

Deep learning, neural networks and frameworks

The Academy of AI 2018/19

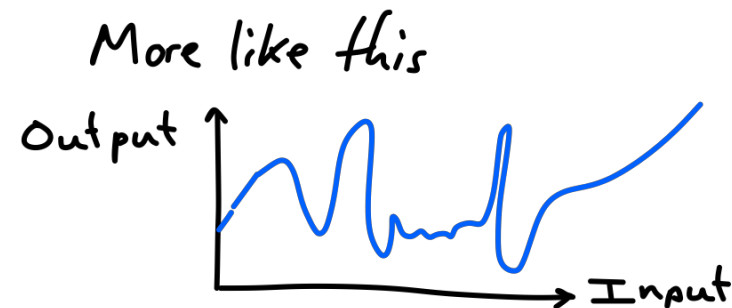
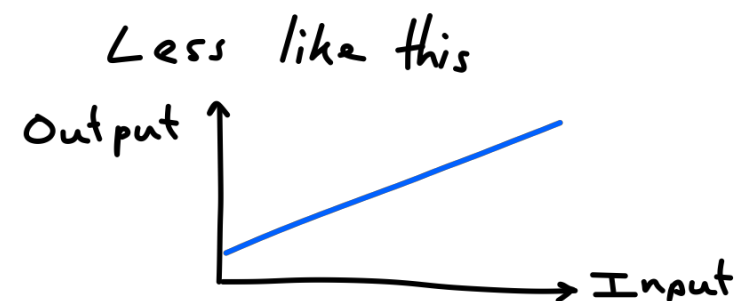
Session 7

Students for AI

Problems motivating deep learning

Modelling complex functions

Useful input-output relationships are too complex to guess what model might fit them



sin wave?
high order polynomial?
gaussian?

Assuming any particular function as the model makes it biased

How can we hope to guess a function that can learn a sufficient model?

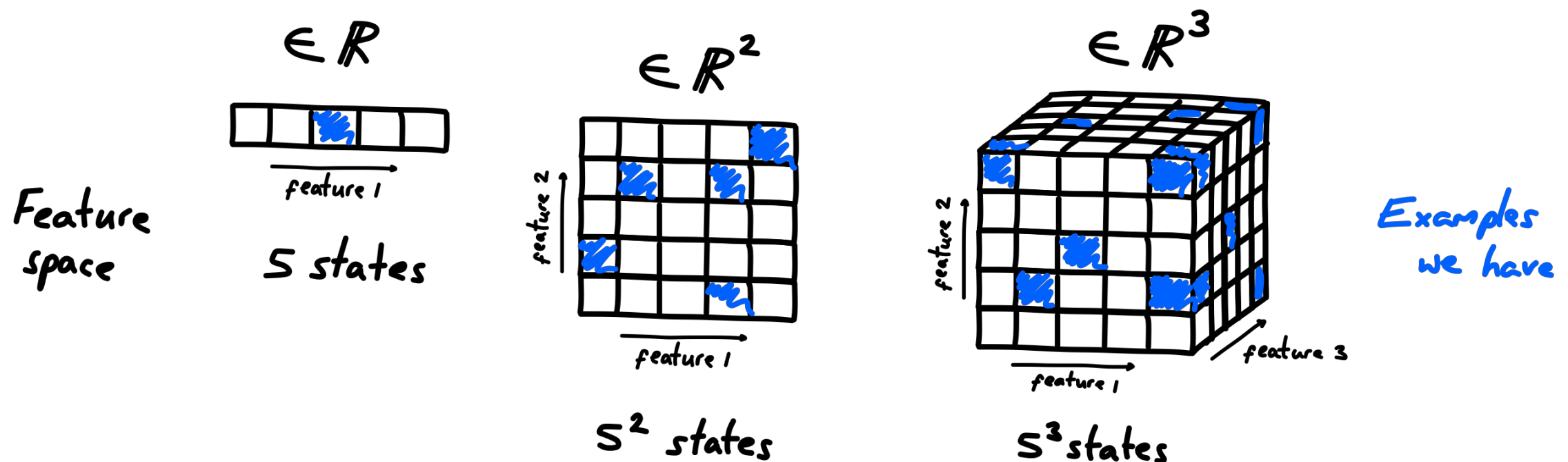
We need our model to be able to fit any shape

Problems motivating deep learning

The curse of dimensionality

As the number of dimensions (features) of a problem increase, there become exponentially more possible states that each example can take.

With more possible states, we need to see more examples to cover the space so that we can form a good representation of how an example at that point is mapped to an output.



Problems motivating deep learning

Abstraction

Us humans we don't pay attention to every individual raw input.

Low level feature

Each photon that hits your retina

Frequency and amplitude of
sound over time

Locations and pressures at
each point on skin

High level feature

"I see a car"

"I hear a voice"

"I feel the apple in my hand"

abstraction



Forming abstract, hierarchical representations of low level data alleviates the problem of the curse of dimensionality.

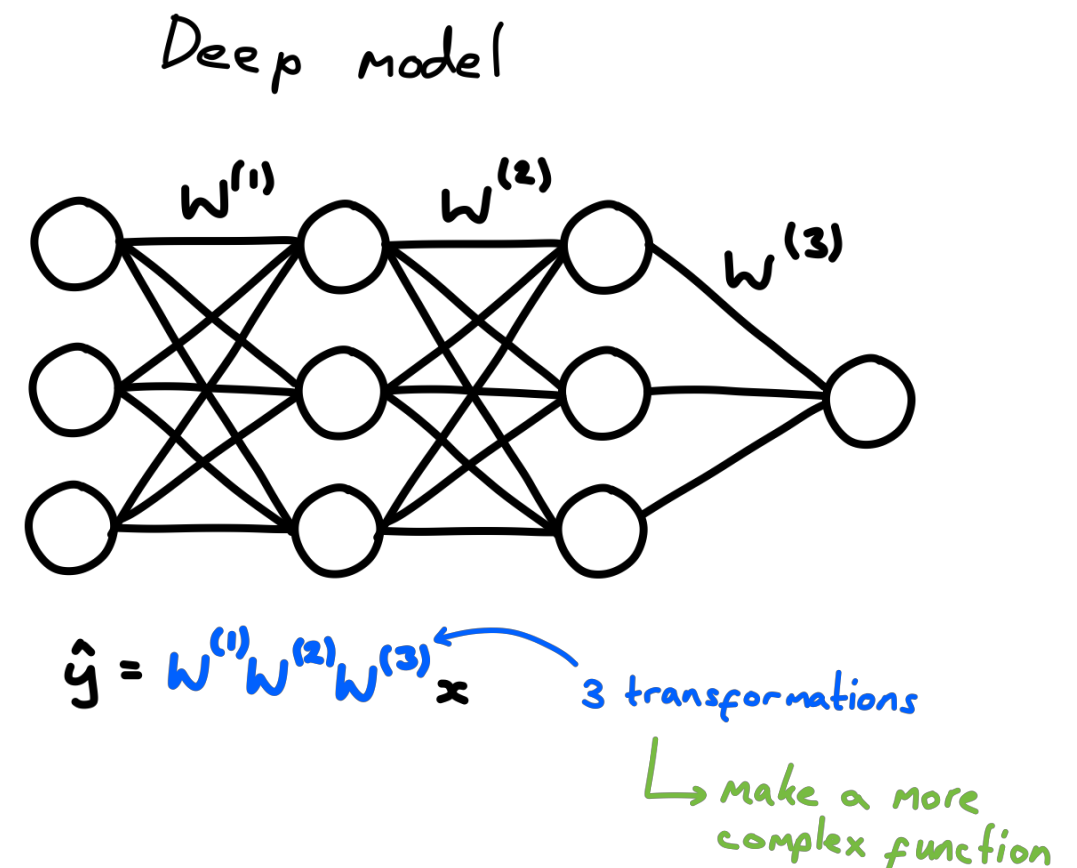
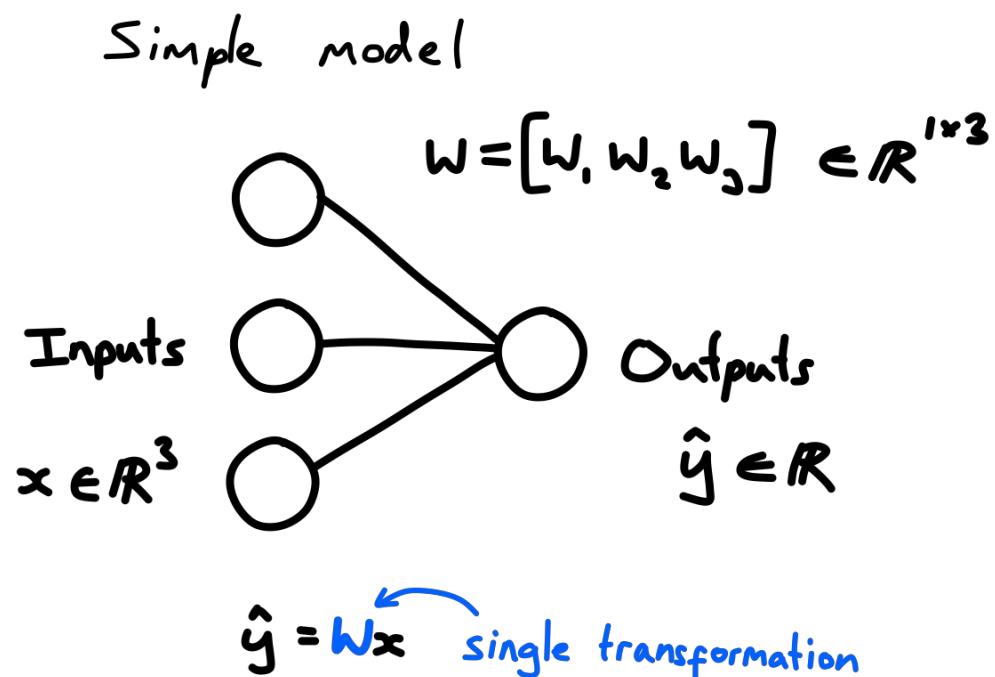
As such, we want our models to be able to make abstractions

Modelling complex transformations

The models that we've seen so far only apply one function to transform the data

But we need to be able to model non-linear relationships

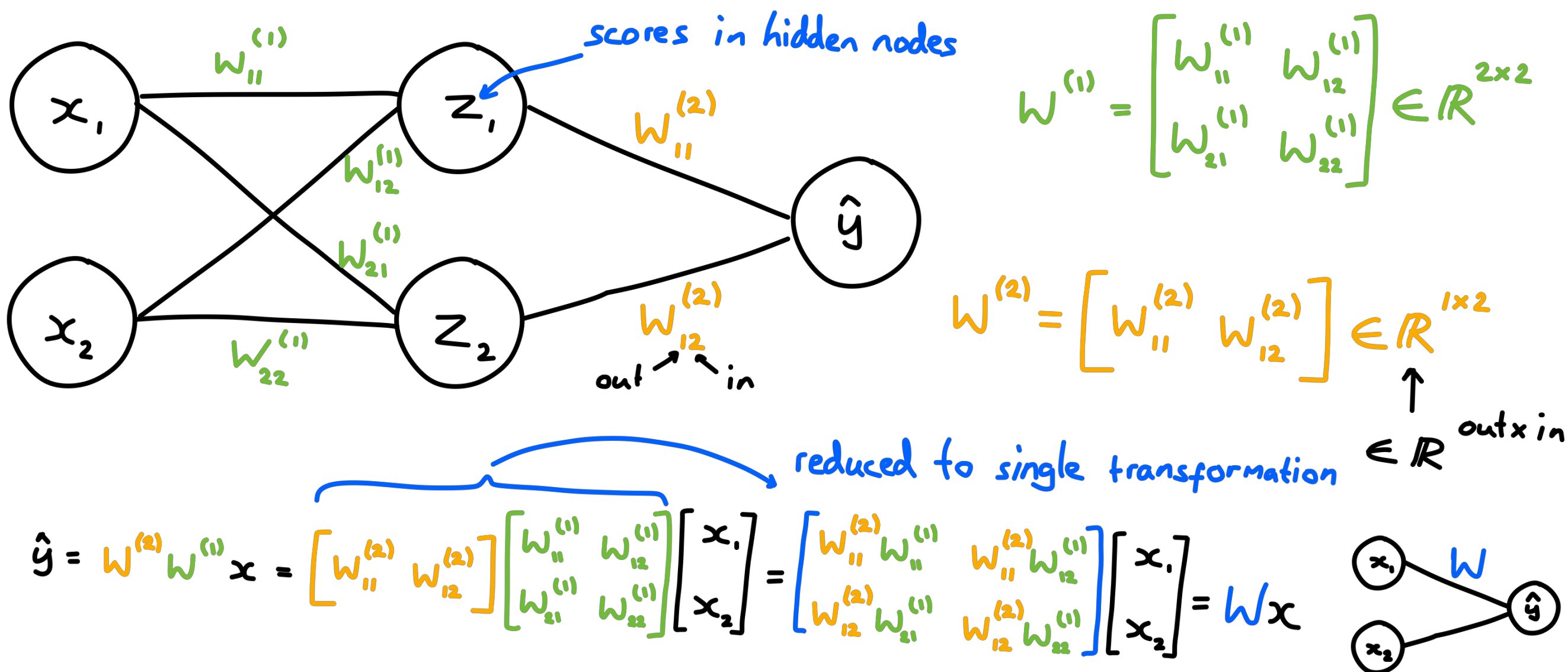
We can represent complex functions by repeatedly applying simple functions.



Linear factorisation

However...

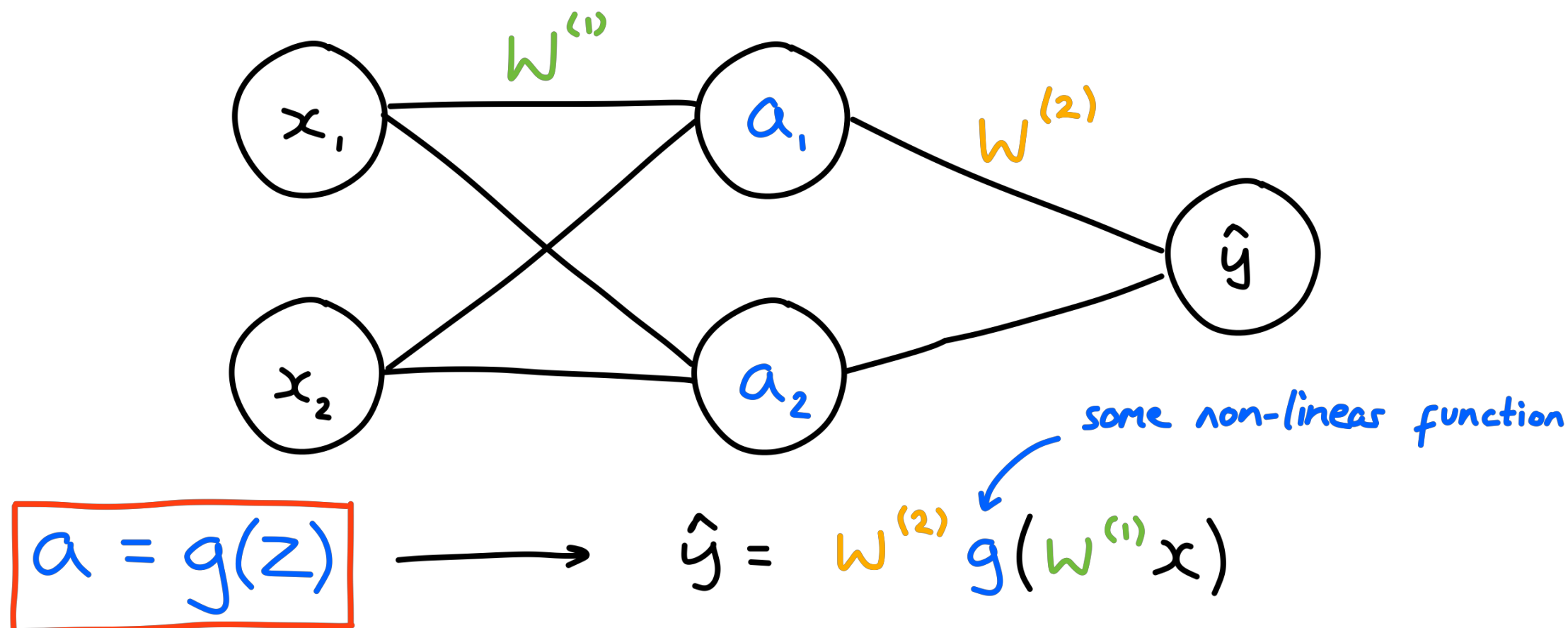
...repeatedly applied linear transformations are equivalent to a single linear transformation



In order to prevent this happening, we use **activation functions**...

Activation functions

Activation functions are **non-linear** transformations that we apply to the input of each hidden node

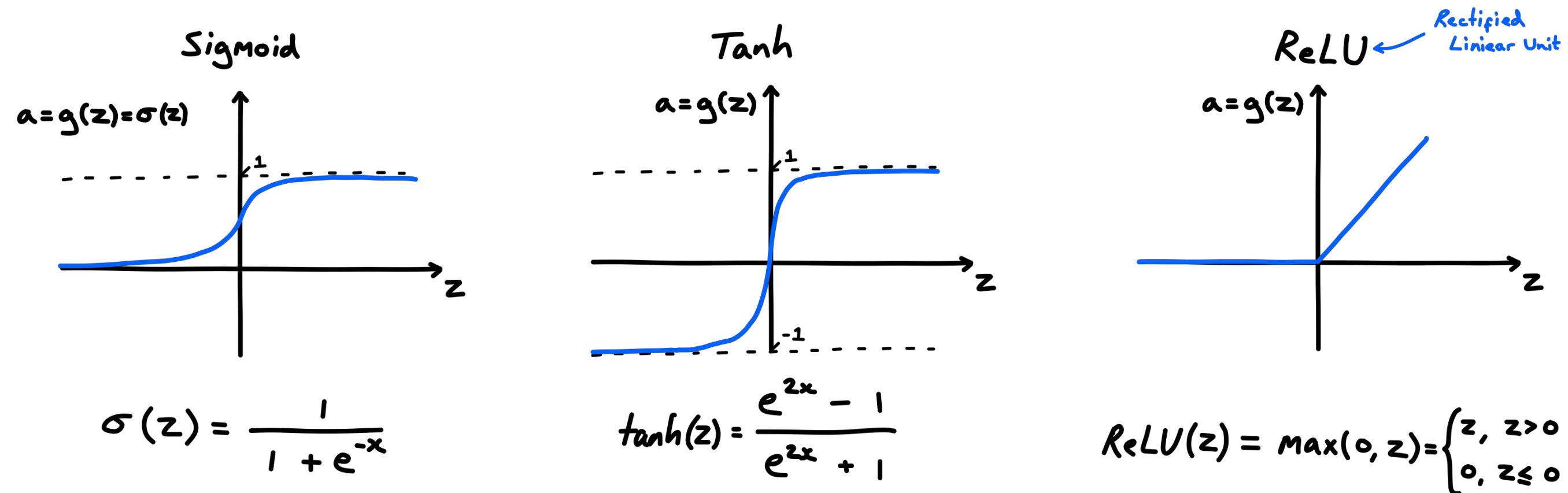


The function can now no longer be reduced to a single linear transformation

We have modelled a non-linear function

Activation functions

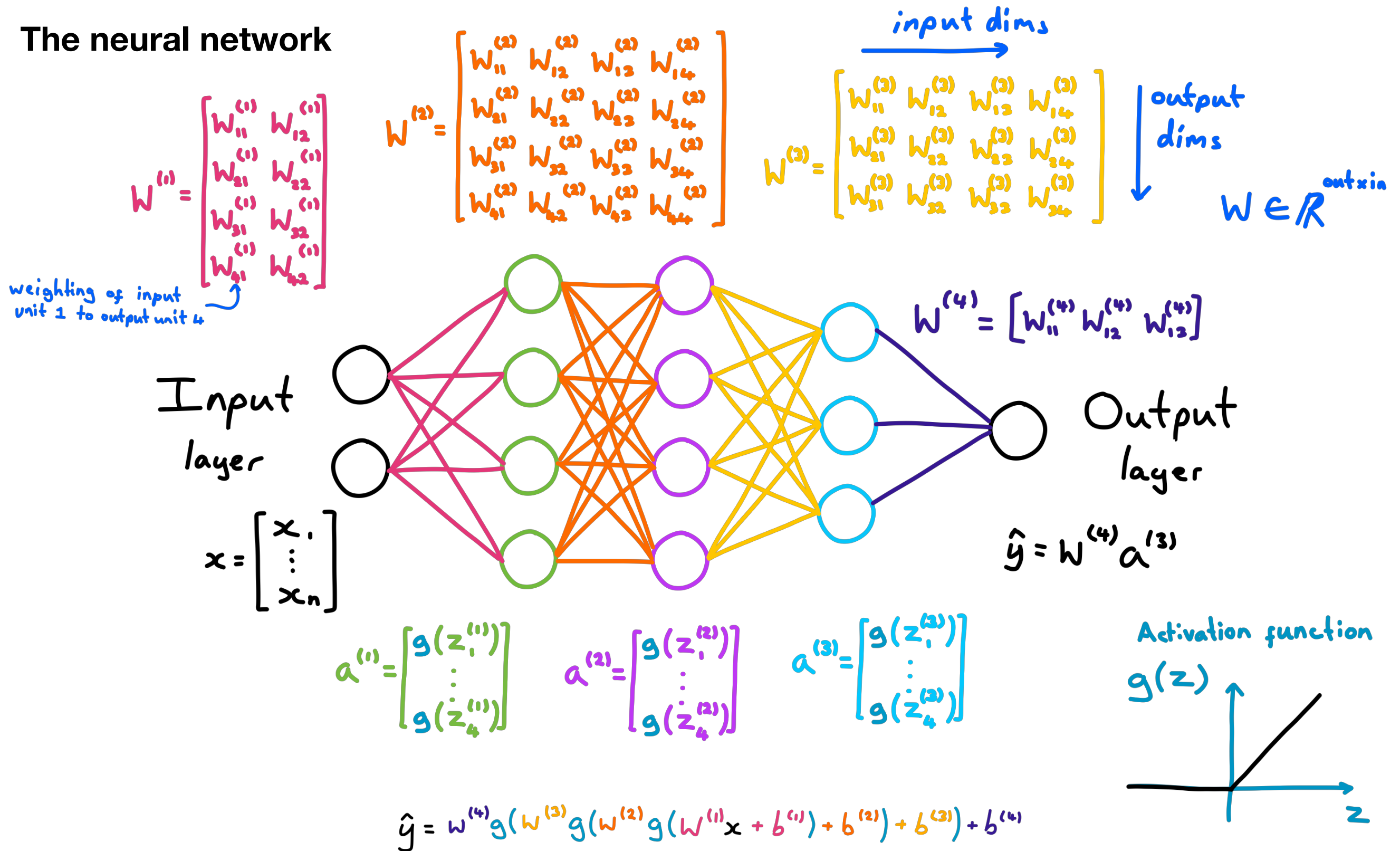
What functions do we use as activation functions?



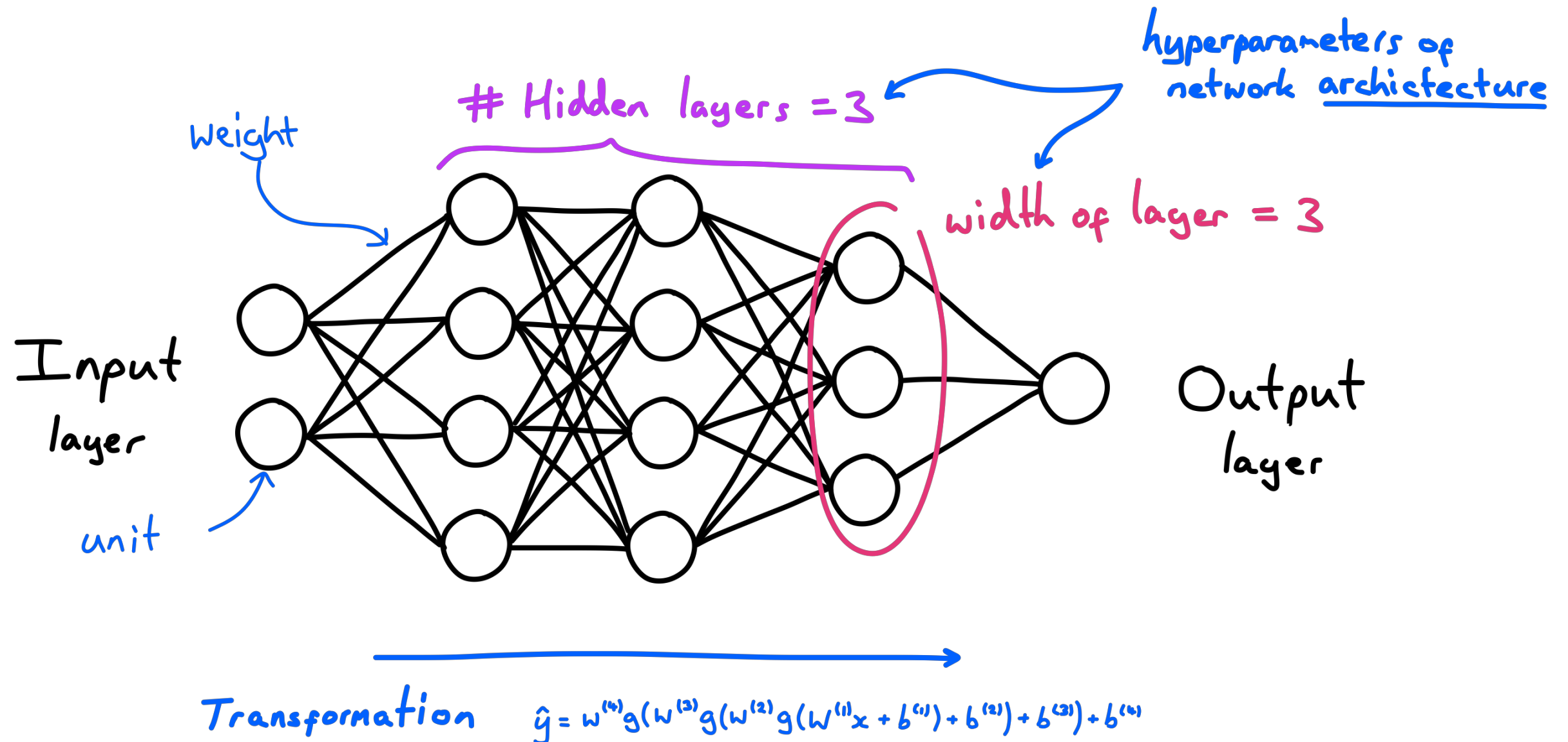
All an activation function needs to do, is to add nonlinearity so that we can learn nonlinear mappings

Activation functions that perform comparably to the best ones we know are so common as to be uninteresting

The neural network



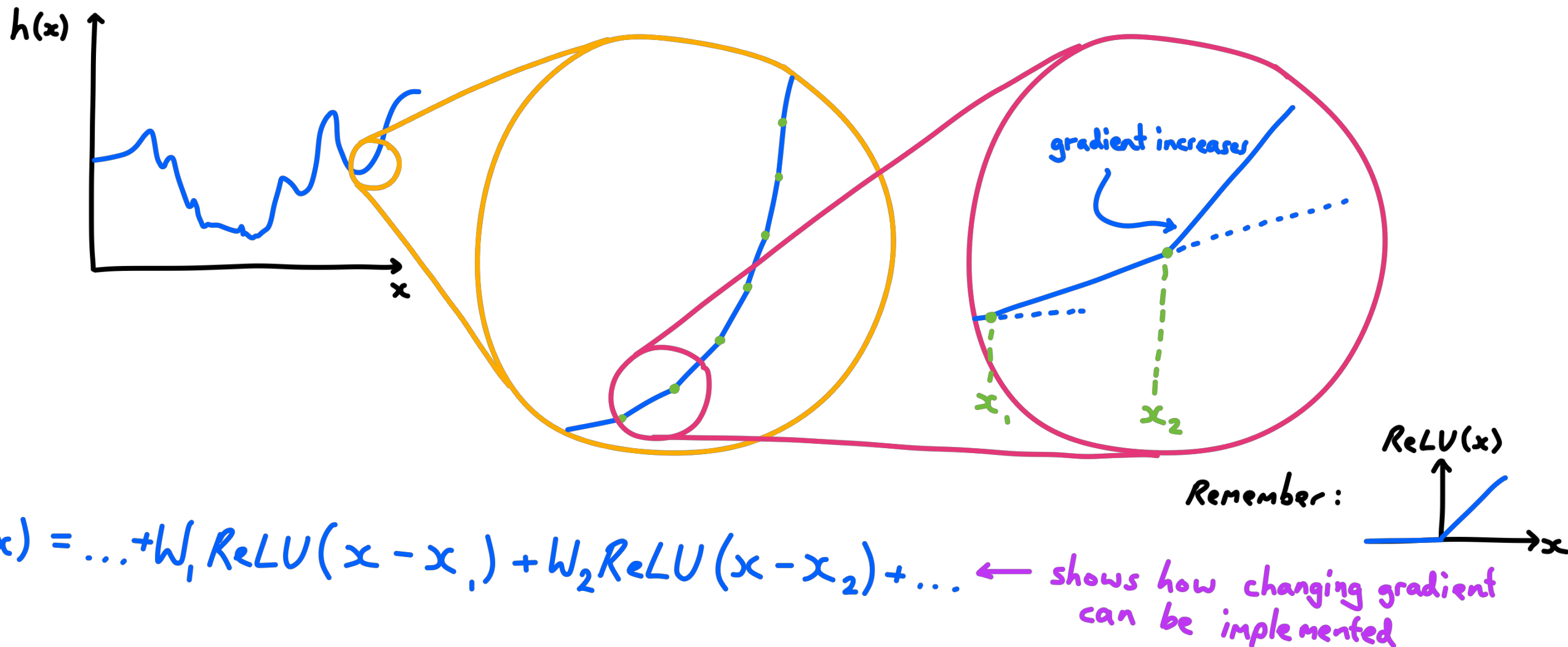
The neural network



Neural networks can be trained to model almost ANY function, no matter how ridiculous...

...but they need to be deep and wide enough to do so.

Simple demonstration of universal approximation



Each ReLU adds a linear contribution to the model

This contribution is not added over the whole domain, and hence the model is piece-wise linear

With an increasing number of hidden units, the model becomes increasingly continuous

Hierarchical feature learning

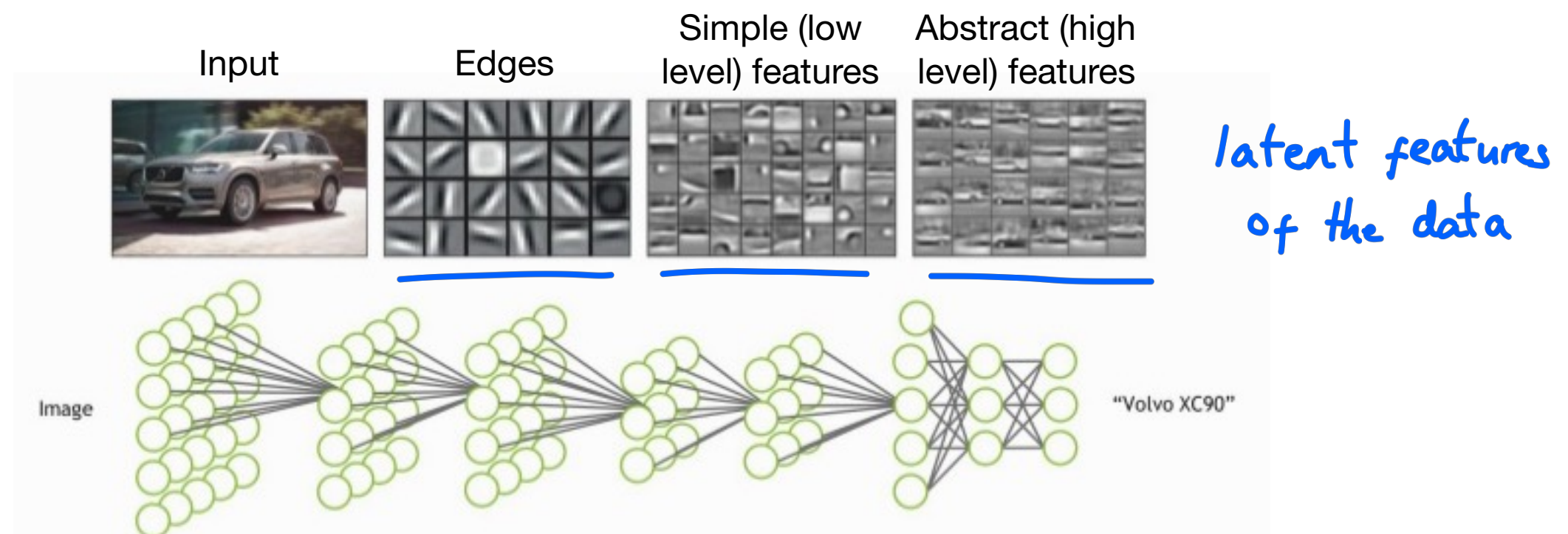
The output of each layer is latent features of the input.

Latent features means features that we don't observe directly

E.g. the edges in an image rather than the raw pixels

Because each layer builds on the features output by the previous layer, each layer can learn to output more abstract features than those input.

This allows neural networks to learn **hierarchical representations**



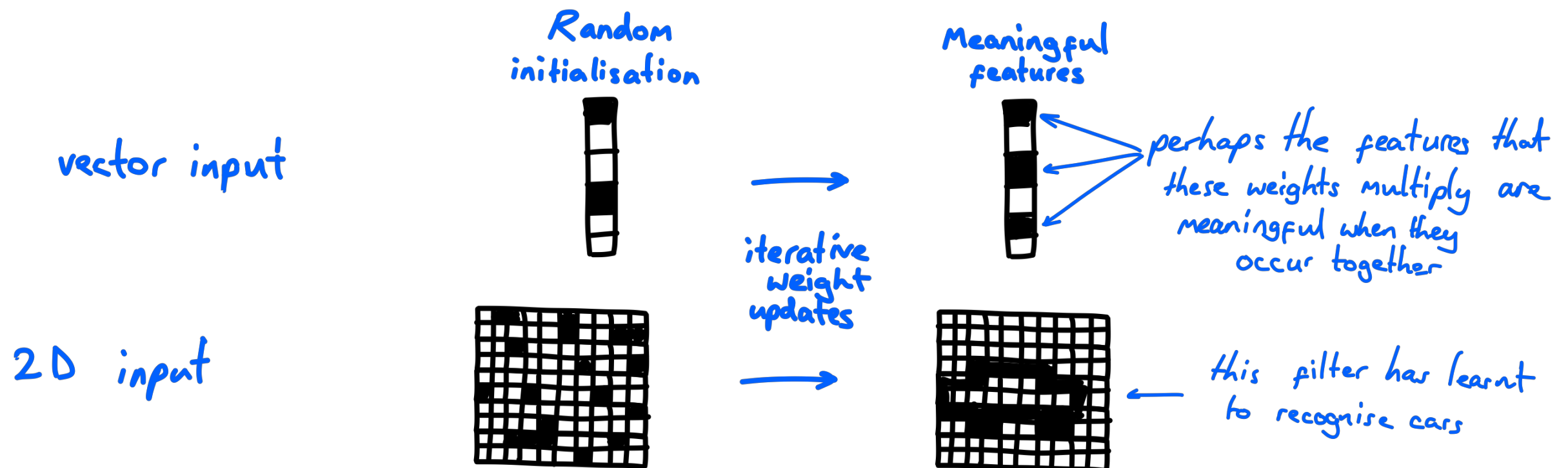
Automatic feature extraction

Gradient based optimisers (like SGD) update model parameters in a way that reduces the loss (improves performance) for the input examples

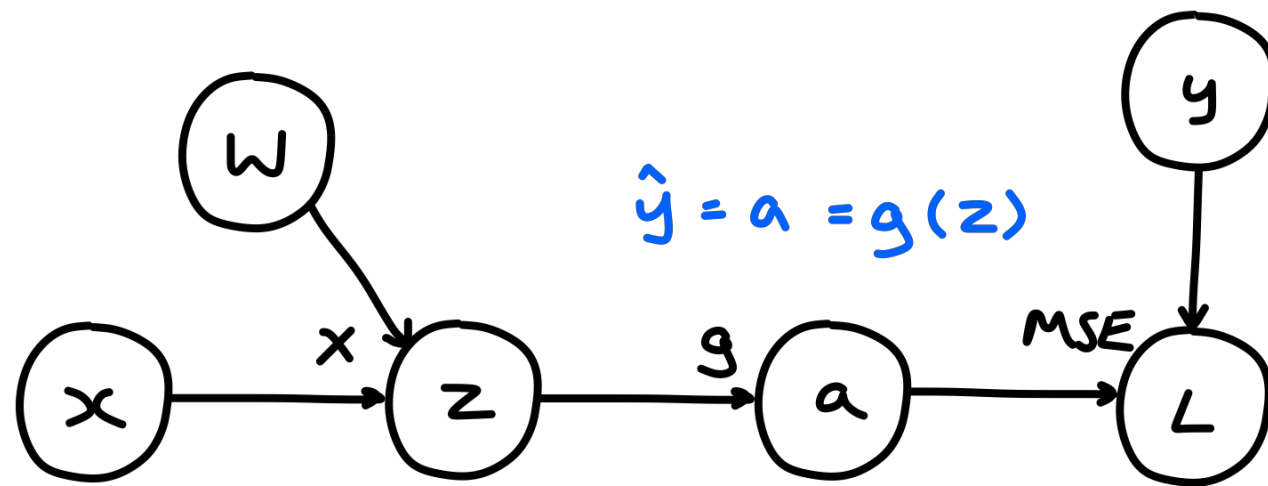
Intuitively, a model will perform better if it can identify meaningful features

So as the model improves, it does so by learning weights that activate units in the network when **meaningful features** are present in the input.

This is called **automatic feature extraction**



Computational graphs



This is a simple computational graph of a logistic regression model.

It shows the relationships between variables that appear in the model and makes it easy to compute derivatives, which we need for gradient based learning.

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial g} \frac{\partial g}{\partial z} \frac{\partial z}{\partial w}$$

We can obtain the derivatives of any variable with respect to another using the chain rule

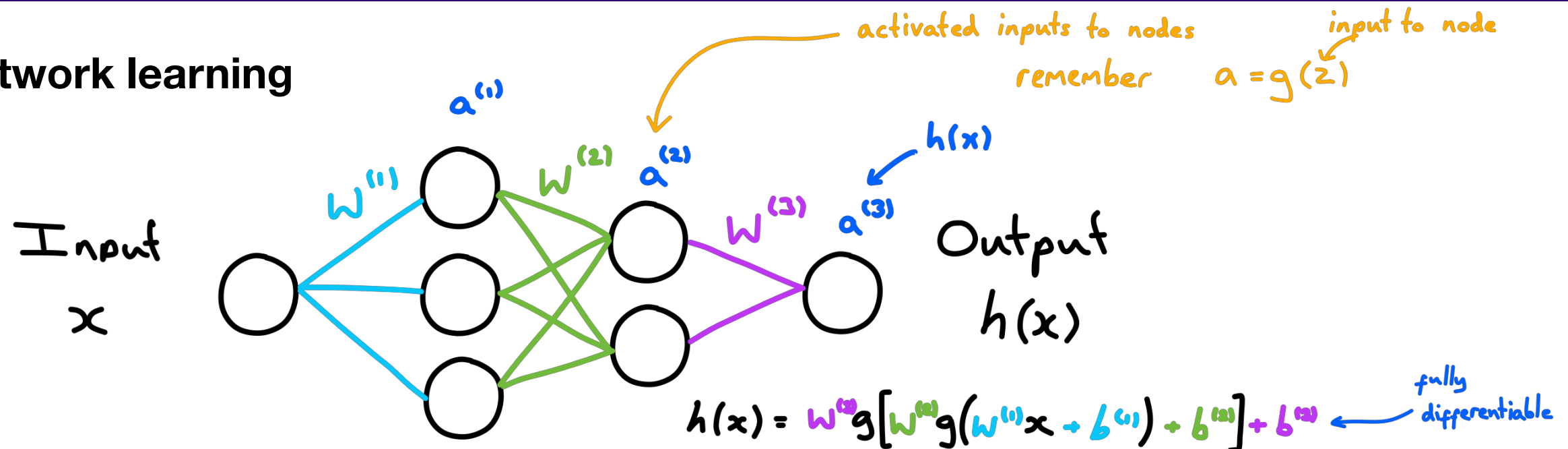
The derivative of each operation is simply a function of its inputs

$$\text{e.g. } a = g(z) = \frac{1}{1 + e^{-x}} \quad \frac{\partial a}{\partial z} = \frac{e^{-x}}{(1 + e^{-x})^2} = g(x)(1 - g(x))$$

So we can just program functions that output the derivative of each particular operation in the graph.

Then apply them to any variable that appears in our graph, as long as we store how they were computed and from what variables.

Neural network learning



To learn, we need to find how each weight affects the error (find the derivatives)

$$\frac{\partial L}{\partial w^{(3)}} = \frac{\partial L}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial w^{(3)}}$$

$$\frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial w^{(1)}}$$

Parts of the equation reappear

so we can reuse them to save unnecessary computation

We can propagate the error backwards through the model to find the gradients required for gradient based optimisers to make it learn

This is called **backpropagation**