

Neural Networks

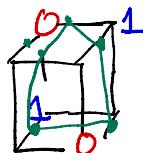
- Classification & regression
- uses ideas from perceptrons, logistic regression, ensemble learning ...

Perceptrons: Can't perform XOR

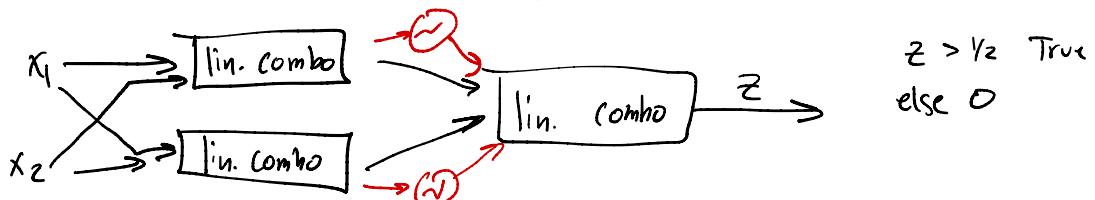
		x_1	x_2
		0	1
x_1	0	1	0
x_2	1	0	1

Can't separate these 4 points in 2D w/ a perceptron
 ↳ perceptron research halted

If you add a quadratic feature $x_1 x_2$, then XOR is linearly separable in 3D

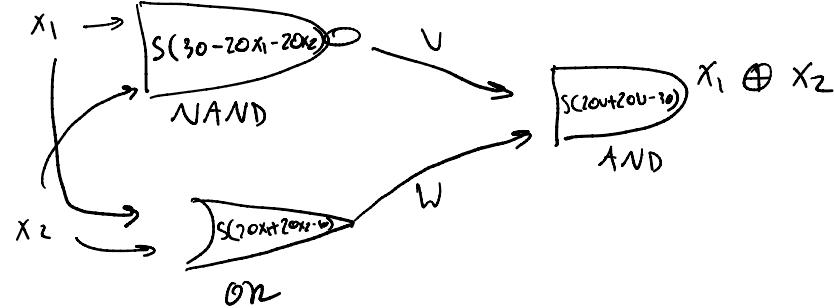


Idea: Use linear classifiers whose outputs are inputs to other lin. classifiers



Still can't do $XOR \rightarrow \text{lin. combo}$ of lin. combos is still linear

Solution: add a Nonlinearity



Neural Network w/ 1 hidden layer

Typical architecture:

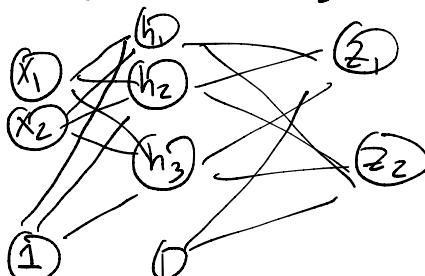
- Input layers: $x_1 \dots x_d$, $x_{d+1} = 1$
- hidden units: $h_1 \dots h_m$, $h_{m+1} = 1$
- output layers: $z_1 \dots z_k$

We add weights between each layer:

- between input & h : $m \times (d+1)$ matrix V , V_i^T is row i
- between h & output: $k \times (m+1)$ matrix W

recall $S(v) = \frac{1}{1+e^{-v}}$ activation function

$$S_1(\vec{v}) = [S(v_1) \dots S(v_n) \ 1]^T$$



$$\vec{h} = S_1(V \vec{x}) \rightarrow h_i = S(v_i^T \vec{x})$$

$$\vec{z} = S(W \vec{h}) = S(W S_1(V \vec{x}))$$

- usually, NNs have more than one output

advantage of NNs: When training multiple classifiers through a single neural network, sometimes these classifiers take advantage of useful hidden units that support other classifiers.

Training Neural Networks

- SGD or BGD
- loss function $L(\vec{z}, \vec{y})$ most common: Squared error
- cost function $J(h) = \frac{1}{n} \sum L(h(\vec{x}_i), \vec{y}_i)$
- goal: find W & V that minimize J

Usually, cost fn. isn't convex at all \rightarrow many local minima

Issue: NNs have Symmetry - one hidden unit is indistinguishable from another
↳ gradient descent can't break the symmetry
↳ might be in situation in which all weights of h_1 & h_2 are same
Solution: random initialization of weights

Batch gradient descent:

$W \leftarrow$ random vector

repeat

$$W \leftarrow W - \frac{\epsilon \nabla J(W)}{\epsilon \nabla L(z_i, h_i, V)}$$

make sure random initialization isn't too big since $s'(\cdot) \rightarrow 0$ if $w_i \gg 0$

Computing Gradients for Arithmetic Expressions

goal: compute $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial a} \\ \frac{\partial f}{\partial b} \\ \frac{\partial f}{\partial c} \\ \frac{\partial f}{\partial d} \end{bmatrix}$

The diagram shows an arithmetic expression tree for $f = (a+b)(c+d)$. The root node is f , which has two children: c and d . Node c has two children: a and b . Node d has two children: e and f . Node e has two children: x and y . Red annotations show the computation of gradients for each node:

- For node a : $\frac{\partial f}{\partial a} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial a} = 1$
- For node b : $\frac{\partial f}{\partial b} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial b} = 1$
- For node c : $\frac{\partial f}{\partial c} = \frac{\partial f}{\partial e} \cdot \frac{\partial e}{\partial c} = 1$
- For node d : $\frac{\partial f}{\partial d} = \frac{\partial f}{\partial e} \cdot \frac{\partial e}{\partial d} + \frac{\partial f}{\partial f} \cdot \frac{\partial f}{\partial d} = 2e$
- For node e : $\frac{\partial f}{\partial e} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial e} + \frac{\partial f}{\partial f} \cdot \frac{\partial f}{\partial e} = 1$
- For node x : $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial e} \cdot \frac{\partial e}{\partial x} = 1$
- For node y : $\frac{\partial f}{\partial y} = \frac{\partial f}{\partial e} \cdot \frac{\partial e}{\partial y} = 1$

Below the tree, it is shown that $d = a+b \Rightarrow \frac{\partial d}{\partial a} = 1, \frac{\partial d}{\partial b} = 1$.

What if a unit's output goes to more than 1 unit?

Consider:

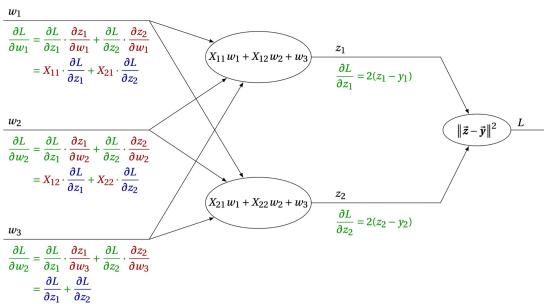
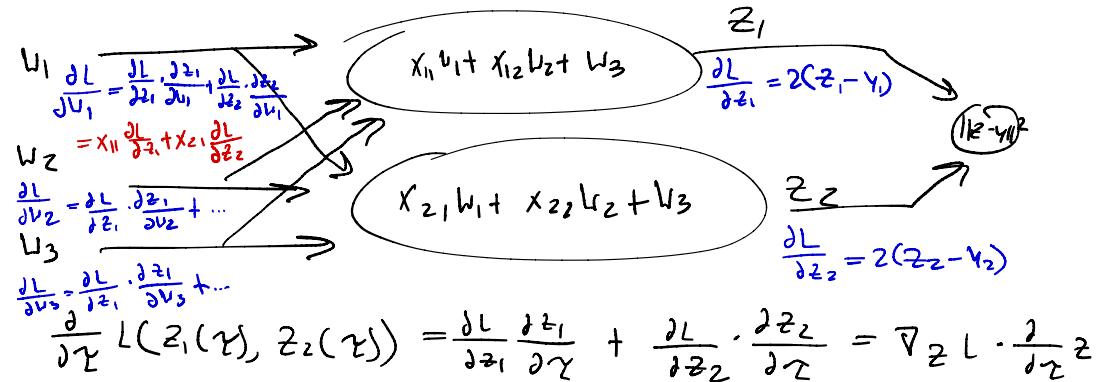


Figure 17.6: Gradient of an expression from least squares linear regression. Backpropagation arrows are omitted here, as they'd make the diagram quite messy.

The backprop Algo.

V_i^T is row i of weight matrix V , same for W

$$\text{recall } S'(r) = S(r)(1-S(r))$$

$$h_i = S(V_i^T \vec{x}) \rightarrow \nabla_{V_i} h_i = S'(V_i^T \vec{x}) \vec{x} = h_i(1-h_i) \vec{x}$$

$$= z_j(1-z_j) \vec{h}$$

