

Artificial Neural Networks

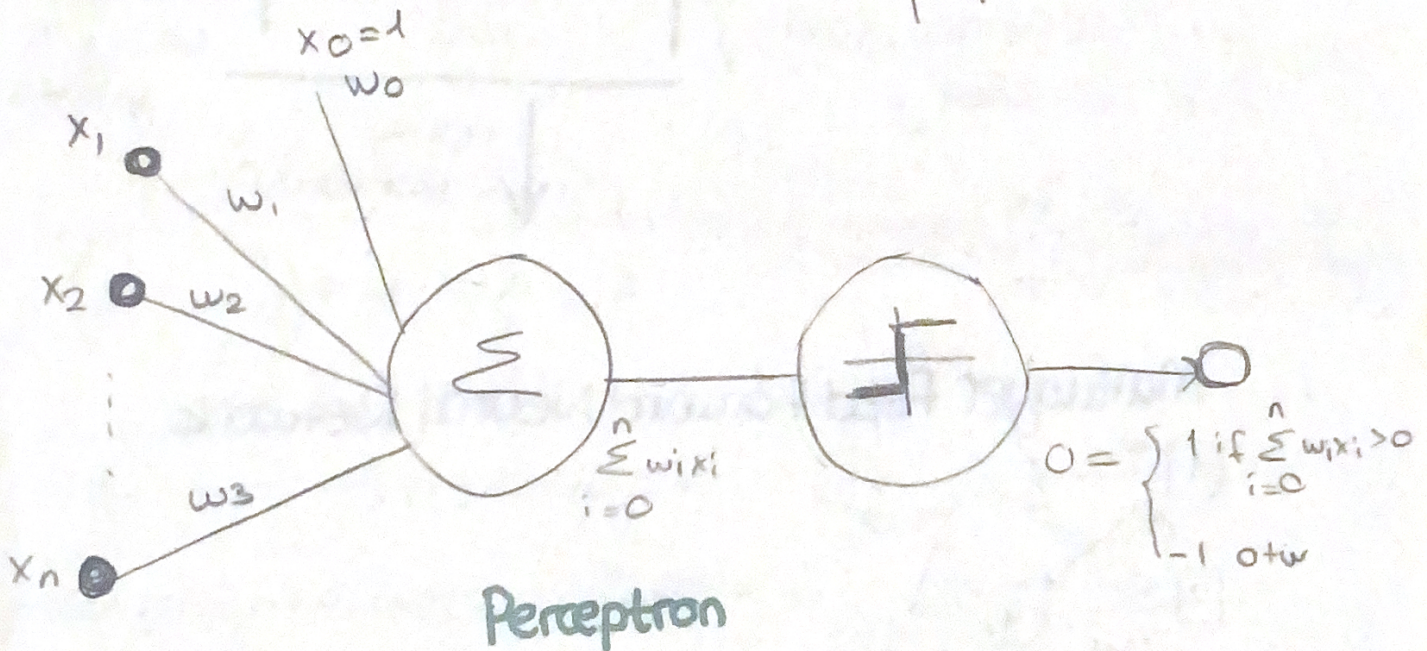
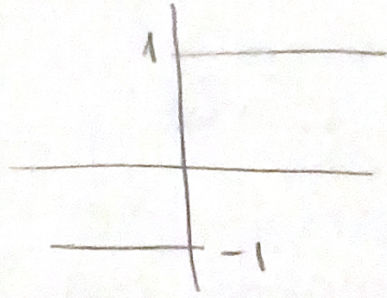
Perceptrons: Takes inputs, calculates a linear combination, prompts 1 if result > threshold & -1 if otherwise.

$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

inputs
output

↳ signum function (sgn)

- $w_0 \rightarrow$ threshold



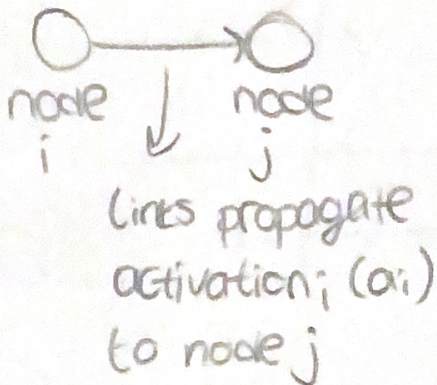
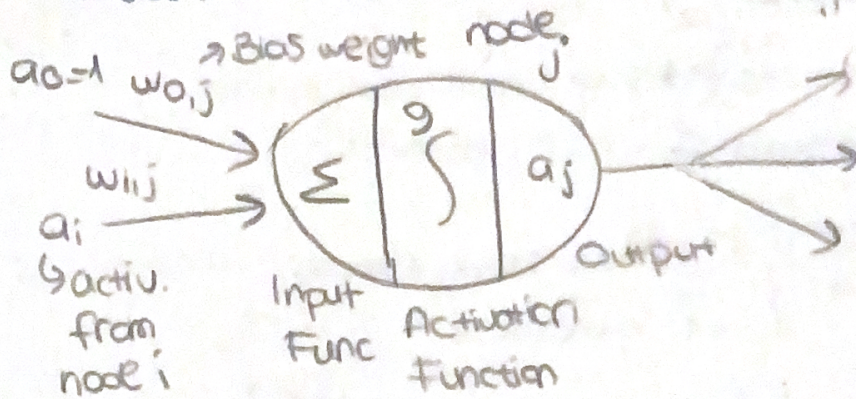
Weights for AND $\rightarrow w_0 = -0.8 \quad w_1 = w_2 = 0.5$

x_1	x_2	res
0	1	0.5 \rightarrow $< 0.8 = 0$
1	1	1 \rightarrow 1
0	0	0 \rightarrow $< 0.8 = 0$

OR $\rightarrow w_0 = 0.3 \quad w_1 = w_2 = 0.5$

x_1	x_2	res
0	1	0.5 \rightarrow $\geq 0.3 = 1$
1	1	1

Neural Networks



w_{ij} : weight of link i, j

a_0 : dummy (bias)

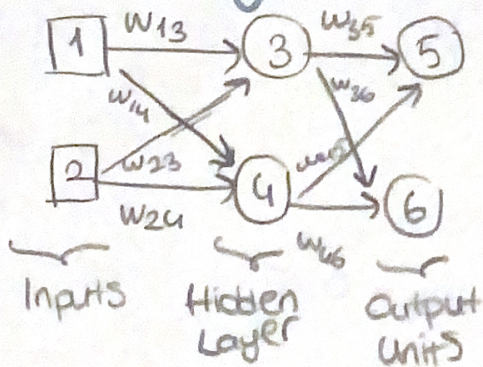
$w_{0,j}$: weight of bias

$$in_j = \sum_{i=0}^n w_{i,j} a_i \rightarrow \text{weighted sum of unit } j\text{'s inputs}$$

↓ apply activation

$$a_j = g(in_j) = g\left(\sum_{i=0}^n w_{i,j} a_i\right)$$

Multilayer Feed Forward Neural Network



$$a_5 = g(w_{0,5} + w_{3,5}a_3 + w_{4,5}a_4)$$

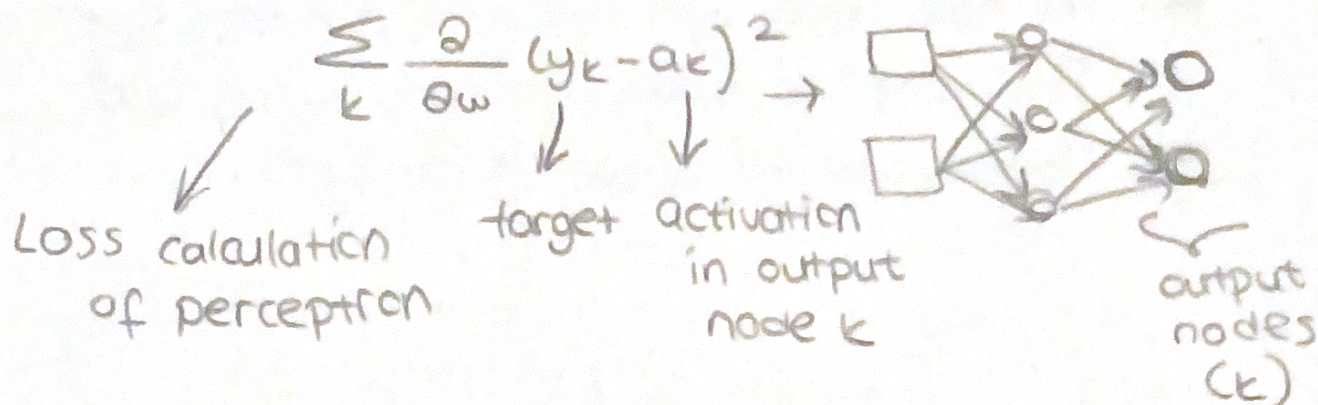
$$= g(w_{0,5} + w_{3,5}(g(w_{1,3}a_1 + w_{2,3}a_2 + w_{0,3}))$$

$$w_{4,5}(g(w_{1,4}a_1 + w_{2,4}a_2 + w_{0,4}))$$

↙ ↘
replace
with inputs
 x_1, x_2

Backpropagation

$$\frac{\partial}{\partial w} \text{Loss}(w) = \frac{\partial}{\partial w} |y - hw(x)|^2 = \frac{\partial}{\partial w} \sum (y_k - a_k)^2 =$$



- In Multilayer networks we take derivative of overall error gradient.

With multiple output units (k):

Err_k : k^{th} component of error vector $y - hw$

$\Delta_k = \text{Err}_k \times g'(in_k)$ } modified error

$w_{j,k} \leftarrow w_{j,k} + \alpha \times a_j \times \Delta_k$ } weight update rule

- To update connections between input units & hidden units we need to define a quantity analogous to the error term for output nodes. (where we need error backprop)
- Hidden node j is responsible for some fraction of Δ_k
- Δ_k values are divided according to strength of connection between hidden node & output node. They are propagated back to provide Δ_j for hidden layer.

→ Propagation rules

for Δ values : $\Delta_j = g'(in_j) \sum_k w_{j,k} \Delta_k$

- Backprop in a nutshell

1. Compute Δ values for output units, use observed error

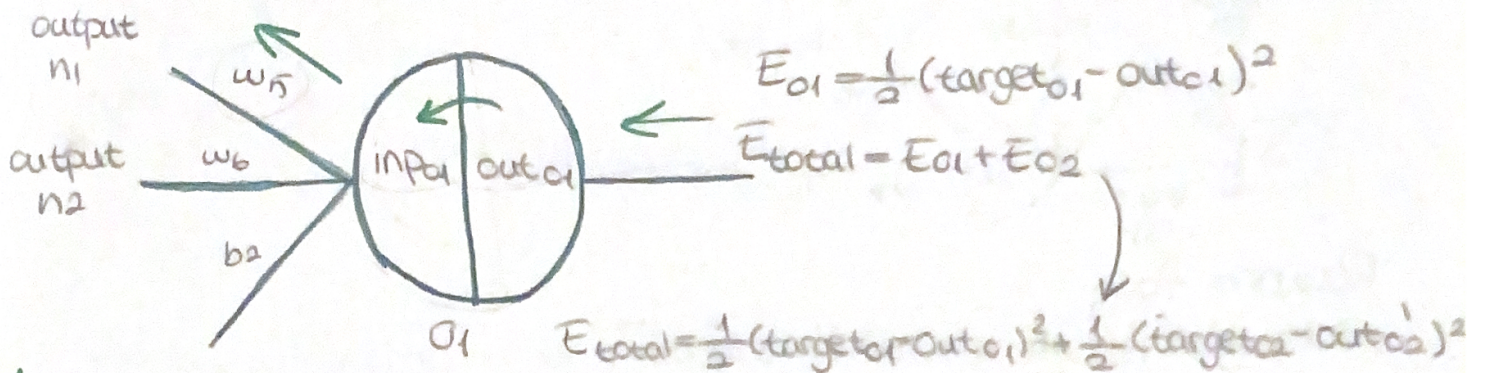
2. For each layer until the earliest hidden layer:

Propagate Δ values back to previous layer

Update weights between two layers

Backprop + Chain Rule Visual

$$\frac{\partial \text{inp}_{o1}}{\partial w_5} * \frac{\partial \text{out}_{o1}}{\partial \text{inp}_{o1}} * \frac{\partial E_{\text{total}}}{\partial \text{out}_{o1}} = \frac{\partial E_{\text{total}}}{\partial w_5}$$



$$1 \quad \frac{\partial E_{\text{total}}}{\partial \text{out}_{o1}} = 2 * \frac{1}{2} (\text{target}_{o1} - \text{out}_{o1})^{2-1} * (-1) + 0$$

$$= -(\text{target}_{o1} - \text{out}_{o1}) = \underline{\underline{0.741}}$$

$$2 \quad \frac{\text{out}_{o1}}{\text{inp}} = \frac{1}{1 + e^{-\text{inp}_{o1}}} \rightarrow \frac{\partial \text{out}_{o1}}{\partial \text{inp}_{o1}} = \text{out}_{o1} (1 - \text{out}_{o1}) = \underline{\underline{0.186}}$$

$$3 \quad \text{inp}_{o1} = w_5 * \text{out}_{n1} + w_6 * \text{out}_{n2} + b_2 * 1$$

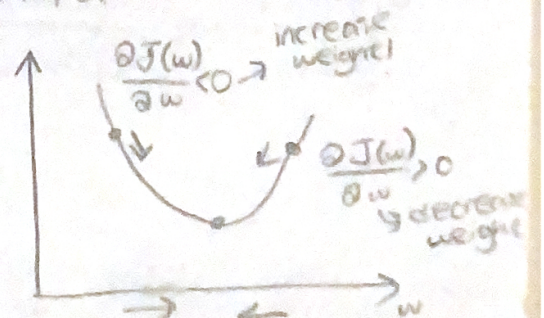
$$\frac{\partial \text{inp}_{o1}}{\partial w_5} = \text{out}_{n1} = \underline{\underline{0.593}}$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial \text{out}_{o1}} * \frac{\partial \text{out}_{o1}}{\partial \text{inp}_{o1}} * \frac{\partial \text{inp}_{o1}}{\partial w_5} = 0.741 * 0.18 * 0.593 = \underline{\underline{0.082}}$$

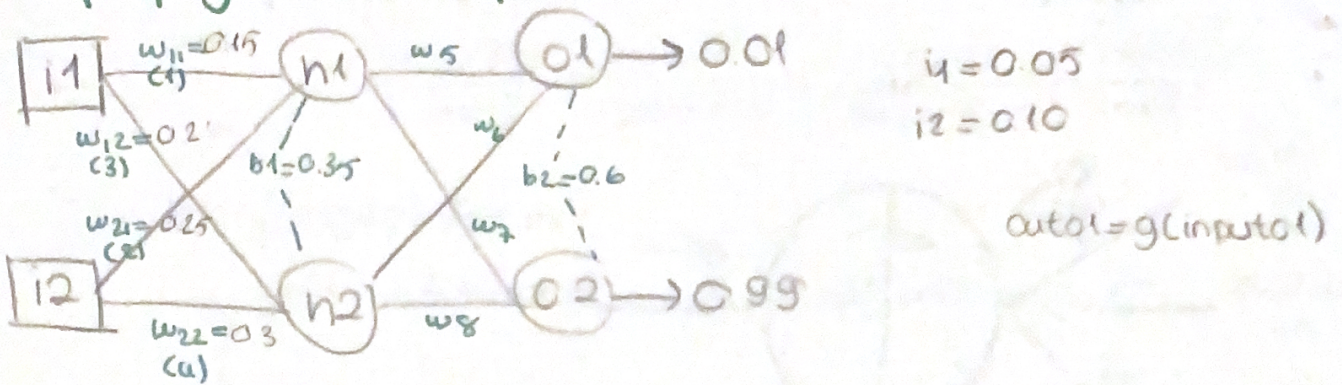
Delta Rule

$$\frac{\partial E_{\text{total}}}{\partial w_5} = - \underbrace{(\text{target}_{o1} - \text{out}_{o1})}_{\frac{\partial E_{\text{total}}}{\partial \text{out}_{o1}}} * \underbrace{\text{out}_{o1} (1 - \text{out}_{o1})}_{\frac{\partial \text{out}_{o1}}{\partial \text{inp}_{o1}}} * \underbrace{\text{out}_{n1}}_{\frac{\partial \text{inp}_{o1}}{\partial w_5}}$$

Then update $w_i = w - \alpha \frac{\partial J(w, b)}{\partial w}$
 $b_i = b - \alpha \frac{\partial J(w, b)}{\partial b}$



Backpropagation Example



Forward Pass

→ Input of n_1 : $w_{11} \cdot i_1 + w_{12} \cdot i_2 + b_1 = 0.15 \cdot 0.05 + 0.2 \cdot 0.1 + 0.35 = 0.3475$

← Output of n_1 : $\frac{1}{1 + e^{-0.3475}} = 0.593$

← Output of n_2 : 0.596

→ Input of o_1 : $w_6 \cdot 0.596 + w_5 \cdot 0.593 + 0.6 = 1.105$

← Output of $o_1 = \frac{1}{1 + e^{-1.105}} = \boxed{0.751}$

← Output of $o_2 = \boxed{0.772}$

Error Calculation

$$E_{total} = \sum \frac{1}{2} (\text{target} - \text{output})^2$$

$$= \frac{1}{2} \left((0.01 - 0.751)^2 + (0.1 - 0.772)^2 \right) = \underline{\underline{0.298}}$$

Total Error

Backward Pass

We want to know how much change in w_5 affects total error $\left(\frac{\partial E_{total}}{\partial w_5} \right)$

Chain Rule

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial input_{o1}} * \frac{\partial input_{o1}}{\partial w_5}$$

Autoencoders

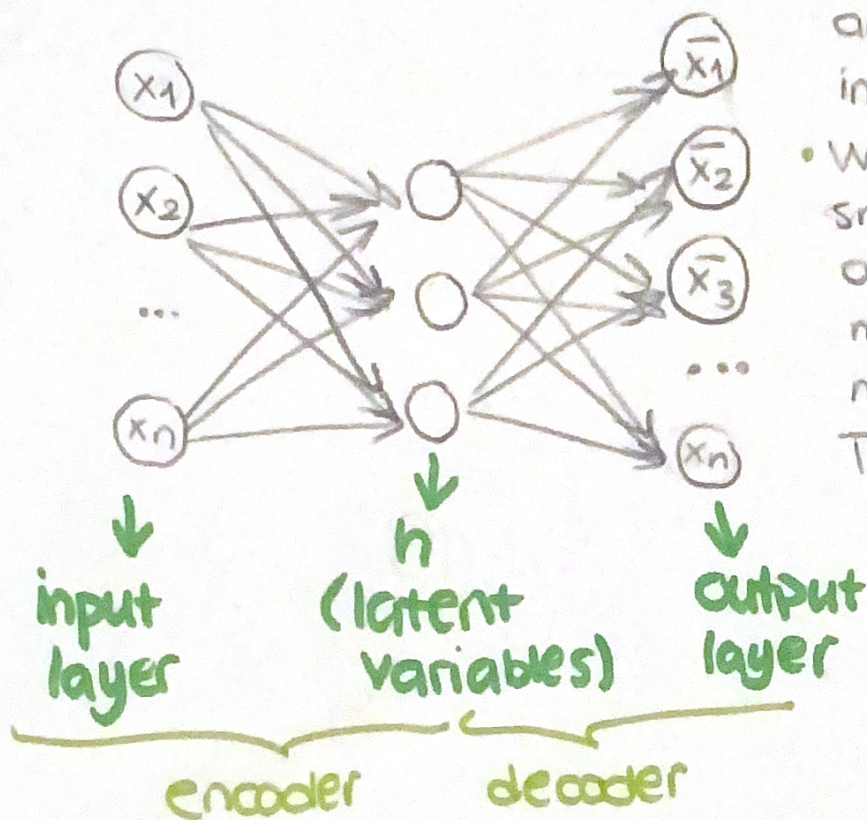
Idea: Autoencoders are trained to reconstruct input.

- When hidden layer is smaller than input & output layers, the model only learns the most salient features. This is called undercomplete autoencoder.

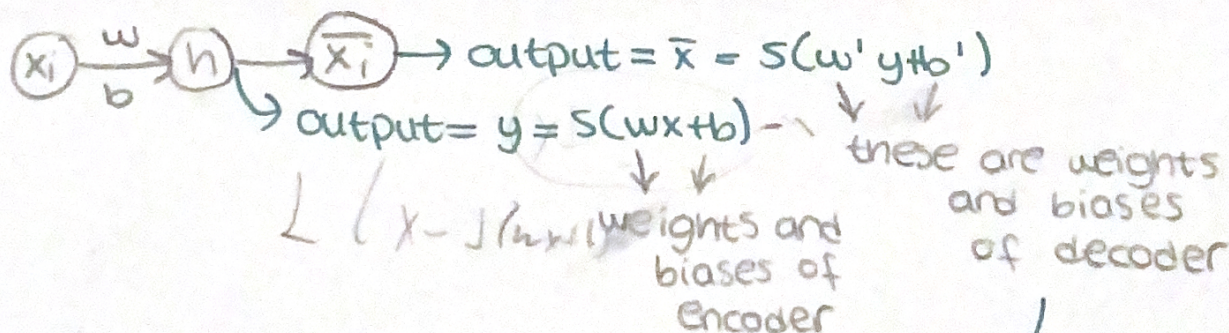
- Autoencoders learn unsupervised.

- $h_{w,b}(x) \approx x$

↓
identity function
(we're trying to learn this)

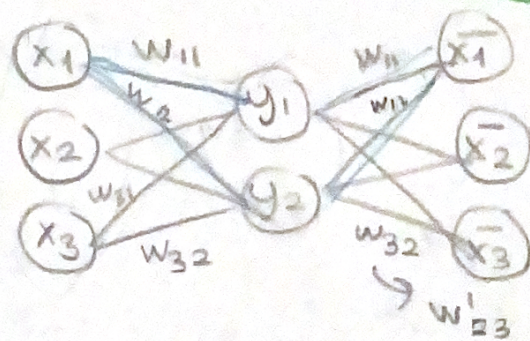


LOSS → $(\bar{x}_i - x_i)$



$W' = W^T$

one is transpose of another!



$$W = \begin{bmatrix} x_1 & x_2 & x_3 \\ w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$W' = W^T \rightarrow W' = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

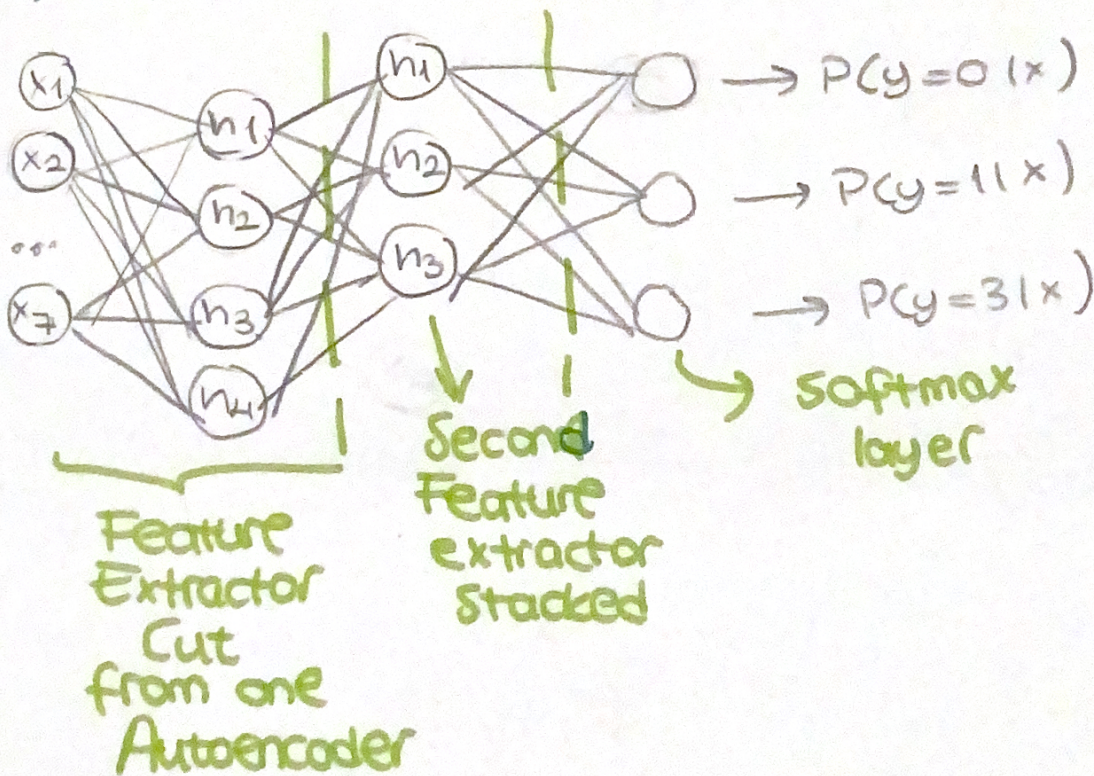
$L(x, \bar{x}) = \|x - \bar{x}\|^2$

Autoencoders (cont'd.)

Stacked Autoencoder

1. Train autoencoders to be feature extractor,
2. Add couple of more layers (optional)
3. Add a classifier layer on top \rightarrow Learning becomes supervised
4. Train with obtain specific data

\rightarrow This is used to avoid vanishing gradients.



Denoising Autoencoder

Changing the reconstruction error term.

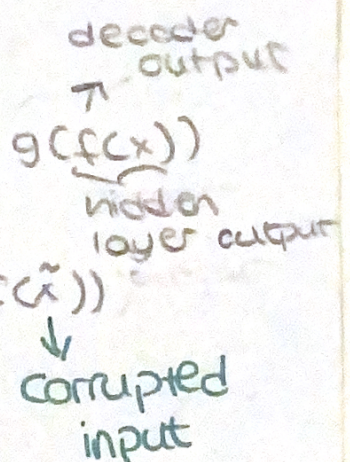
Traditional reconstruction loss $\rightarrow L(x, \hat{x})$

Denoiser autoencoder loss $\rightarrow L(x, g(f(\tilde{x})))$

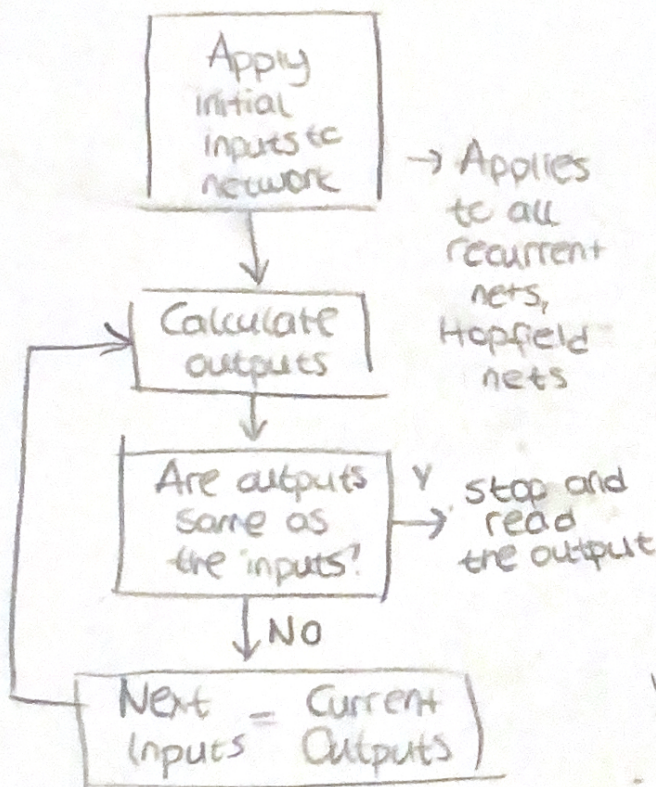
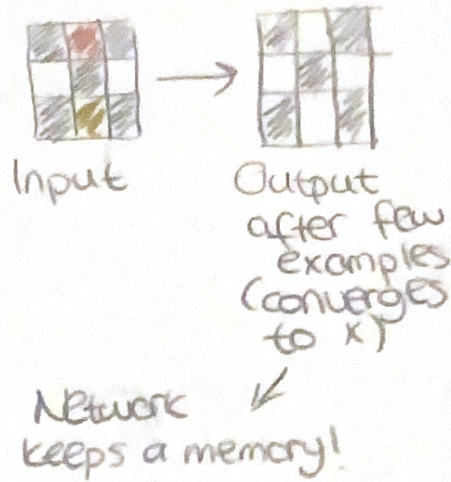
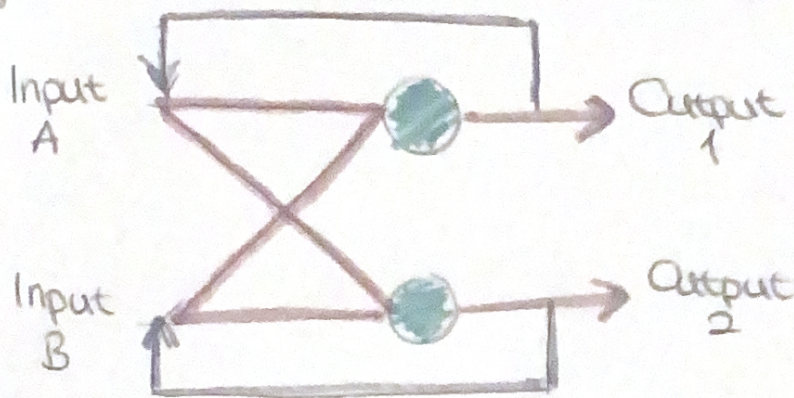
It receives a corrupted datapoint as input and tries to predict uncorrupted original data point.

$C(\tilde{x}|x) \rightarrow$ conditional distr over corrupted samples (\tilde{x})
 Preconstruct $(\tilde{x}|\tilde{x}) \rightarrow$ reconstruction distribution estimated from training pairs.

Preconst $(x|\tilde{x}) = \underbrace{p_{\text{decoder}}(x|n)}_{\text{gen}}$ \rightarrow output of encoder $f(\tilde{x})$



Recurrent Neural Network



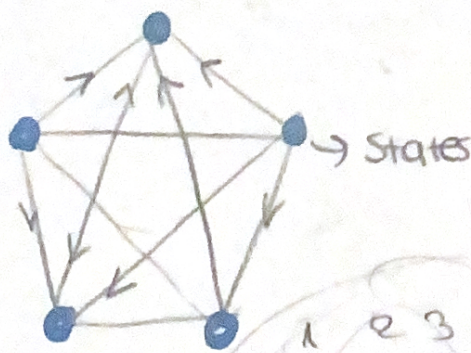
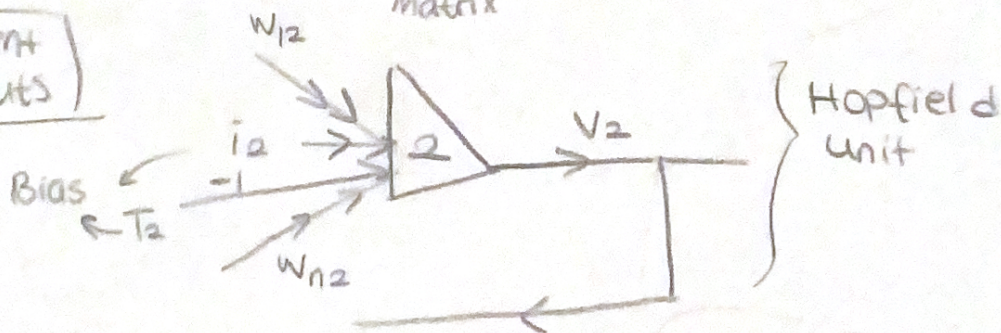
Hopfield Networks

- Weight matrix is symmetrical
 $W_{ij} = W_{ji}$
total input coming to neuron
- Output $v_i = -1$ if $net_i < 0$
 $v_i = 1$ if $net_i \geq 0$

Activation function used is the sign function.

$$net_i = w_i \cdot v + i_i - T_i$$

\downarrow weight matrix \downarrow output matrix \downarrow Bias inputs

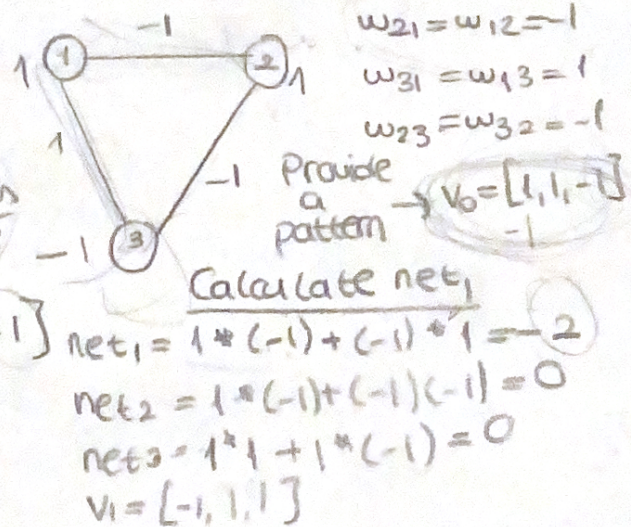


$$W_i = \begin{bmatrix} 1 & 2 & 3 \\ 0 & -1 & 1 \\ 2 & -1 & 0 & -1 \\ 3 & 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

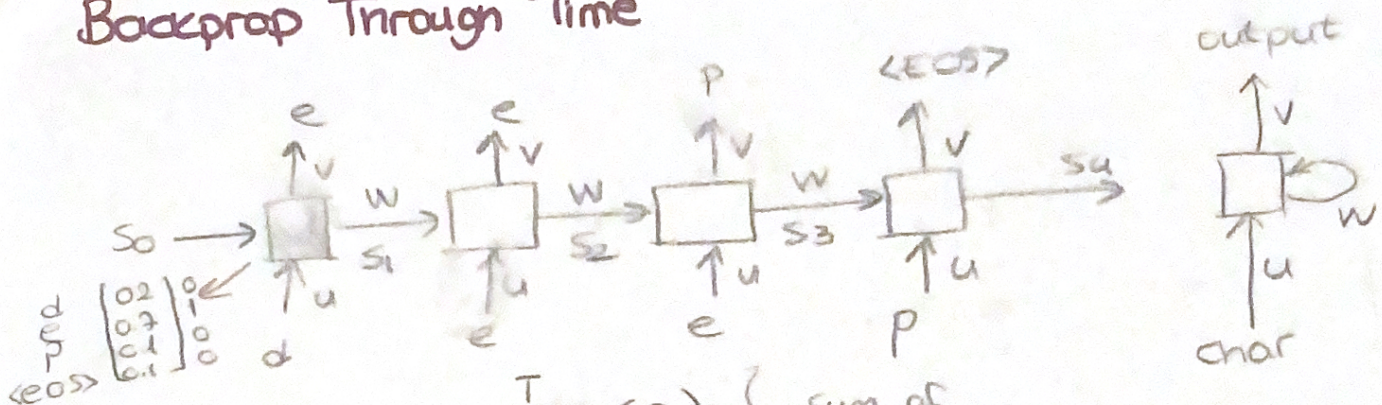
HL

$$W_{12} \cdot v_{12} + W_{13} \cdot v_{13}$$

ex 11



Backprop Through Time



Total Loss $\rightarrow \sum_{t=1}^T d_t(\theta)$ } Sum of loss over all steps

$d_t(\theta) = -\log(y_{tcc})$
 ↳ prob produced for true class at time step t (0.7 in above example, log will decrease if it gets closer to 1)

For backprop we compute gradients for w, u, v

For v : $\frac{\partial d(\theta)}{\partial v} = \frac{\partial d_1(\theta)}{\partial v} + \frac{\partial d_2(\theta)}{\partial v} + \dots$ } standard backprop

↳ this is not being affected by states

v is a matrix

s_4 depends on s_3 & w
 s_3 depends on s_2 & w
 ...

For w : $\frac{\partial d(\theta)}{\partial w}$

↳ this is calculated with regular BP

$\frac{\partial d_4(\theta)}{\partial w} = \frac{\partial d_4(\theta)}{\partial s_4} \cdot \frac{\partial s_4}{\partial w}$

we can't calculate this
 $s_4 = G(w s_3 + b)$

Total derivative not a part.

Explicit $\rightarrow \frac{\partial^* s_4}{\partial w}$ } treats all inputs as constant

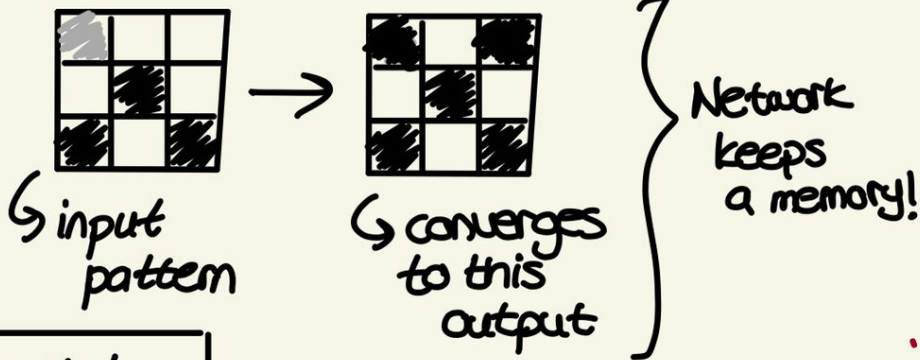
Implicit \rightarrow Sum all indirect paths from s_4 to w .

$$\frac{\partial s_4}{\partial w} = \frac{\partial^* s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \cdot \frac{\partial s_3}{\partial w} = \frac{\partial^* s_4}{\partial w} + \frac{\partial s_4}{\partial s_3} \left[\underbrace{\frac{\partial^* s_3}{\partial w}}_{\text{exp}} + \frac{\partial s_3}{\partial s_2} \frac{\partial s_2}{\partial w} \right]$$

$$\dots$$

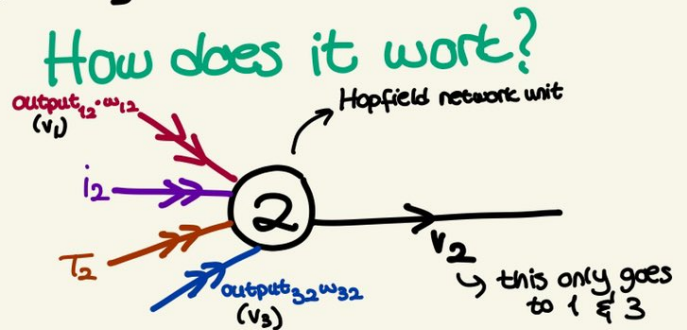
$$\frac{\partial^2 s_2}{\partial s_1} \cdot \frac{\partial s_1}{\partial w}$$

Hopfield Networks

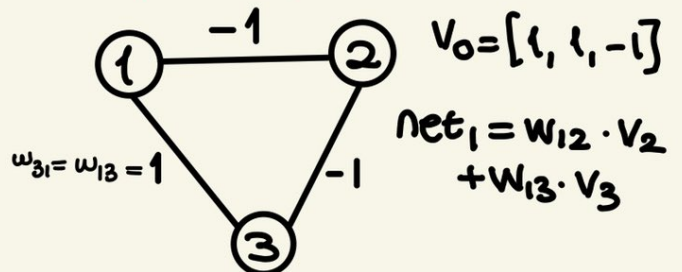


```

    graph TD
      A[Apply initial inputs to network (y)] --> B[Calculate output (y-hat)]
      B --> C{Is output same as the input?}
      C -- Yes --> D[Stop]
      C -- No --> E[Next Input ← Current Output]
      E --> B
  
```



Example Hopfield Network



Calculate net_i:

$$net_1 = 1 \cdot (-1) + (-1) \cdot 1 = -2 \rightarrow -1$$

v₁ w₁₂ v₃ w₁₃

$$net_2 = 1 \cdot (-1) + (-1) \cdot (-1) = 0 \rightarrow 1$$

$$net_3 = 1 \cdot 1 + 1 \cdot (-1) = 0 \rightarrow 1$$

$$V_1 = [-1, 1, 1]$$

Properties:

- Weight matrix is symmetrical. ($w_{ij} = w_{ji}$)
- Activation function is sign function.

Output $v_i = -1$ if $net_i < 0$
 $v_i = 1$ if $net_i \geq 0$
 ↳ total input coming to neuron

Conv Nets

→ Used for problems that have grid-like topology.

* 1D → Time series

* 2D → Image recognition

Convolution

$$S(t) = \underbrace{(x * w)}_{\int x(a)w(t-a)da} (t) \quad \left. \begin{array}{l} \text{input} \\ \text{kernel} \end{array} \right\} \text{feature map}$$

$$2D \rightarrow S(i,j) = (K * I)(i,j) = \sum_m \sum_n (i-m, j-n) K(m,n)$$

↓ Kernel ↓ Input

Why?

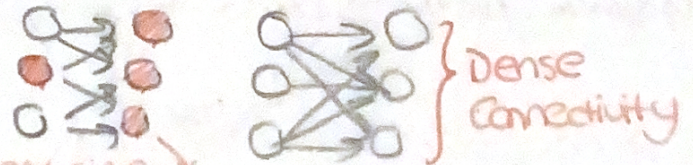
→ Sparse Interactions: Features can be detected with less params

→ Parameter Sharing

→ Equivariant Representations

→ We can work with inputs of variant size

receptive field



• Sparsity: we store less info thanks to kernels.

$\left. \begin{array}{l} m \text{ input} \\ n \text{ output} \end{array} \right\} \begin{array}{l} m \times n \text{ params} \\ O(m \times n) \end{array}$

$k \rightarrow$ number of connections each input can have
 $O(k \times n) \rightarrow$ sparsely connected approach

• Parameter Sharing: Using same parameter for more than one function in a model.

ex11 Dense Nets → Each element of weight matrix is multiplied for once with input

• CNNs → Has tied weights, value of the weight applied to one input is tied to another value elsewhere. Each element in kernel is used at every position of the input. (Except for boundary pixels)

Parameter Sharing



CNN



DNN

→ Connections that use a particular parameter.
(Reduces storage requirements!)

Equivariance

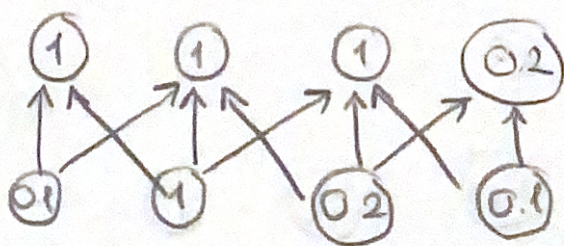
An equivariant function means if the input changes, output changes in the same way.

$$f \xleftrightarrow{\text{equivariant}} g \quad \text{if } f(g(x)) = g(f(x))$$

- For convolution → If we move object in the input, it's representation will move the same amount in the input. Thus we can create 2D map of where certain features appear.

Pooling

- Pooling replaces outputs of the net at a certain location with a summary statistic of the nearby outputs.
- Max pooling reports the maximum outputs within a rectangular nbhd. (Other pooling func: L2 norm, average, weighted average)
- The representation becomes invariant to small translations. (Useful when some feature is present than exactly where it is.)
e.g. we don't need to know exact location of eyes, we need to know that a face has eyes on right & left.
- Pooling is essential to handle varying input sizes. (Varying size of an offset between pooling regions such that classification layer always receives the same number of summary stats regardless of input size.)



Max Pooling
(Invariance)

