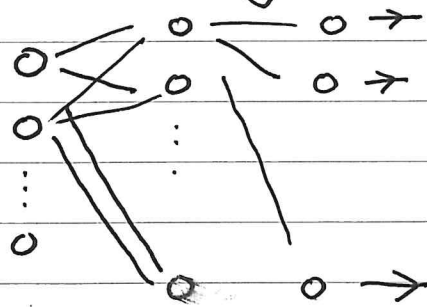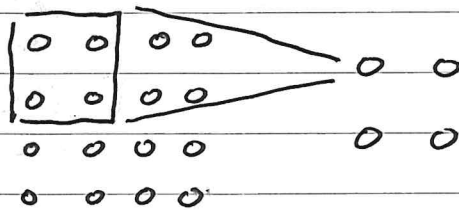Artificial
Neural Networks (ANN)

Various types of ANNs are
used for many different ML tasks:

- Feed-forward neural networks (FFNN)
  (a.k.a multilayer perceptron although
  not necessarily step-function activation)



- Convolutional NN (CNN)
                            - not fully connected
                            - spatial structure
                              of input

                            - reduced # weights
                              e.g. image
                                recognition

- Recurrent NN (RNN)

    - not just feed-forward

    - output depends on previous cycles
      => sequential information

    - used for e.g. text and speech
      recognition

- Other types; also for unsupervised
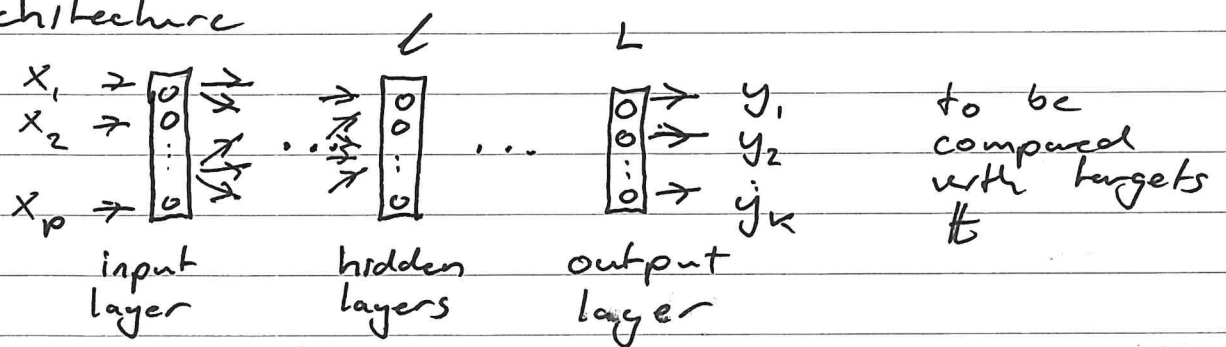  learning. Deep Boltzmann Machines

## FFNN

Can be described by its

1. Architecture
2. Activation rule(s)
3. Learning algorithm

Universal approximation theorems imply
that ANN's are able to represent various
functional relationships.

E.g. a FFNN with a single hidden
layer and a finite # of neurons
can approximate a continuous
multidim. function to arbitrary
accuracy assuming that the
activation function is non-constant,
bounded, monotonically increasing, continuous.
( => non-linear )'

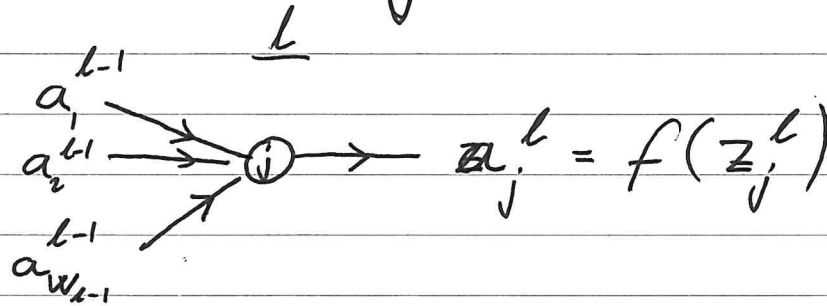1/ Architecture



| input | hidden | output |
| layer | layers | layer |

Each layer has a number of
nodes/units/neurons, each with its
set of weights and bias.
The width $W_l$ of layer $l$

## 2/ Activation rule

consider a single neuron $j$ in hidden Layer $\ell$



outputs from $W_{\ell-1}$ nodes in layer $\ell-1$

Note that the output $a_j^\ell$ becomes an input to all nodes in layer $(\ell+1)$

The activation $$z_j^{\{\ell\}} = \sum_i w_{ij}^{\{\ell\}} a_i^{\ell-1} + b_j^\ell$$

or in matrix form

$$z^\ell = \left[ w^\ell \right]^T a^{\ell-1} + b^\ell$$

$[W_\ell \times 1]$     $\underbrace{W_{\ell-1} \times W_\ell}$     $\underbrace{W_{\ell-1} \times 1}$     vector $[W_\ell \times 1]$

several different activation functions can be considered

- Sigmoid
(logit)

$$f(z) = \sigma(z) \equiv \frac{e^z}{1+e^z}$$

· suffers from vanishing
gradients => slow learning



- Hyperbolic tangent

$$f(z) = \tan h(z) = 2\sigma(2z) - 1$$

Both of these imply that signals
will be non-zero everywhere
=> inefficient.
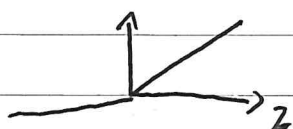
- Rectified linear unit (ReLU)

$$f(z) = \max(0, z)$$



Many neurons will
be quiet. Good!
But $\frac{\partial f}{\partial w} = 0$ if $z < 0$

=> dying ReLUs

=> Leaky ReLU



$$f(z) = \max(0, z) + \alpha \min(0, z)$$
with $\alpha \sim 0.01$

or Exponential Linear Unit (ELU)

$$f(z) = \begin{cases} z & \text{for } z \geq 0 \\ \alpha e^z - 1 & \text{for } z \leq 0 \end{cases}$$

Note that the output layer often has a linear activation function to give a continuous output, or softmax to give classification probabilities.

The propagation of signals through the ANN is known as feed-forward

3) Learning algorithm

- Model (cost function)
    MSE, MAE, cross-entropy
    i.e. how to compare output with targets. (Physics knowledge can be incorporated.
- Regularization (priors)

- Back-propagation (using the chain rule to obtain the gradients to adjust the parameters)

- Gradient descent
    ∘ standard SGD
    ∘ momentum -11-
    ∘ AdaGrad
    ∘ RMS prop
    ∘ Adam

- Data
  - training data
  - ~~test data~~
  - validation data (adjust hyperpars)
  - test data (final test of performance. <u>cannot</u> be used for tuning)

Since learning is stochastic we typically feed the training data many times (feed-forward + backprop. One complete pass is known as an <u>epoch</u>.

The data is often split into <u>batches</u> doing the gradient for a batch of data.

The validation score is evaluated for each ~~bat.~~ epoch.

Tune # of epochs (eventually overfitting)
batch size
learning rate (monitor both training and val. scores)

Conclusion
  o Scary many options...

  o A physicist also wants to learn about the model itself!

# Backpropagation algorithm

Backprop: Using the chain rule
to obtain the gradient
of the cost function
w.r.t weights and biases.

Learning: E.g. with SGD; training
of weights and biases

## Regression example

The ANN gives $n$ output $\{y^{(i)}\}_{i=1}^{n}$
to be compared with targets $\{t^{(i)}\}_{i=1}^{n}$.

The output is the output of the
last layer $(L)$: $y^{(i)} = a_i^L$

And the cost function

$$C(\theta) = \frac{1}{2} \sum_{i=1}^{n} \left( a_i^L - t^{(i)} \right)^2$$

depends on all weights and biases
of the network.

Consider node $j$ in hidden layer $l$

$$z_j^l = \sum_i w_{ij}^l a_i^{l-1} + b_j^l \qquad \text{(activation)}$$

$$a_j^l = f(z_j^l) = \frac{1}{1 + e^{-z_j^l}} \qquad \left(\begin{array}{l}\text{output;} \\ \text{Logistic} \\ \text{activation}\end{array}\right)$$

Chain rule

$$\frac{\partial z_i^{l}}{\partial w_{ij}^{l}} = a_i^{L-1} \quad ; \quad \frac{\partial z_j^{l}}{\partial a_i^{L-1}} = w_{ij}^{l}$$

$$\boxed{\frac{\partial a_j^{l}}{\partial z_j^{l}} = \frac{\partial f(z_j^{l})}{\partial z_j^{l}} = \left\{ \begin{array}{c} \text{logit} \\ \text{function} \end{array} \right\} = -\frac{-e^{-z_j^{l}}}{(1+e^{-z_i^{l}})^2}}$$

$$= \frac{(1+e^{-z_i^{l}})-1}{(1+e^{-z_i^{l}})^2} = \boxed{a_j^{l}(1-a_j^{l})}$$

Derivatives of the cost function

$$\frac{\partial C}{\partial w_{jk}^{L}} = (q^{L} - t^{(j)}) \frac{\partial a_j^{L}}{\partial w_{jk}^{L}}$$

$$= \frac{\partial a_j^{L}}{\partial z_j^{L}} \frac{\partial z_j^{L}}{\partial w_{jk}^{L}}$$

$$= a_j^{L}(1-a_j^{L}) a_k^{L-1}$$

Let us define

$$\delta_j^{L} \equiv a_j^{L}(1-a_j^{L}) \left[ a_j^{L} - t^{(j)} \right]$$

$$= f'(z_j^{L}) \frac{\partial C}{\partial a_j^{L}} = \frac{\partial C}{\partial z_j^{L}}$$

which gives

$$\boxed{\frac{\partial C}{\partial w_{jk}^{L}} = \delta_j^{L} \cdot a_k^{L-1}}$$

Also

$$\boxed{\frac{\partial C}{\partial b_j^{L}} = \left[a_j^{L} - t^{(j')}\right] \frac{\partial a_j^{L}}{\partial b_j^{L}}}$$

$$= \frac{\partial a_j^{L}}{\partial z_j^{L}} \underbrace{\frac{\partial z_j^{L}}{\partial b_j^{L}}}_{= 1}$$

$$= \left[a_j^{L} - t^{(j')}\right] a_j^{L}(1-a_j^{L}) = \underline{\boxed{\delta_j^{L}}}$$

Now, consider <u>Layer $\ell$</u>

$$\delta_j^{L} = \frac{\partial C}{\partial z_j^{\ell}}$$

can be defined in terms of Layer $(\ell+1)$:

$$\delta_j^{\ell} = \sum_k \frac{\partial C}{\partial z_k^{\ell+1}} \frac{\partial z_k^{\ell+1}}{\partial z_j^{\ell}}$$

$$= \sum_k \delta_k^{\ell+1} \frac{\partial z_k^{\ell+1}}{\partial z_j^{\ell}}$$

where

$$z_j^{l+1} = \sum_{i=1}^{w_l} w_{ij}^{l+1} \underbrace{a_i^l}_{} + b_j^{l+1}$$

such that $= f(z_j^l)$

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} f'(z_j^l)$$

Finally, we can compute all the $\delta_j^l$ quantities from the current set of weights and outputs — starting from the last layer and propagating backwards!

The SGD learning (with rate $\eta$) then implies

$$w_{jk}^l \leftarrow w_{jk}^l - \eta \, \delta_j^l \, a_k^{l-1}$$

$$b_j^l \leftarrow b_j^l - \eta \, \delta_j^l$$