

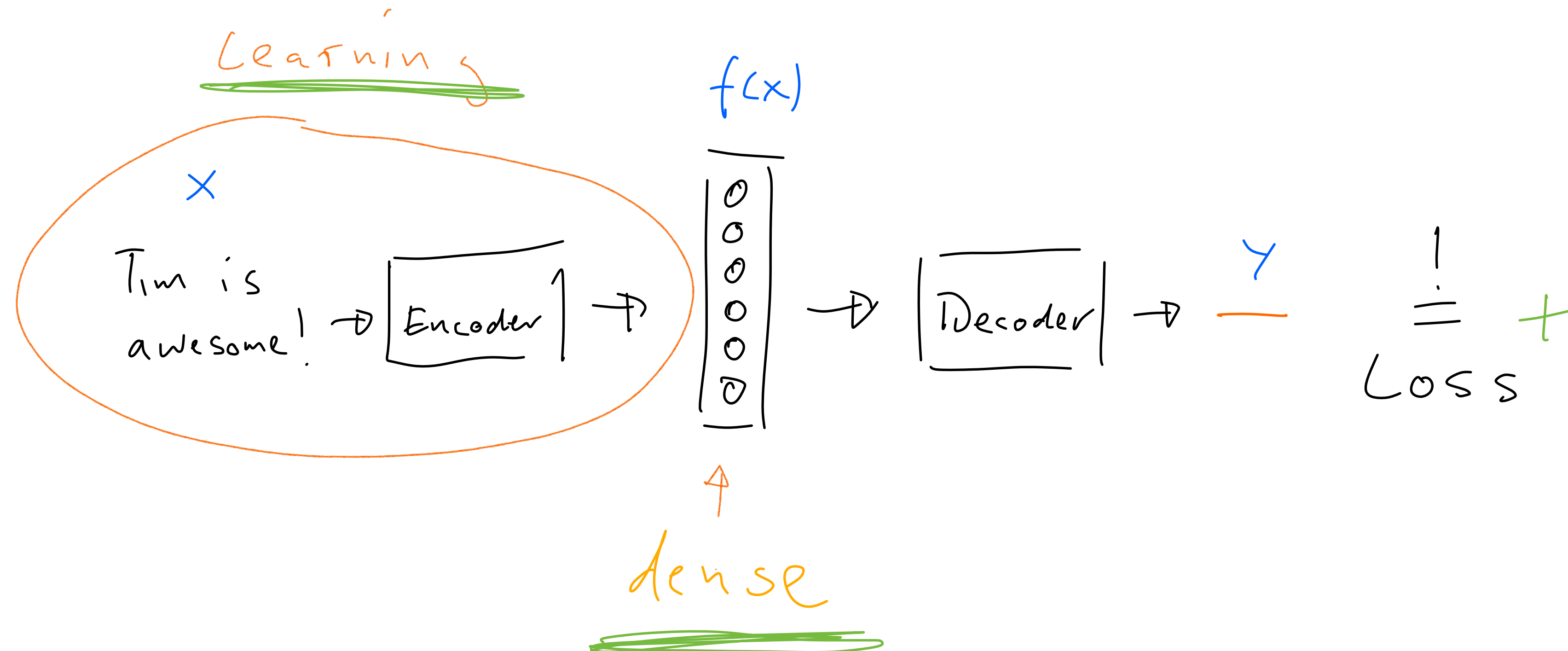
Introduction to Neural Networks and Backpropagation

Tim Rocktäschel & **Sebastian Riedel**
COMP0087 Natural Language Processing



Overview

This Lecture!



Task: Text Classification

St Pauli: the club that stands for all the right things ... except winning

→ **Sports**

Brexit: Are we running out of parliamentary time?

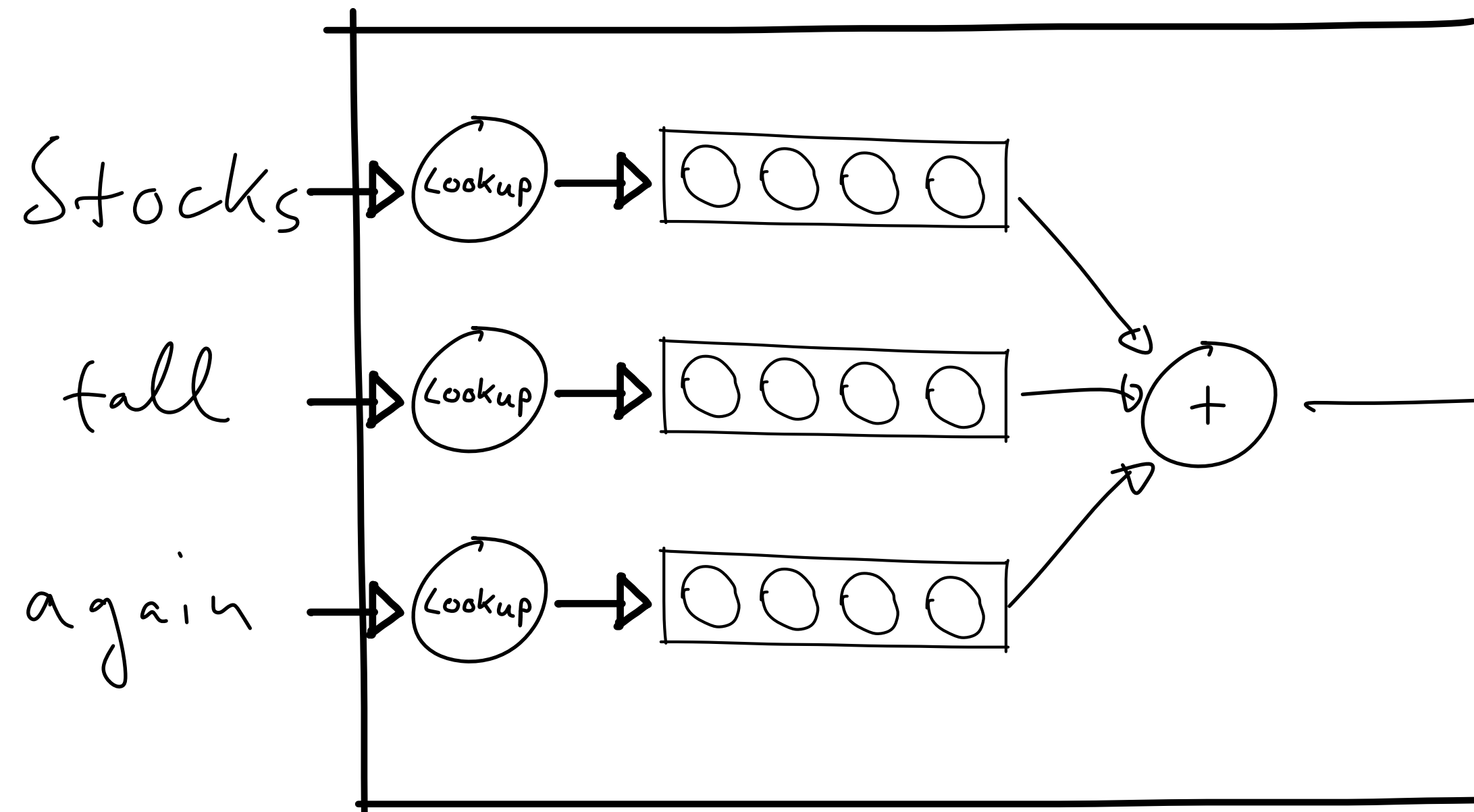
→ **Politics**

Stocks fall after Morgan Stanley earnings miss

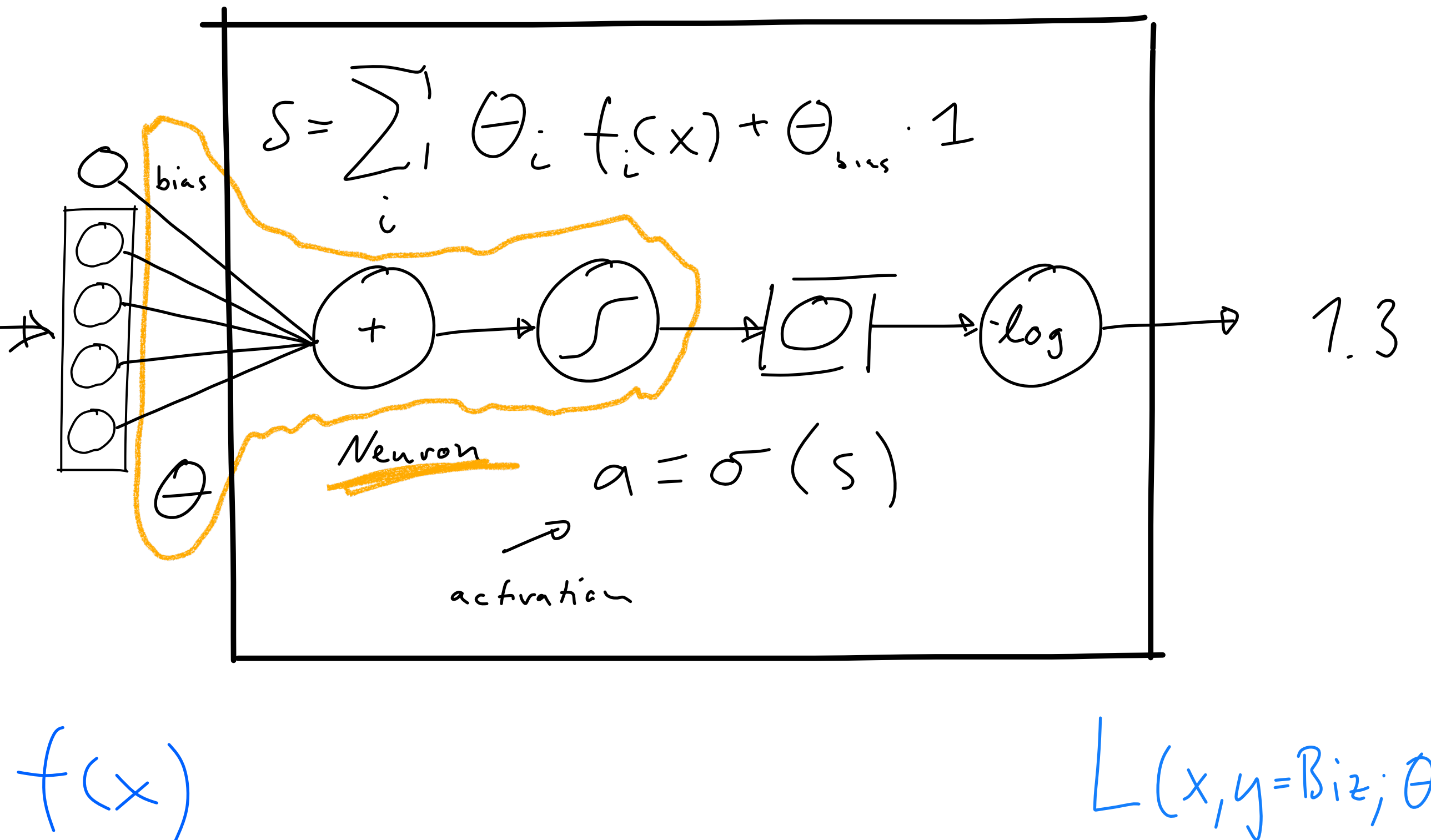
→ **Business**

Computation Graph

Encoder

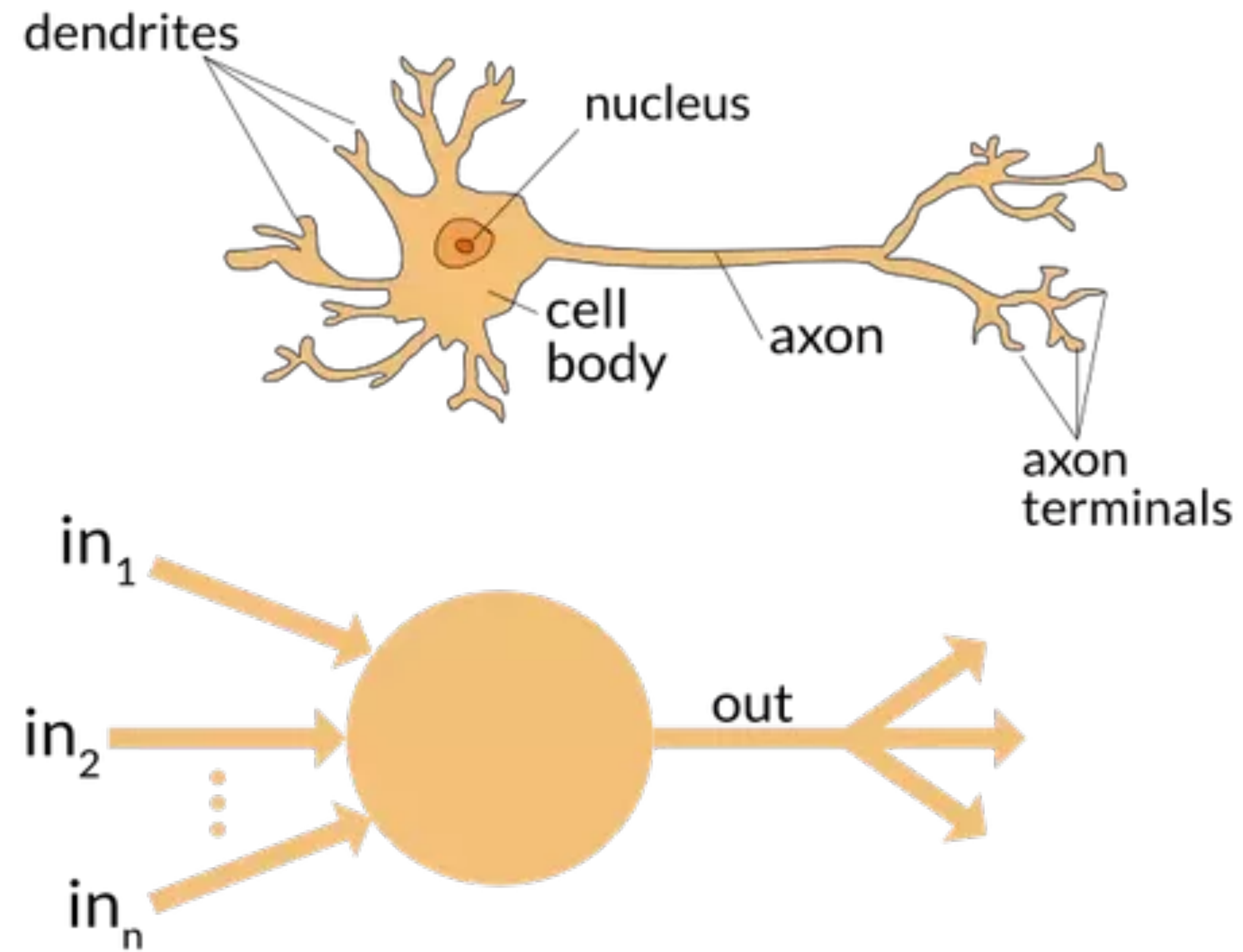


Decoder

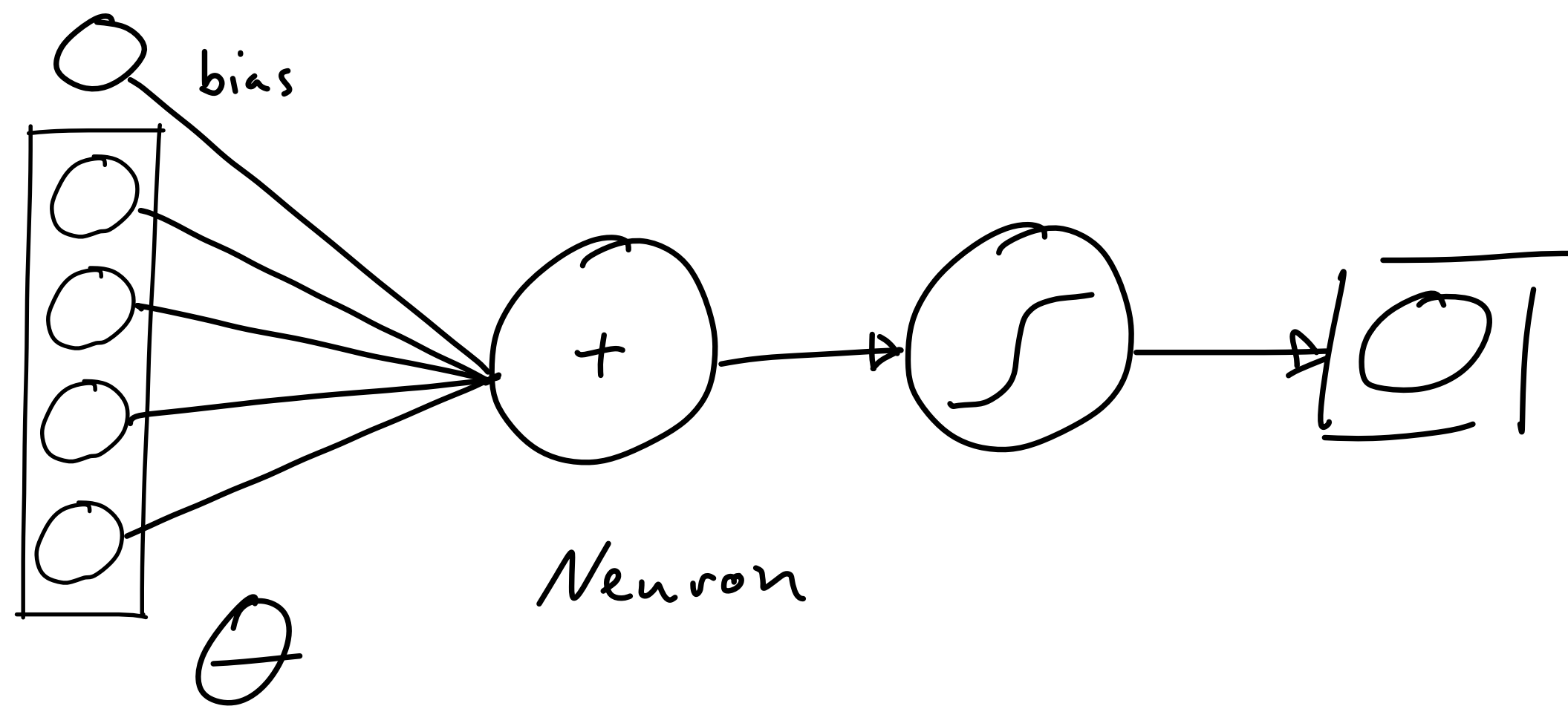


$$L(x, y = \text{Biz}; \theta)$$

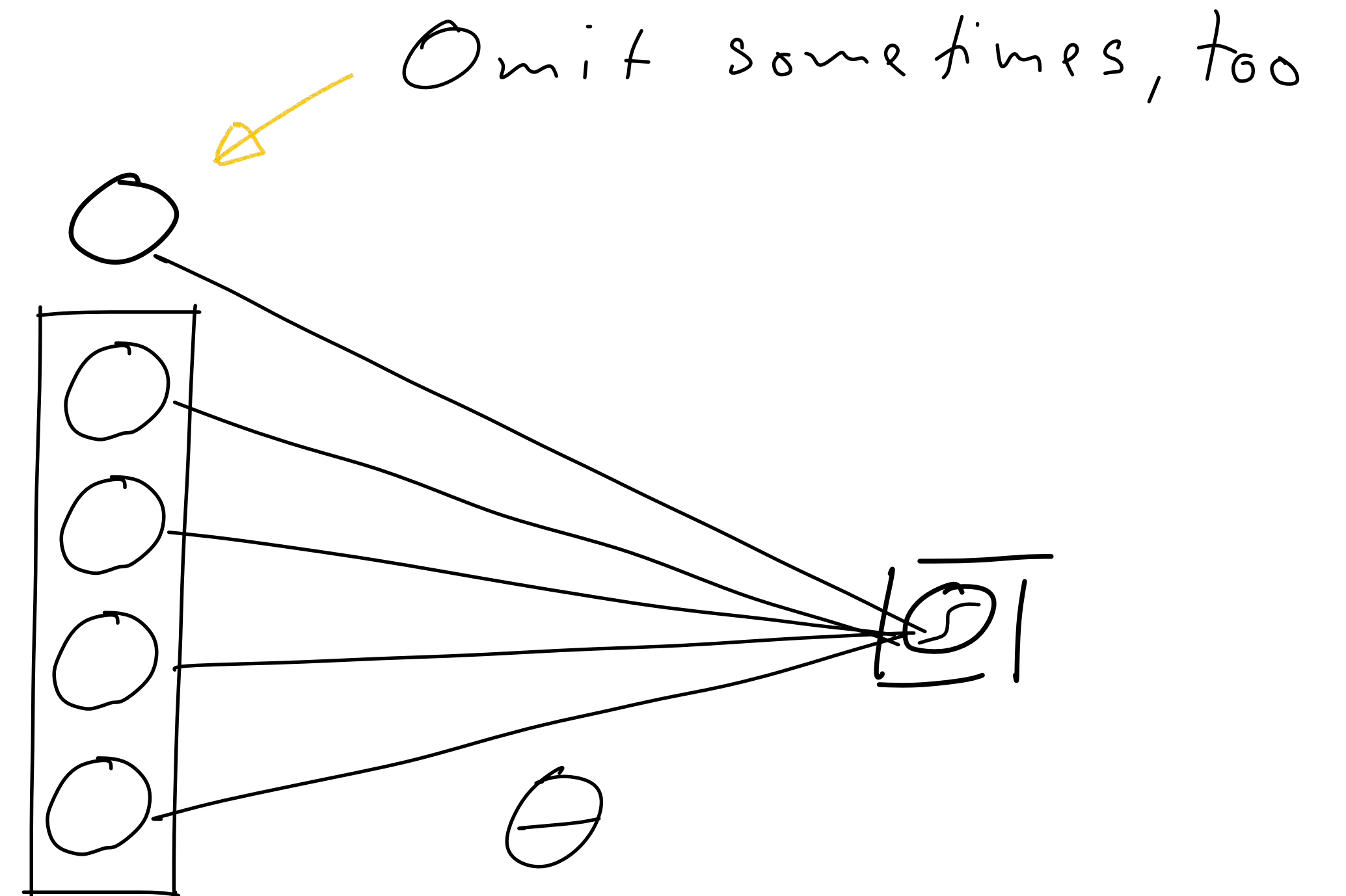
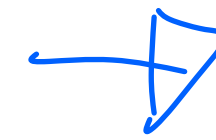
Artificial and Biological Neurons



Simplified Notation

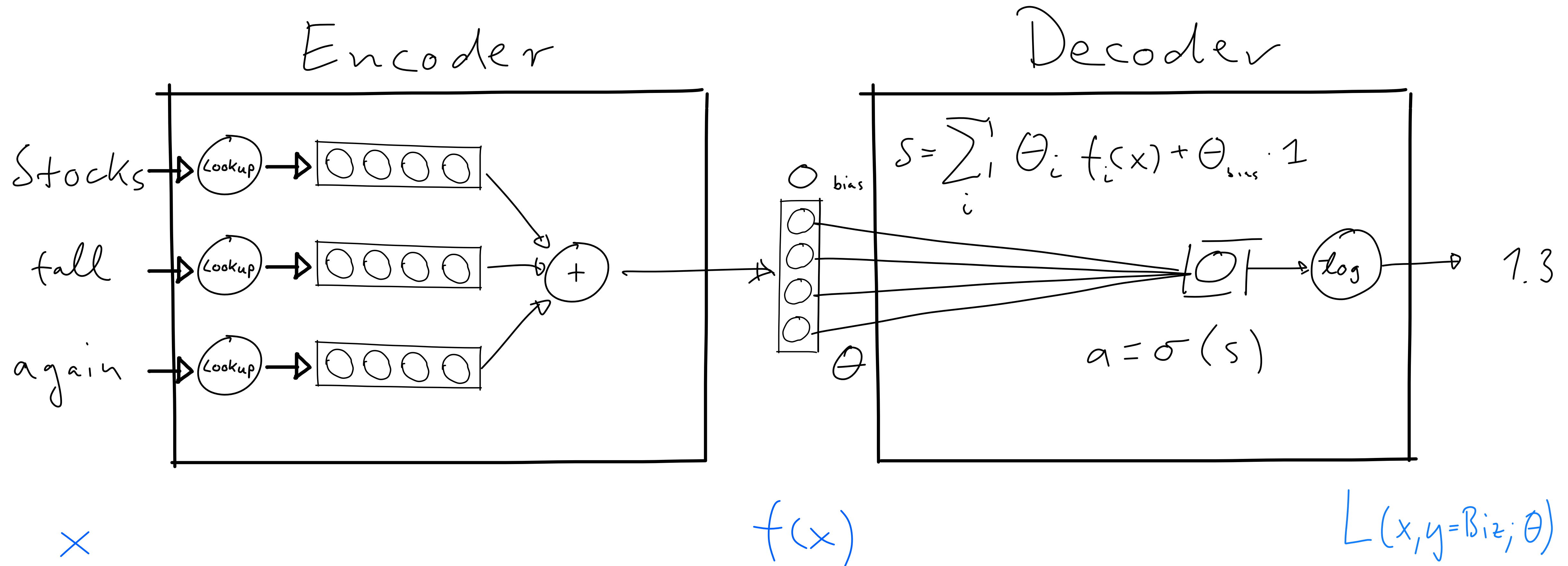


$$a = \sigma(\Theta^T x)$$



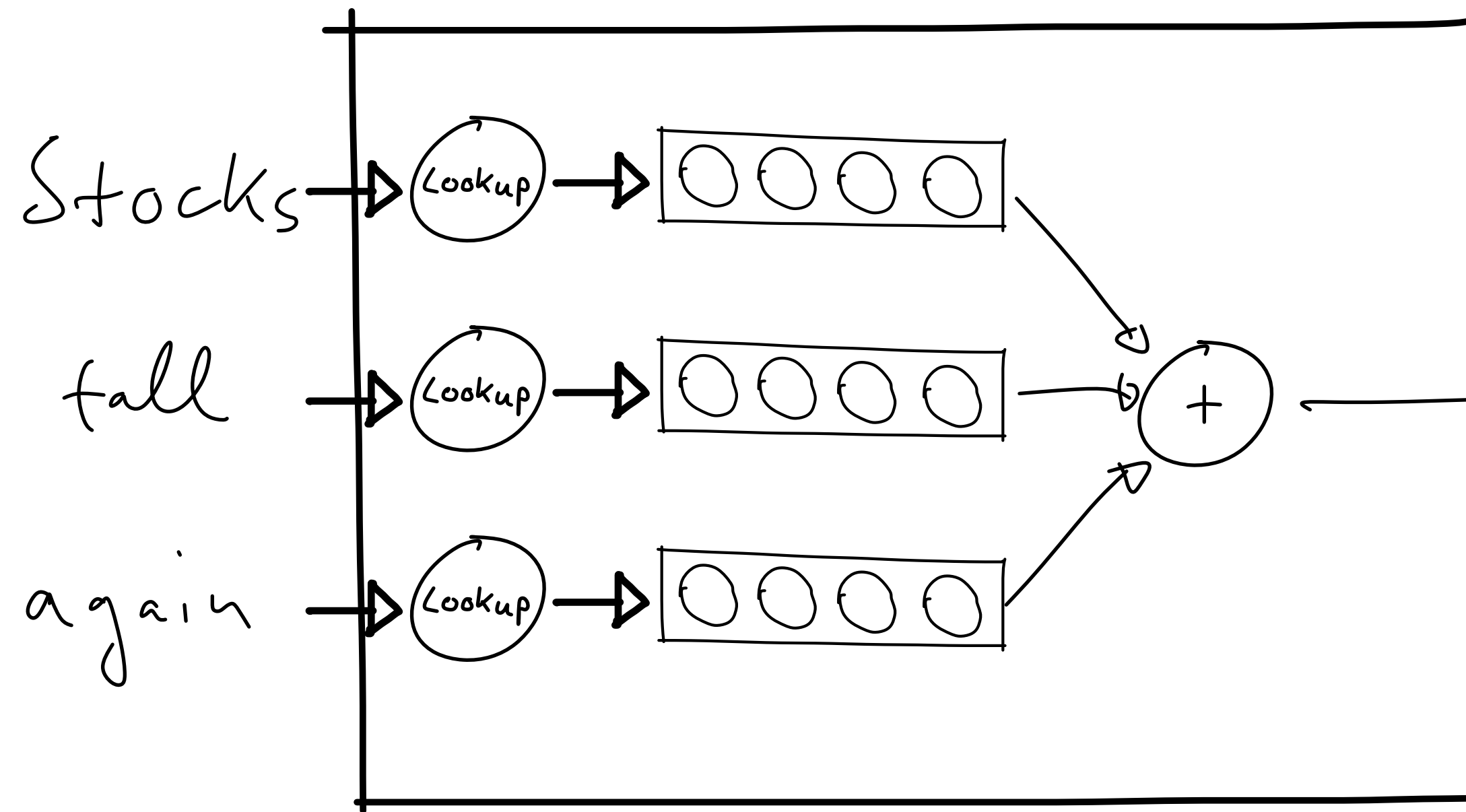
$$a = \sigma(\Theta^T x)$$

Computation Graph

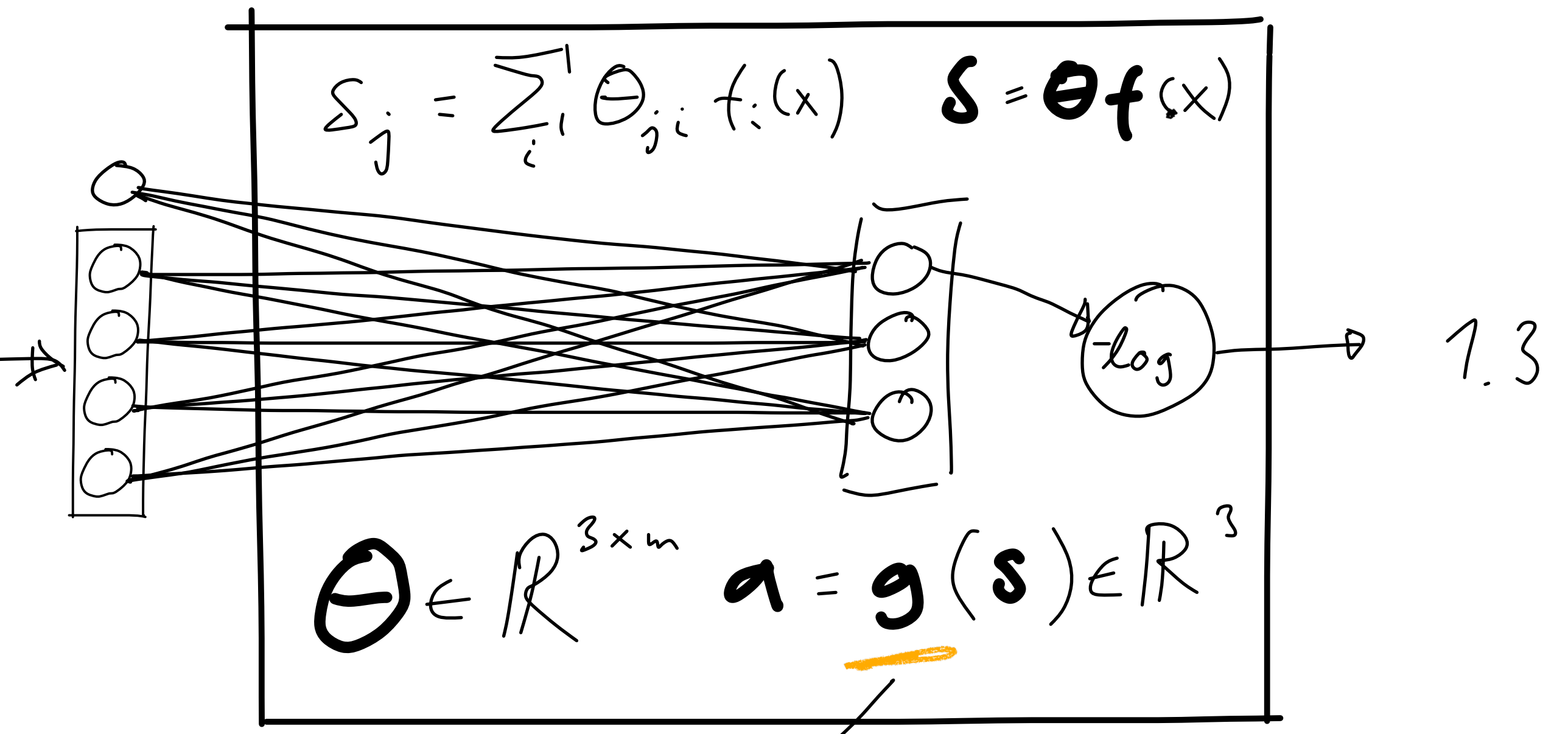


Multiclass

Encoder



Decoder



×

$f(x)$
 $\in \mathbb{R}^m$

How to define g
such that a is
a distribution?

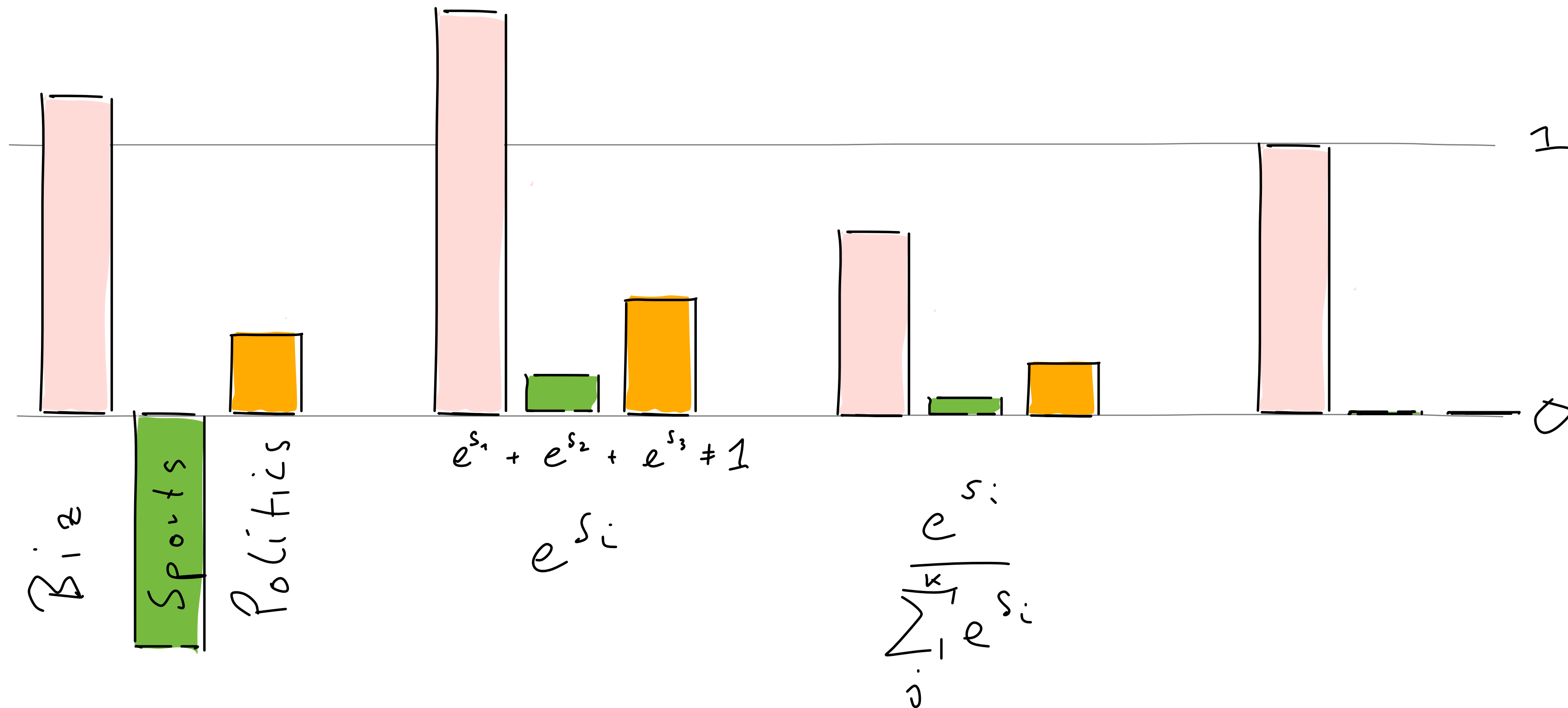
$L(x, y = \text{Biz}; \Theta)$

Softmax

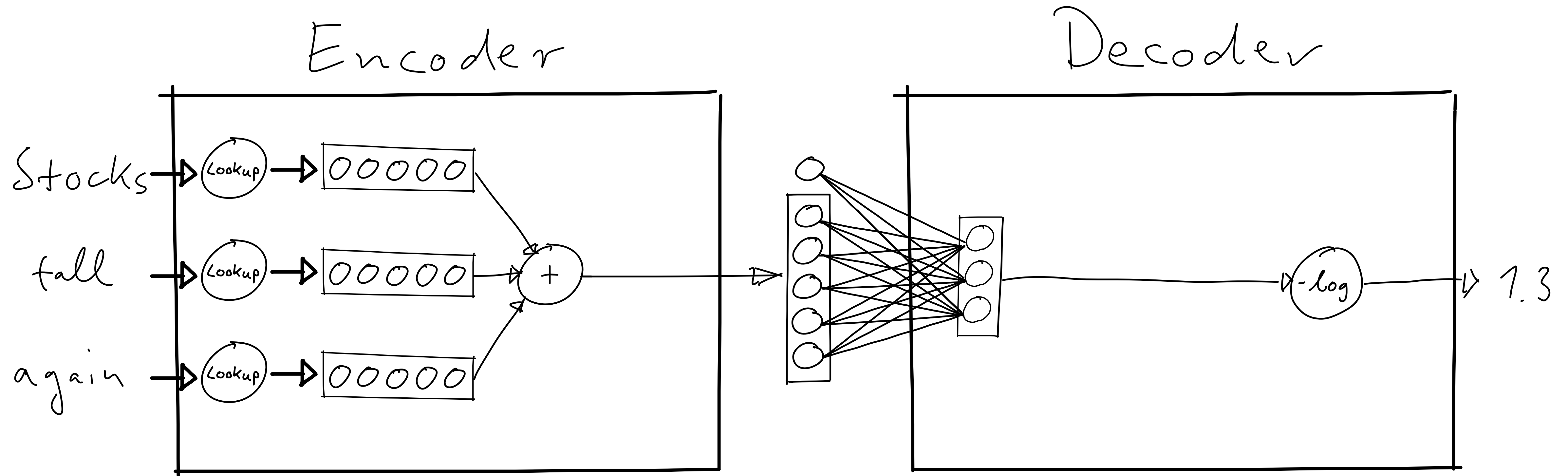
Scores/logits non-negative normalized max

S

$\text{Softmax}(S)$

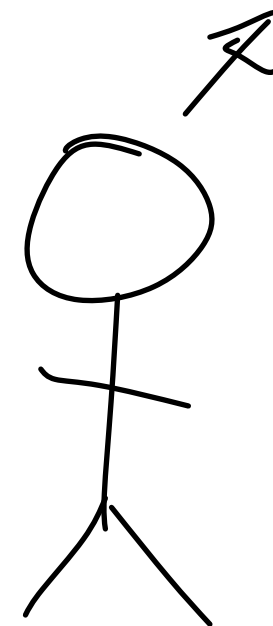


Limits of Feature Engineering



×

engineer



"Boo! Can we just learn this?"

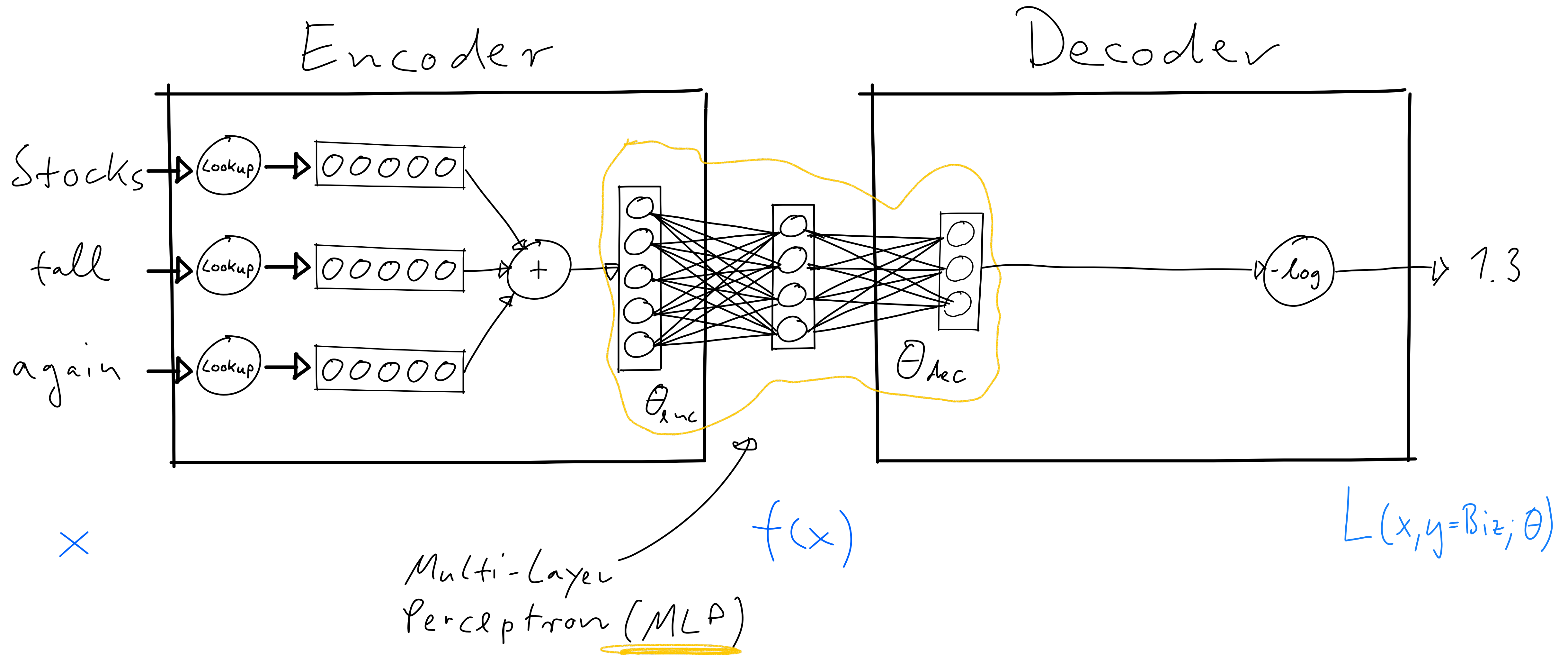
$f(x)$

① tedious

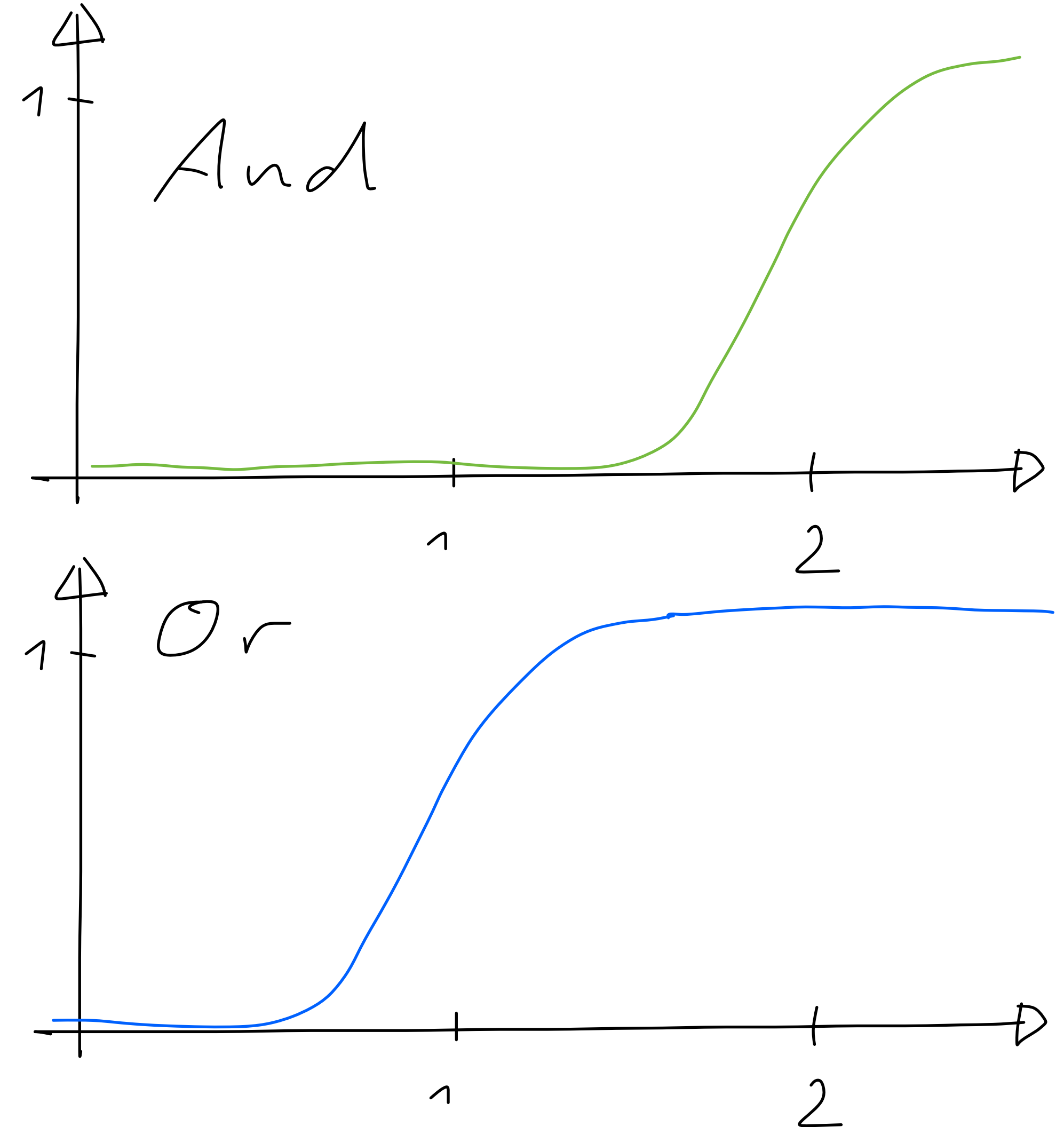
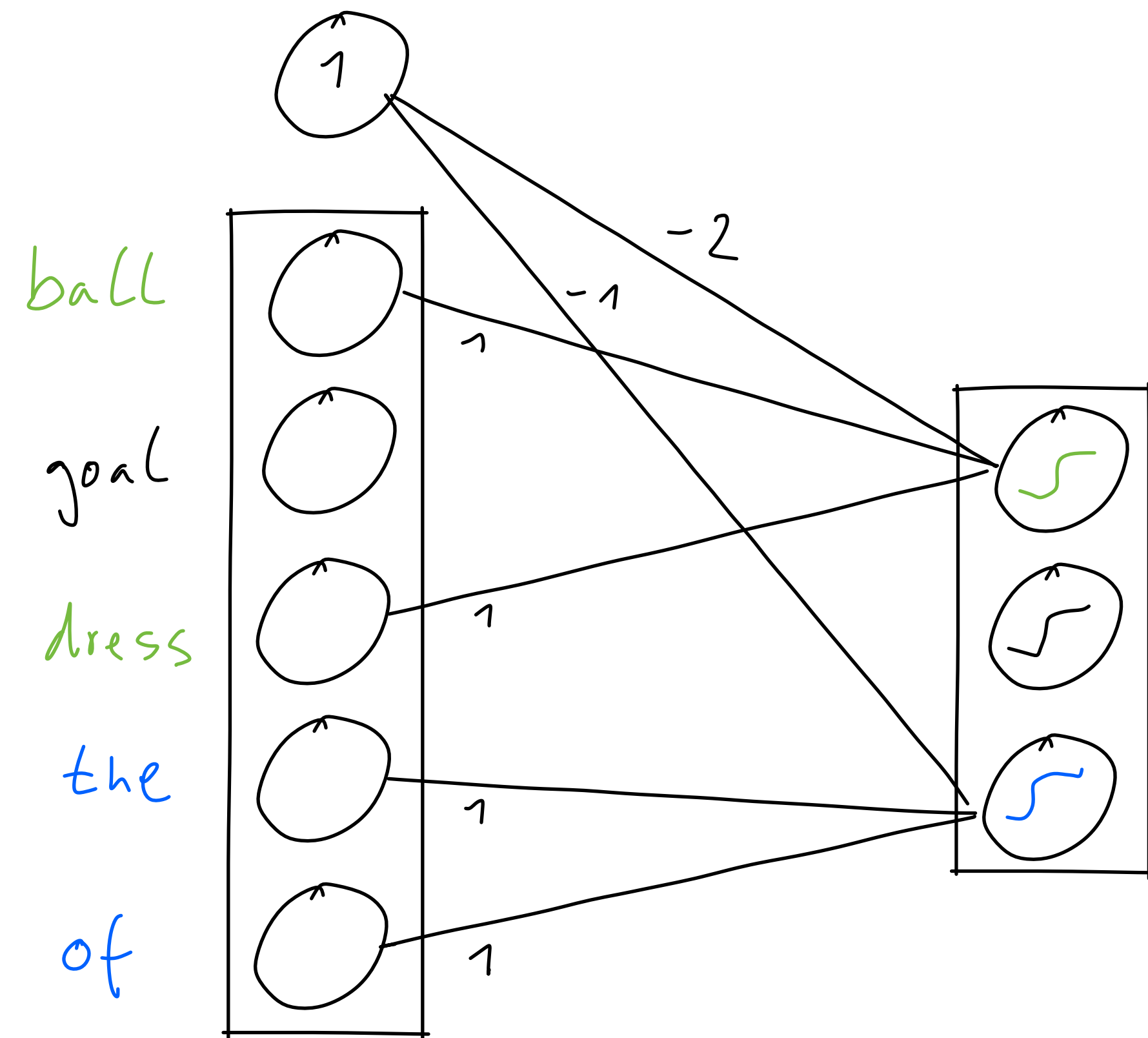
② suboptimal

$L(x, y = Biz; \theta)$

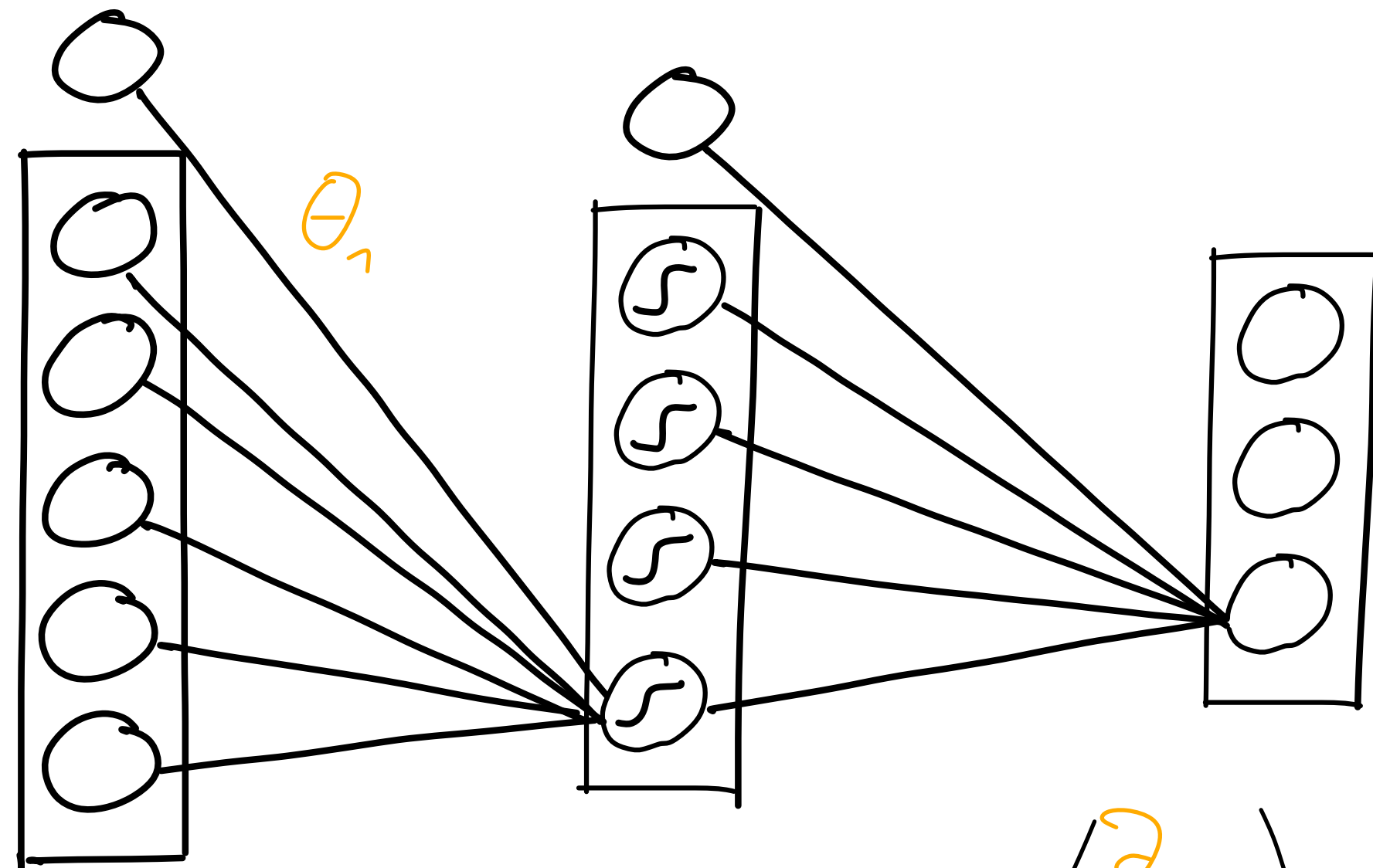
Deep Learning



What can the layer learn?



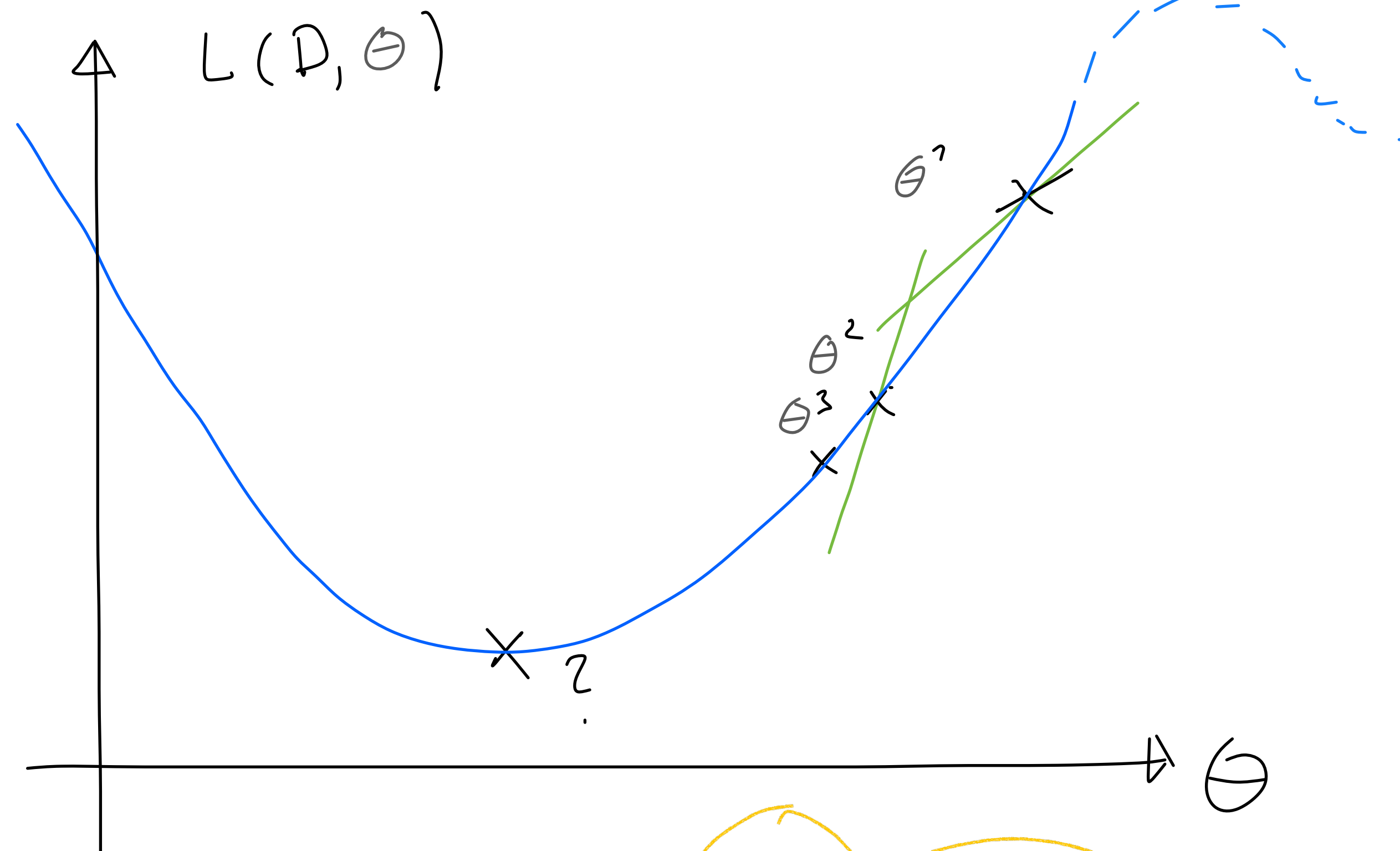
Training an MLP



$$\Theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_m \end{pmatrix}$$

$$\nabla_{\Theta} = \begin{pmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_m} \end{pmatrix}$$

can be calculated
automatically



$$\Theta^{i+1} = \Theta^i - \alpha \nabla_{\Theta} l(x_i, y_i, \Theta^i)$$

need this!

Warning: It will get a little
technical

But: This is (almost) all there is to
Deep Learning

Leibniz Notation

$$x \in \mathbb{R}$$

$$y = f(x) \in \mathbb{R}$$

derivative of y with respect to x

$$\frac{\partial y}{\partial x}(x) = f'(x)$$

often omitted

if I wiggle x a tiny
tiny bit at x , how
does y change?

$$\in \mathbb{R}$$

$$\mathbf{x} \in \mathbb{R}^n$$

$$y = f(\mathbf{x}) \in \mathbb{R}$$

$$\frac{\partial y}{\partial \mathbf{x}}(\mathbf{x}) = \begin{pmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{pmatrix}$$

||

$$\underline{\nabla_{\mathbf{x}} f(\mathbf{x})}$$

$$\in \mathbb{R}^n$$

$$\mathbf{x} \in \mathbb{R}^n$$

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}(\mathbf{x}) = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

$$\in \mathbb{R}^{n \times m}$$

Jacobian

Chain Rule

$$a = f(x)$$

$$wx$$

$$b = g(a)$$

$$\log(a)$$

$$\frac{\partial b}{\partial x}(x) = \boxed{\frac{\partial a}{\partial x}(x)}^{\text{Local}} \cdot \boxed{\frac{\partial b}{\partial a}(a)}^{\text{Upstream}}$$

$$\frac{\partial \log(wx)}{\partial x} \stackrel{?}{=}$$

Multidimensional Chain Rule

$$\mathbf{x} \in \mathbb{R}^l$$

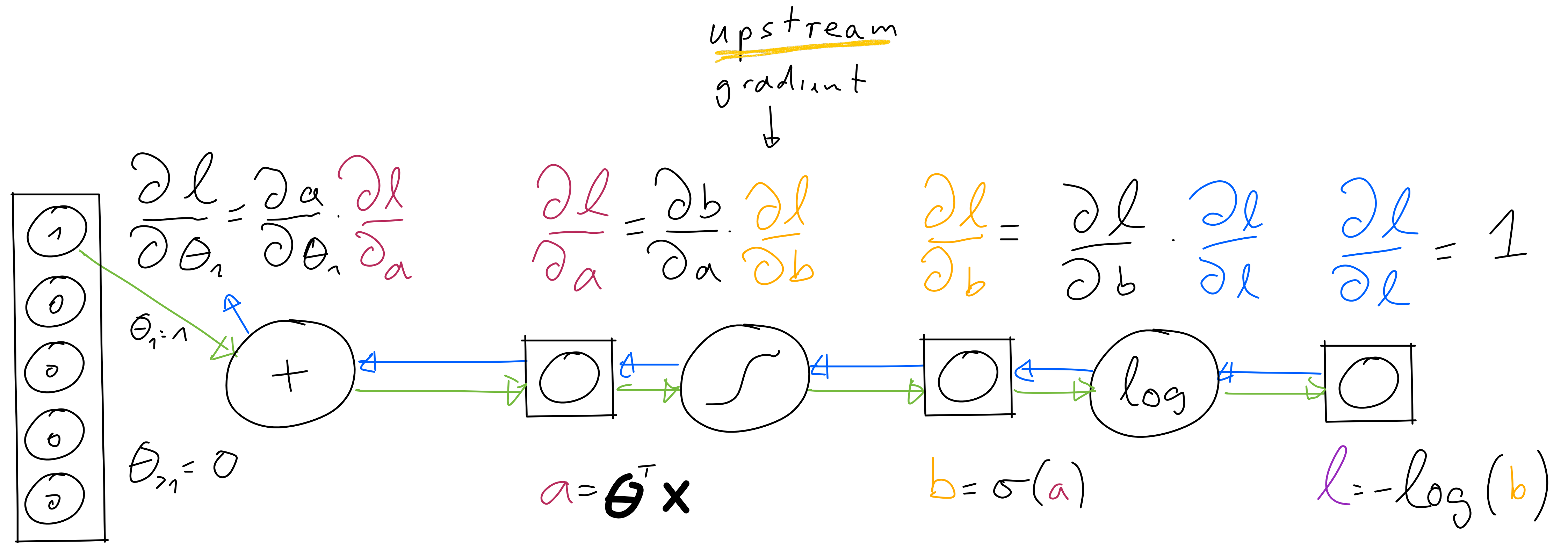
$$\mathbf{a} = \mathbf{f}(\mathbf{x}) \in \mathbb{R}^n$$

$$\mathbf{b} = \mathbf{g}(\mathbf{a}) \in \mathbb{R}^m$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}}(\mathbf{x}) = \boxed{\frac{\partial \mathbf{a}}{\partial \mathbf{x}}(\mathbf{x})}^{\text{Local}} \cdot \boxed{\frac{\partial \mathbf{b}}{\partial \mathbf{a}}(\mathbf{a})}^{\text{Upstream}}$$

$\mathbb{R}^{l \times m}$
 $\mathbb{R}^{l \times n}$
 $\mathbb{R}^{n \times m}$

Forward and Backward Pass

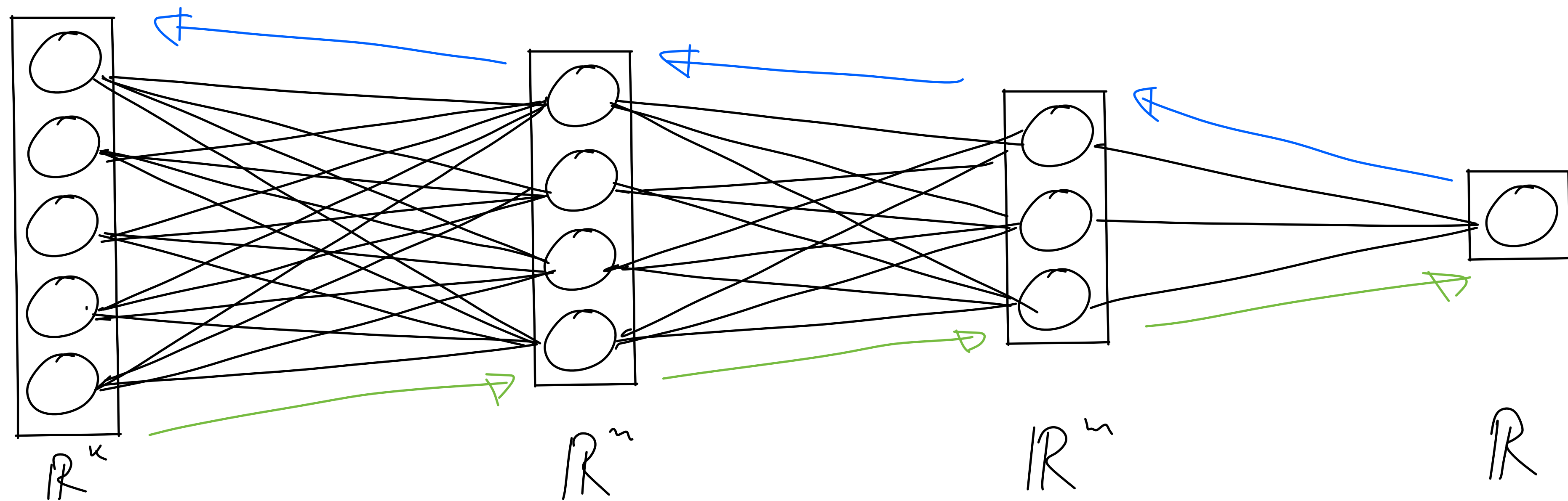


X

Multi-dimensional

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{a}} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{a}} = \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{b}} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{b}}$$

1



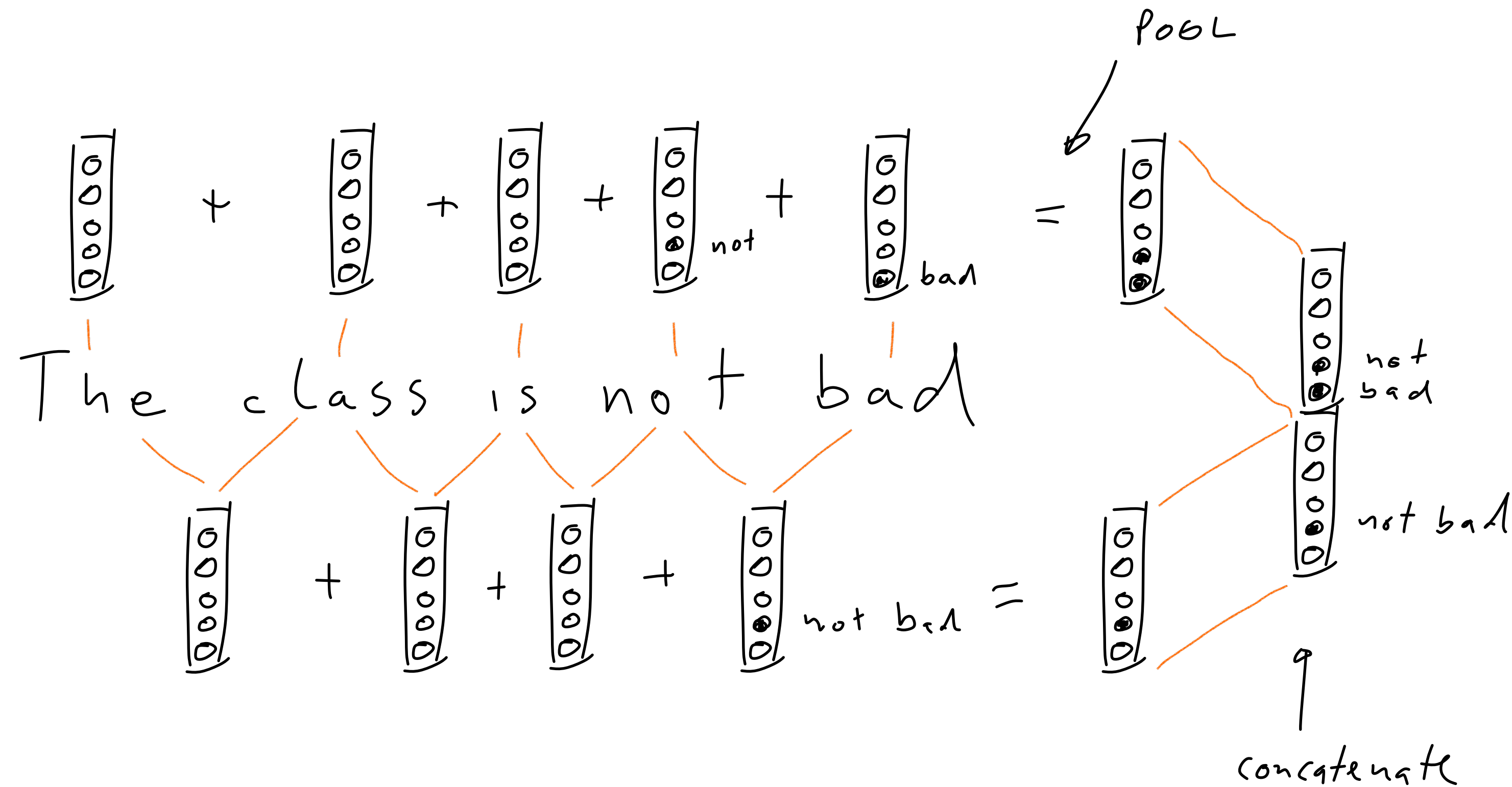
\mathbf{x}

$$\mathbf{a} = \sigma(\mathbf{W}\mathbf{x})$$

$$\mathbf{b} = \sigma(\mathbf{V}\mathbf{a})$$

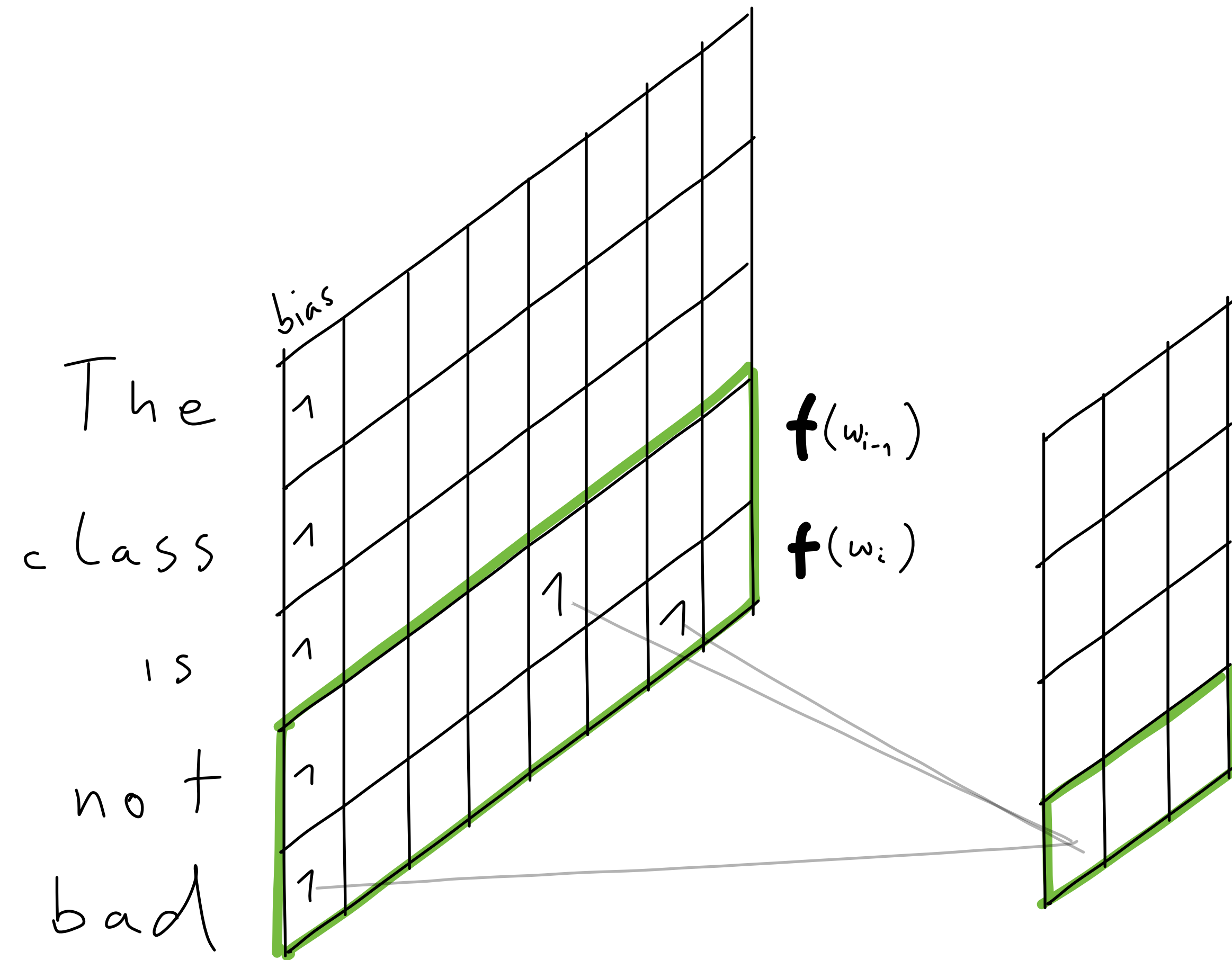
$$c = \sigma(\boldsymbol{\theta}^T \mathbf{b})$$

Remember Bigrams?



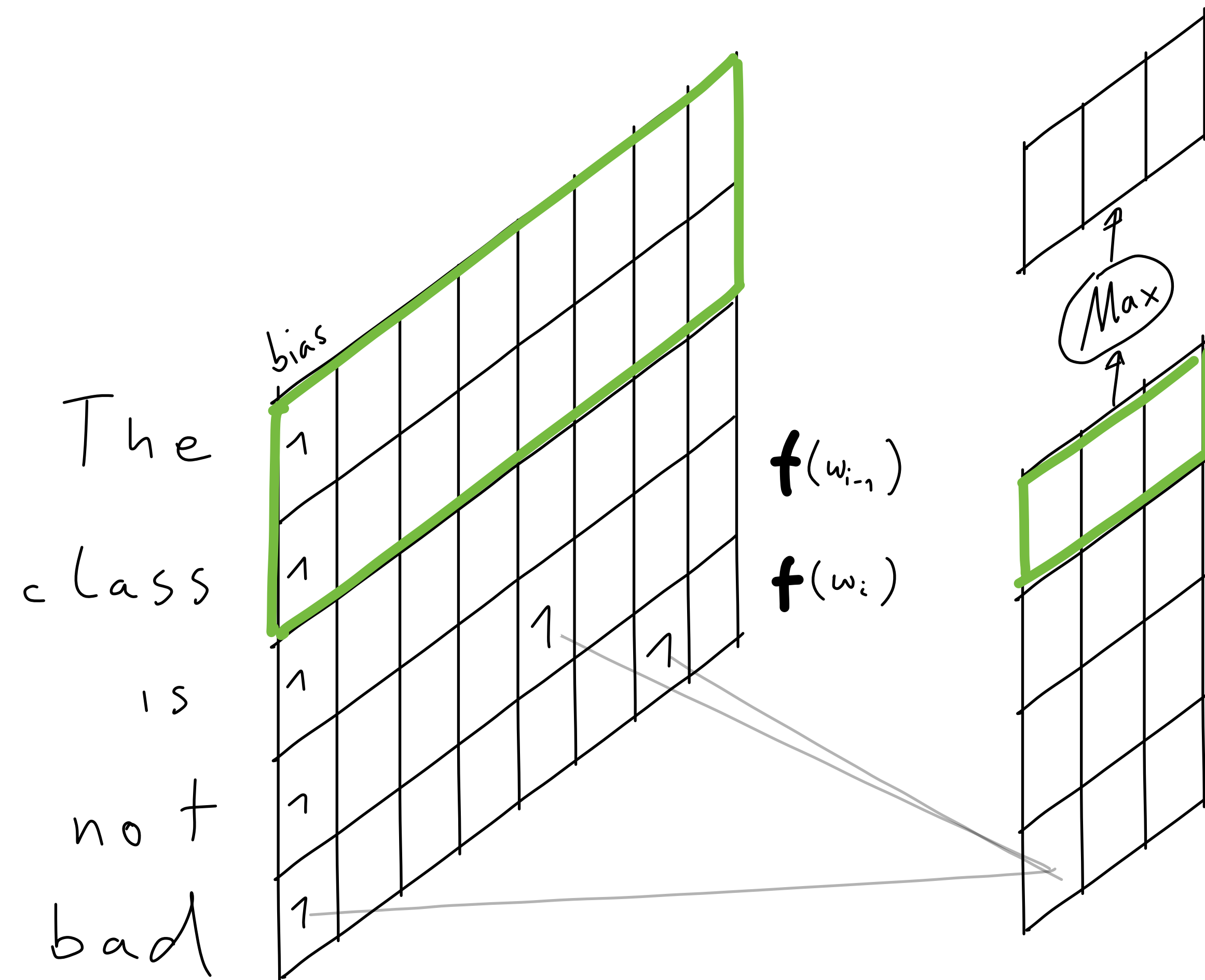
Can the MLP learn this?

Convolutional Neural Networks



$$a_i = \sigma \left(w \begin{bmatrix} f(w_{i-1}) \\ f(w_i) \end{bmatrix} \right)$$

Convolutional Neural Networks

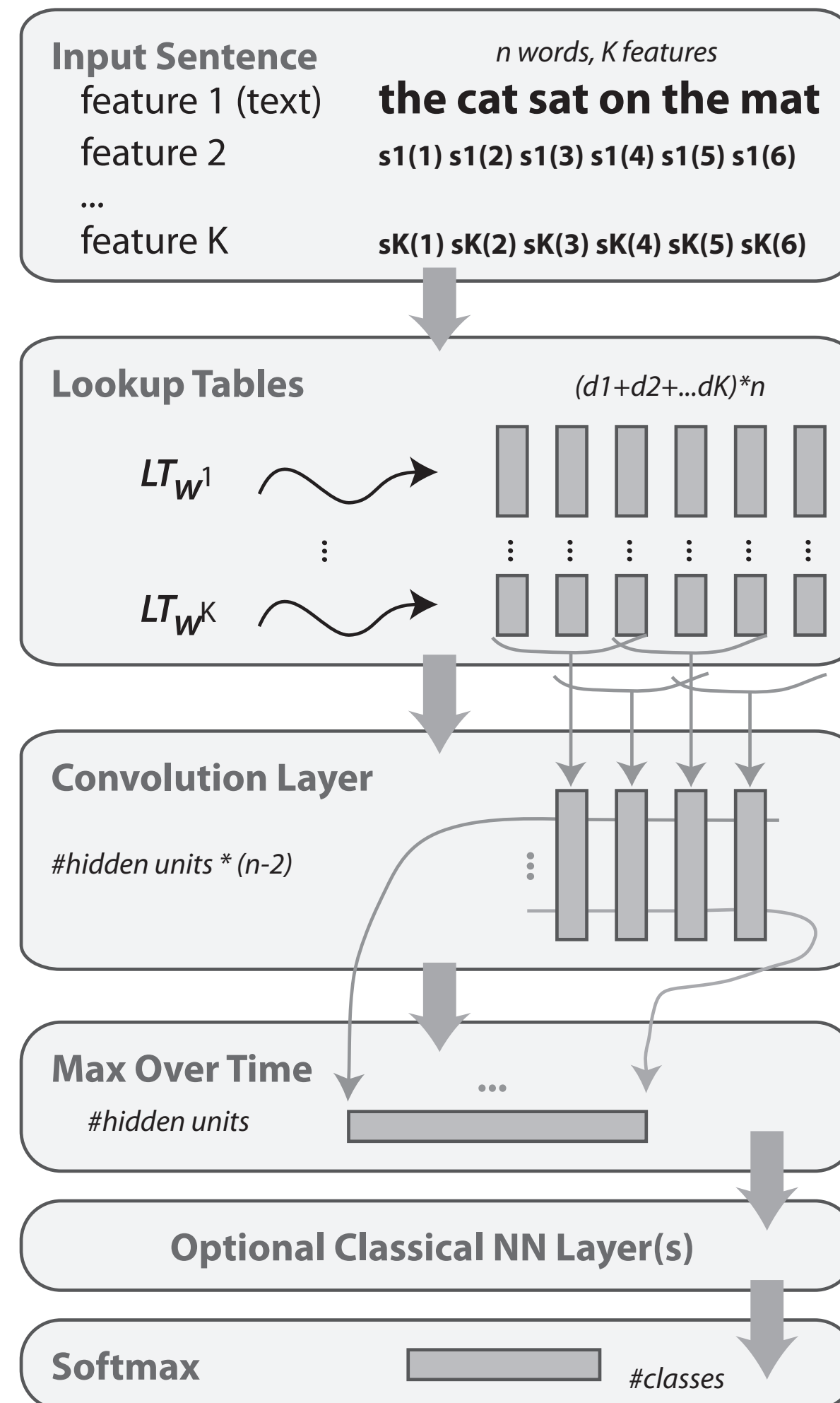


$$a = \max_i a_i$$

$$a_i = \sigma \left(w \begin{bmatrix} f(w_{i-1}) \\ f(w_i) \end{bmatrix} \right)$$

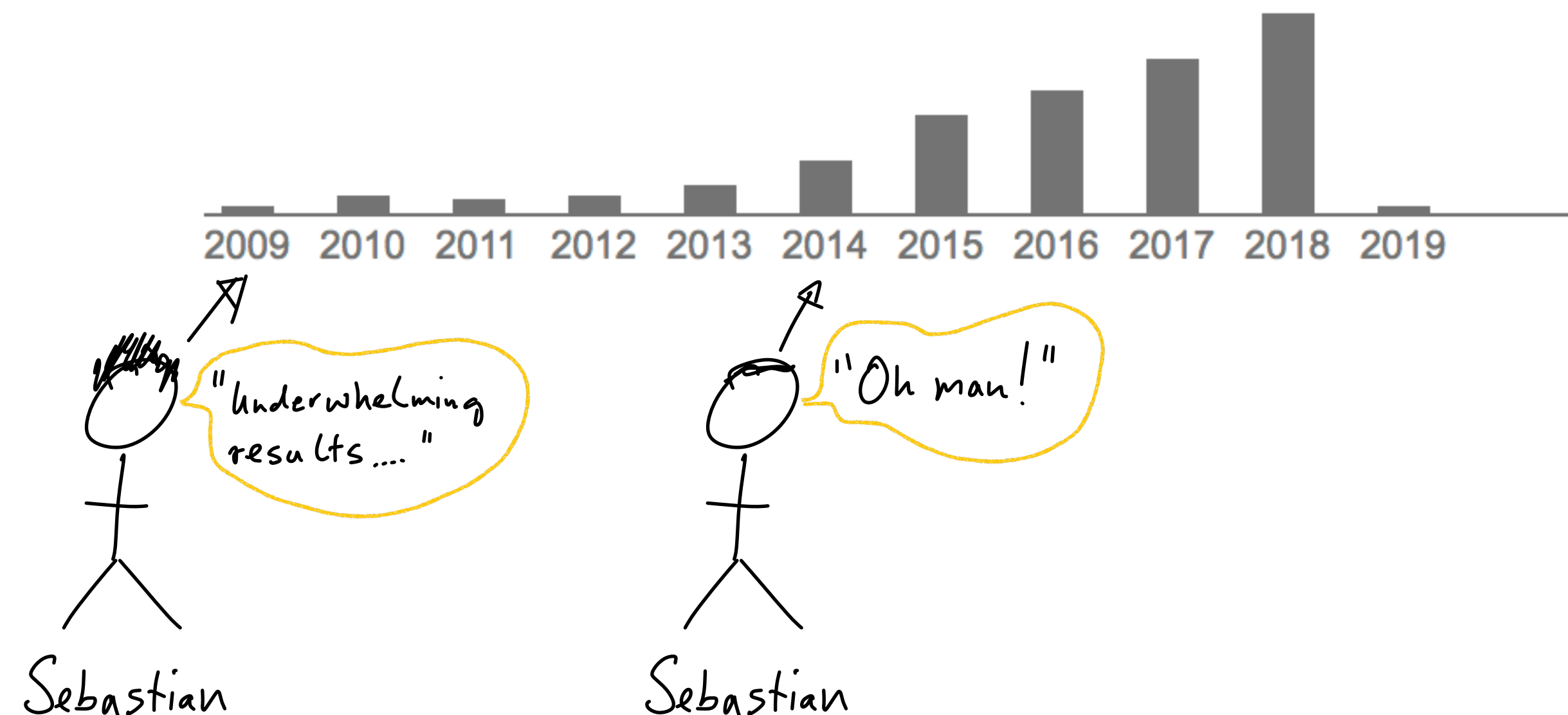
A Unified Architecture for Natural Language Processing

Collobert & Weston 2008



Encoder/Decoder?

Total citations Cited by 2974



Final Words



On backpropagation:

“My view is throw it all away and start again.”

“The future depends on some graduate student who is deeply suspicious of everything I have said.”

Geoffrey Hinton

“Godfather of Deep Learning”

References

- A Unified Architecture for Natural Language Processing, Collobert & Weston, ICML 2008
- Natural Language Processing (Almost) from Scratch, Collobert et al., 2011, JMLR
- Goldberg, Chapter 4 & 6