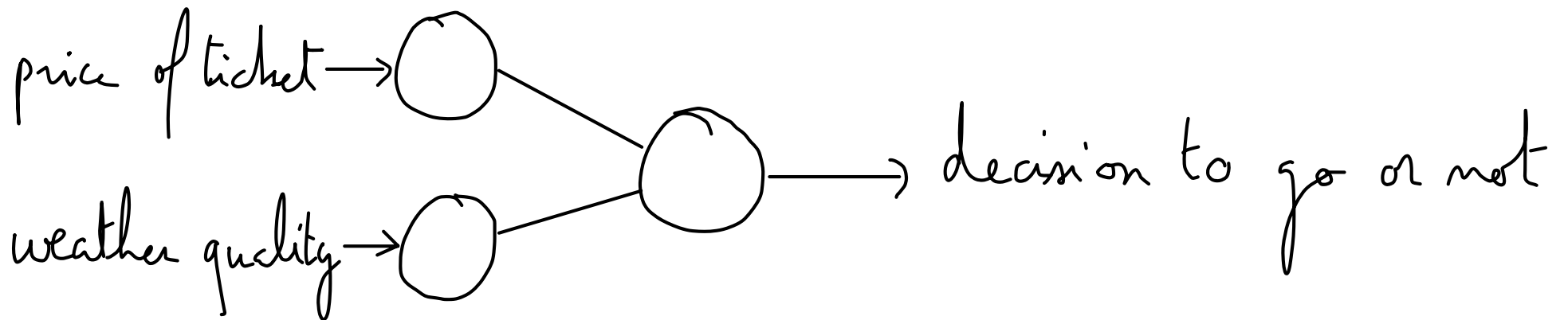
$$L(A \xrightarrow{f} B) =$$



$$GD = \text{Para}(L)$$

# III (Naive) Neural Networks

Ex: I want to go to London if

   a) The plane ticket is cheap enough

   b) The weather is good enough.
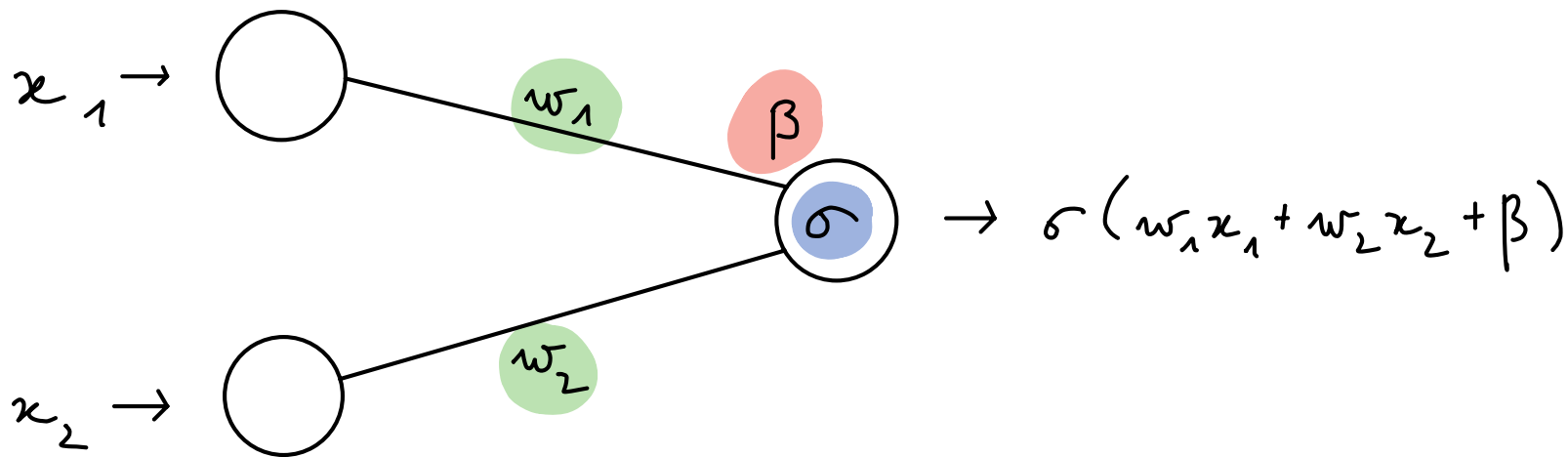
price of ticket $\longrightarrow$ (○)

weather quality $\longrightarrow$ (○) $\longrightarrow$ (○) $\longrightarrow$ decision to go or not

1) Set weights: $w_1 =$ importance of cheapness

$w_2 =$ importance of weather.

2) Set a bias: I'll go if $w_1 x_1 + w_2 x_2 > -\beta$

3) Choose an activation function $\sigma$ $\left(e.g. \quad \sigma(z) = \dfrac{1}{1+e^{-z}}\right)$

$x_1 \rightarrow$ ◯ $\xrightarrow{w_1}$

$\beta$

◯ $\sigma$ $\rightarrow \sigma(w_1 x_1 + w_2 x_2 + \beta)$

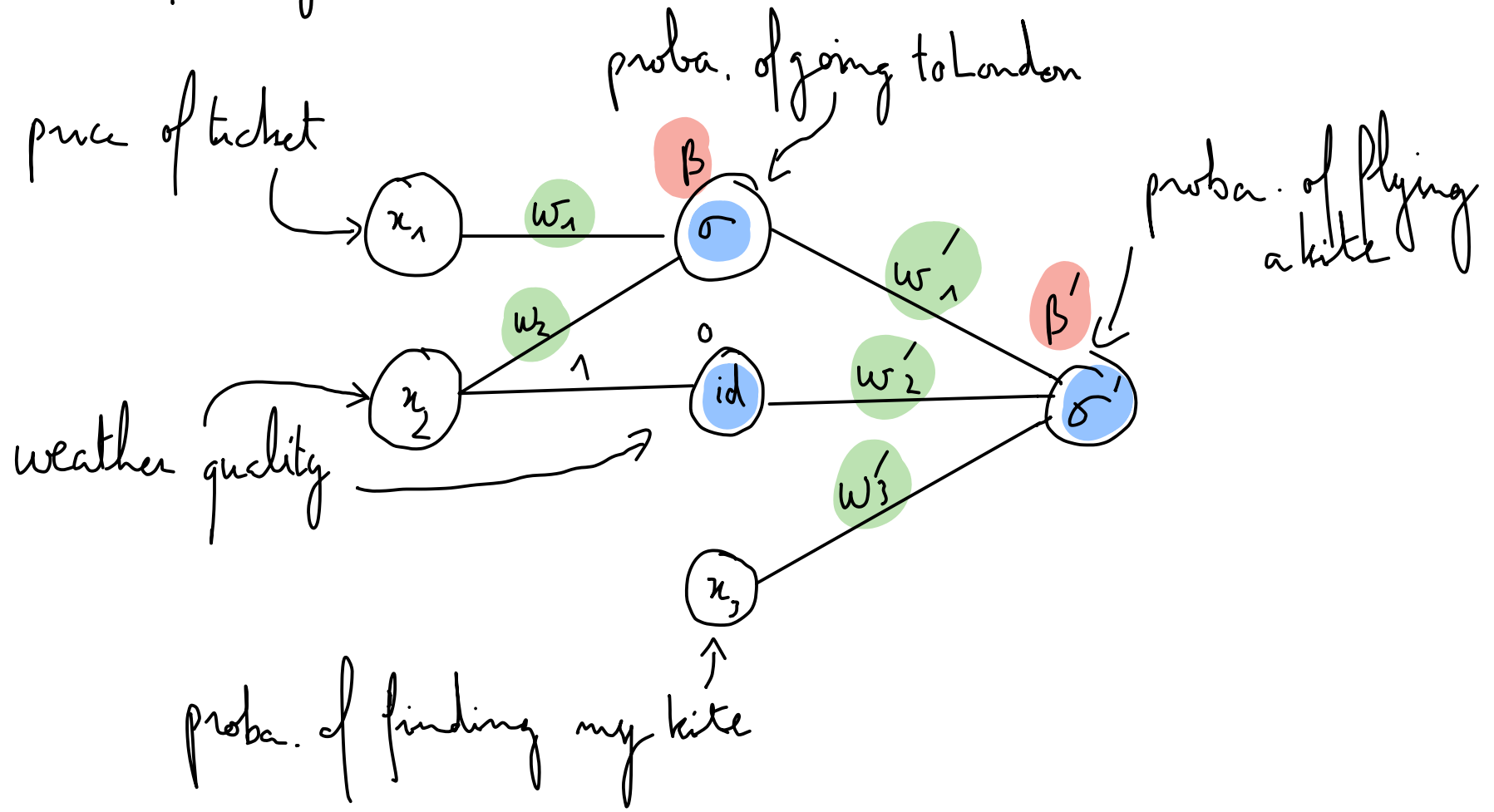$x_2 \rightarrow$ ◯ $\xrightarrow{w_2}$

1 - Layered Neural Network

# Composing Neural Networks

I want to fly my kite     a) If I can find it,
                          b) If the weather is good,
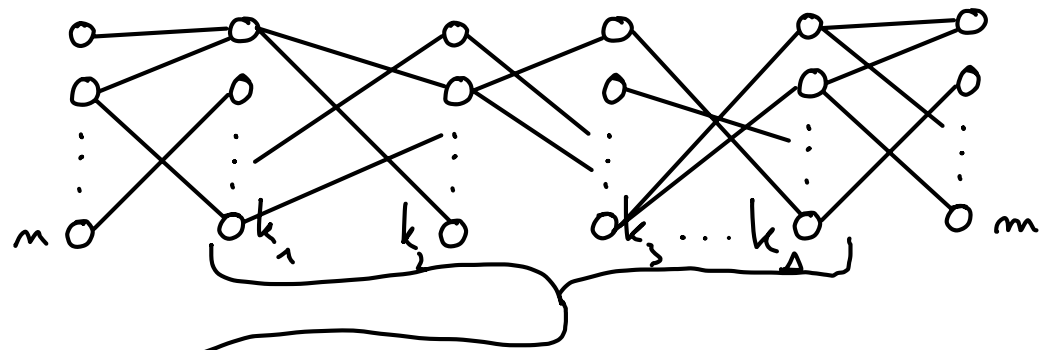
but if I go to London, I may be too busy.

proba. of going to London

price of ticket

proba. of flying a kite

proba. of finding my kite

weather quality

$$x_1 \xrightarrow{W_1} \beta\ \sigma \xrightarrow{W_1'} \beta'\ \sigma'$$

$$x_2 \xrightarrow{W_2} \quad \xrightarrow{\ \ } id \xrightarrow{W_2'}$$

$$x_3 \xrightarrow{W_3'}$$

# IV Categorical Neural Networks

Def: The category (PROP) NNET has:

- objects : $\mathbb{N}$
- morphism $n \longrightarrow m$: freely generated by $\mathcal{P}([n] \times [m])$



1-layered morphism

$$n \longrightarrow m$$

s-layered morphism

$$n \longrightarrow m$$

"hidden layers"

Prop (Fong - Spivak - Tuyéras): Given $\sigma : \mathbb{R} \to \mathbb{R}$

there is a functor

$$I_\sigma : NNET \longrightarrow Para(Smooth) \text{ such that}$$

$$I_\sigma(n) = n, \text{ and for } C \subseteq [n] \times [m],$$

$$I_\sigma(C) : n \xrightarrow{|C| + m, I} m \quad \text{with:}$$

$$I : \mathbb{R}^{|C| + m} \times \mathbb{R}^n \longrightarrow \mathbb{R}^m$$

$$\left( \left( w_{ji} \right)_{(j,i) \in C}, \left( \beta_j \right)_{j \in [m]}, \left( x_i \right)_{i \in [m]} \right) \longmapsto \left( \sigma \left( \left( \sum_{(j,i) \in C} w_{ji} x_i \right) + \beta_j \right) \right)_{j \in [m]}$$

Terminology (Fong - Spivak -Tuyéras): Given $\sigma$, $\varepsilon$, $e$
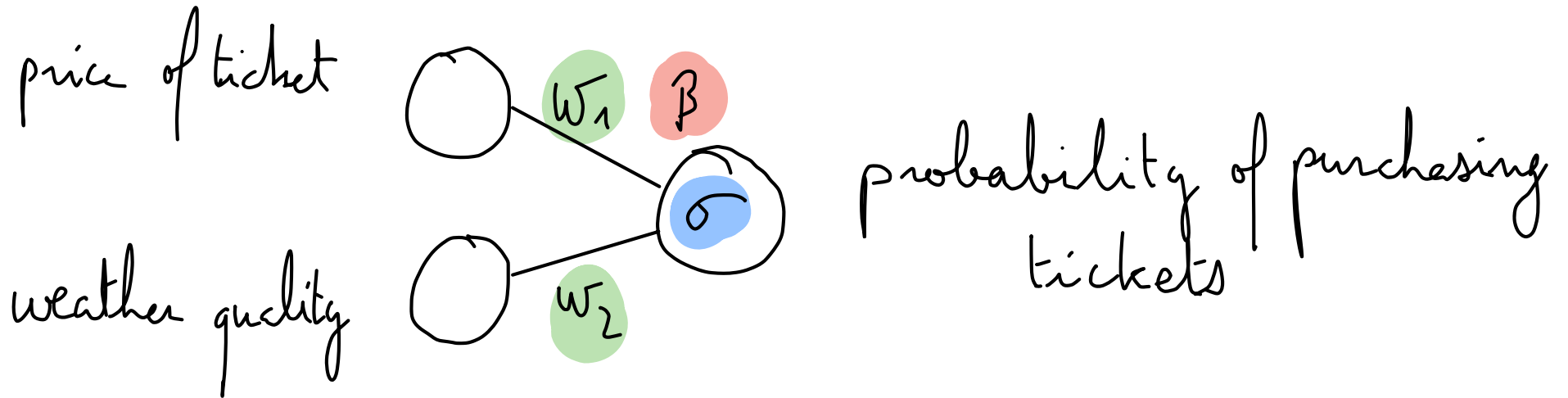
$$G_{\varepsilon, e} \circ \underline{I}_{\sigma} : \text{NNET} \longrightarrow \text{Para}(\text{Optic}(\text{Smooth}))$$

is <u>deep learning</u>. The resulting request function

is <u>deep dreaming</u>.

# V Training Your Network

## "Go to London" example:

price of ticket

weather quality



probability of purchasing tickets

a) Pick a starting guess for $w_1, w_2, \beta$

b) Train your network using a data set:

$$\{(\text{conditions}, \text{outcome})\}, \text{ e.g:}$$

$$((\text{ticket was } 80\text{€}, \text{weather was meh}), \text{did not purchase})$$

For each sample data, compute:

$$\cup \left( W_1, W_2, \beta, \text{ conditions, outcome} \right) = \text{new } W_1, W_2, \beta$$

$$\cdots$$

At the end,

$$I \left( W_1, W_2, \beta, \text{ conditions} \right) \quad \text{predicts the decision.}$$

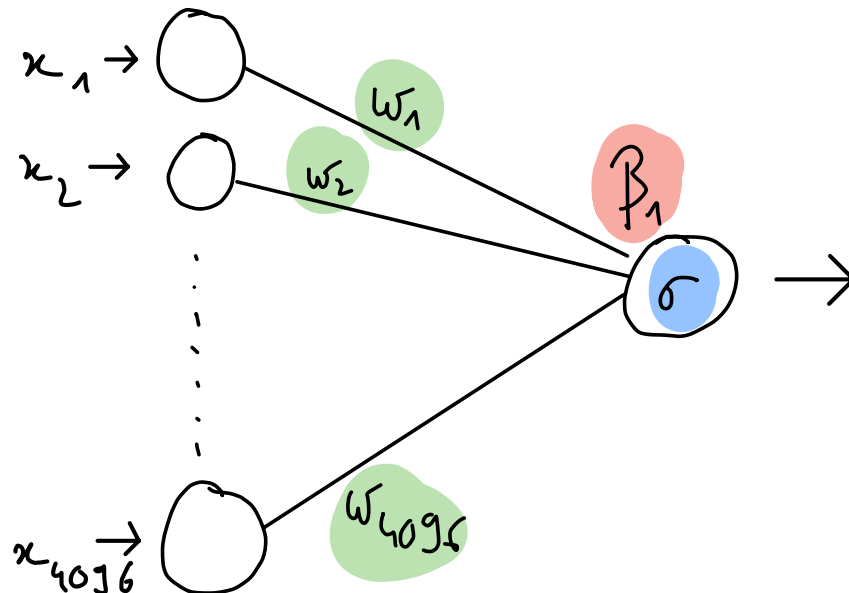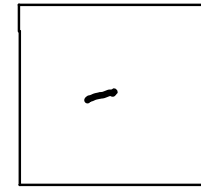# Example : Image recognition.

Aim : Recognise hand-drawn digits in b & w,

$(4096 =) 64 \times 64$ pixels images.

1) Encoding images : $(x_i) \in \mathbb{R}^{4096}$ $]-\infty ; +\infty[$
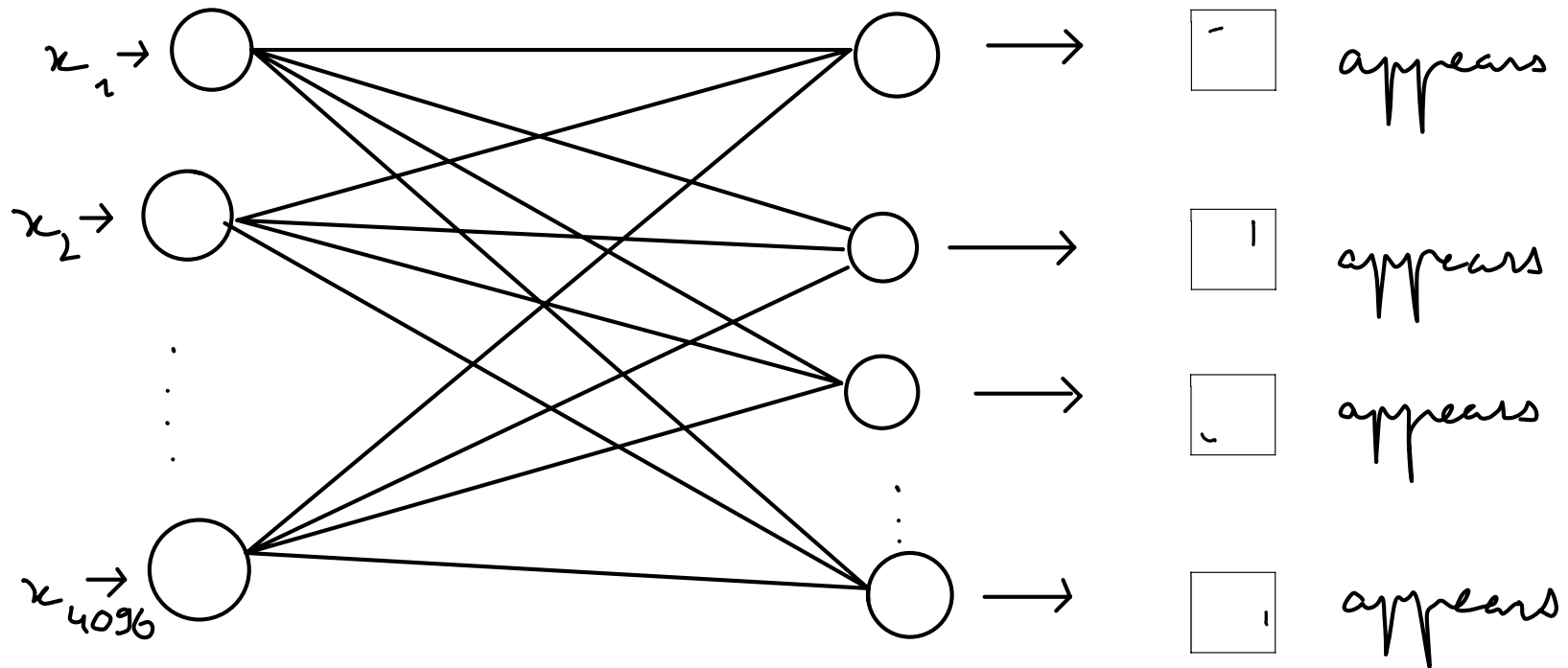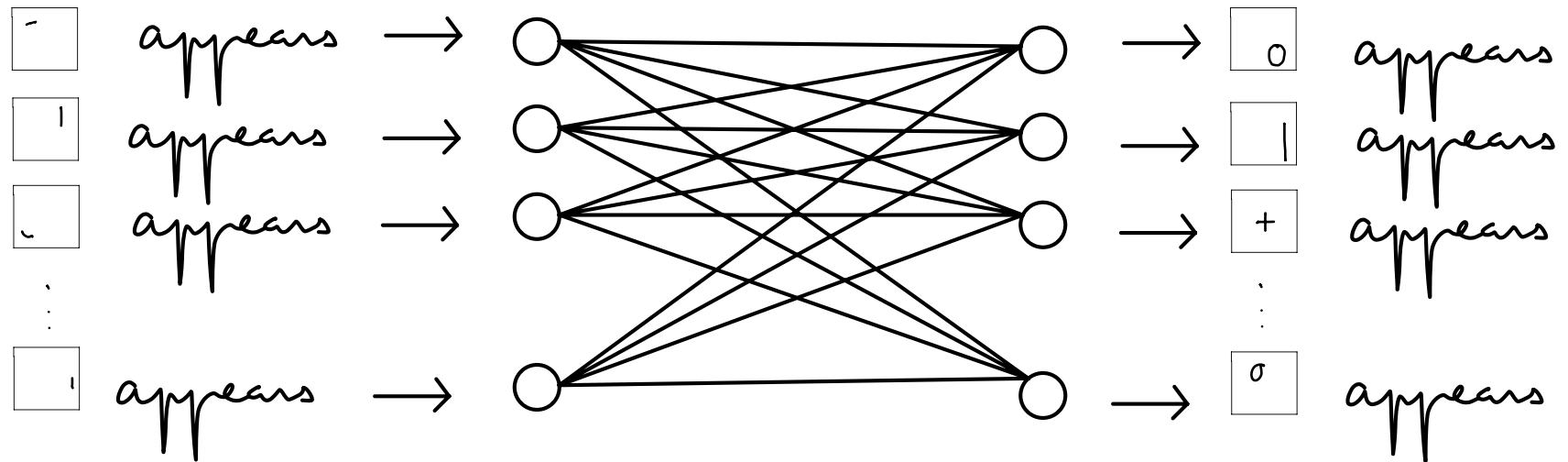
2) First Layers recognise small edges :

3) Train your network with sample data:

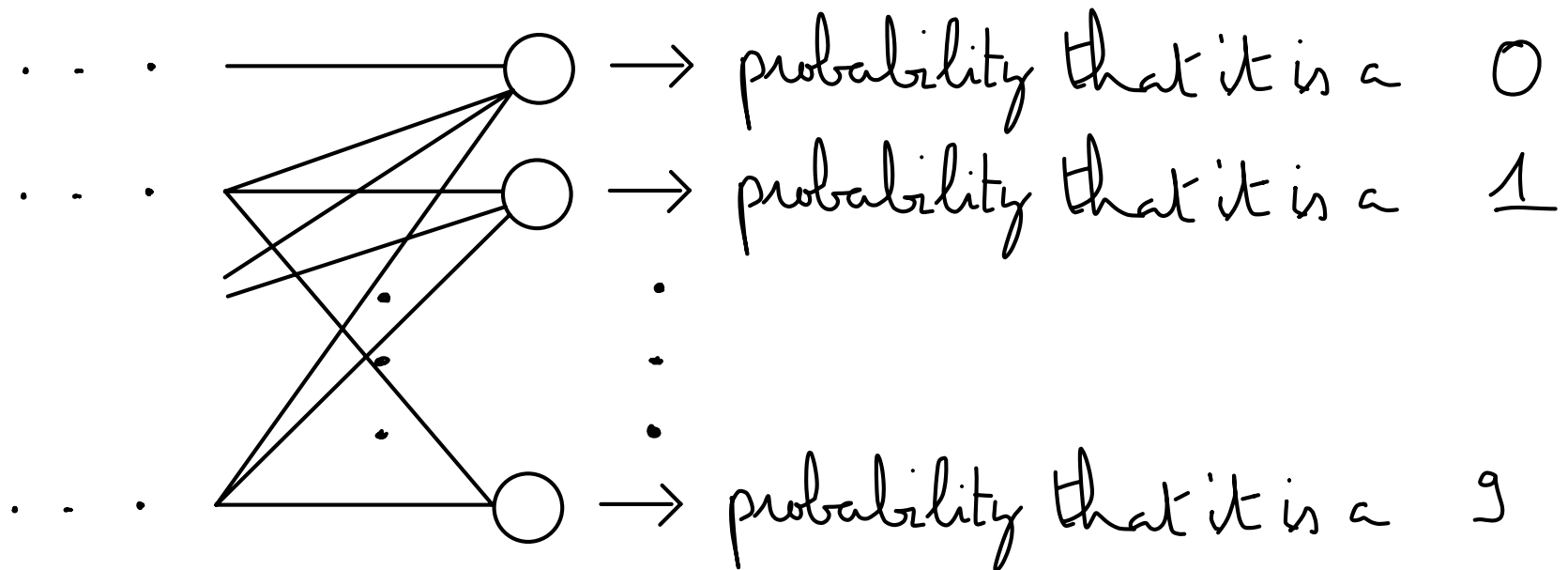$$\{ \; ( \; \boxed{5} \; , \; \text{mostly yes} \; ) \}$$

4) Repeat for <u>many</u> edges to get first layer.



$x_1 \rightarrow$ ⟶ $\boxed{-}$ appears

$x_2 \rightarrow$ ⟶ $\boxed{1}$ appears

⟶ $\boxed{\phantom{x}}$ appears

$x_{4096} \rightarrow$ ⟶ $\boxed{1}$ appears

# 5) Second layers, third layers...: more complex shapes



| | | |
|---|---|---|
| ⊢ appears → | ○ —→ | 0 appears |
| 1 appears → | ○ —→ | 1 appears |
| ∪ appears → | ○ —→ | + appears |
| ⋮ | | ⋮ |
| ⌐ appears → | ○ —→ | σ appears |

# 6) Final Layer



· · · —————○ —→ probability that it is a  0

· · · —————○ —→ probability that it is a  1

⋮

· · · —————○ —→ probability that it is a  9

What does $r$ do?

$$\mathcal{L}\left( (w, \beta), \boxed{\,\varphi\,} , \text{``This is definitely an } 8\text{''} \right)$$

$$=$$

$$\boxed{\,8\,}$$

$$\ldots \text{ Image Generation}$$

# Thank You !

References:

Fong, Spivak, Tuyéras - Backprop as Functor: A compositional perspective on supervised learning

Cruttwell, Gavranović, Ghani, Wilson, Zanasi - Categorical Foundations of Gradient-Based Learning

Cockett, Cruttwell, Gallagher, Lemay, MacAdam, Plotkin, Pronk - Reverse derivative categories

Capucci, Gavranović, Hedges, Rischel - Towards Foundations of Categorical Cybernetics