

Lecture 1 ~ some comments

Mathematically, a DNN is a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$
 $x \in \mathbb{R}^n \quad y \in \mathbb{R}^m$

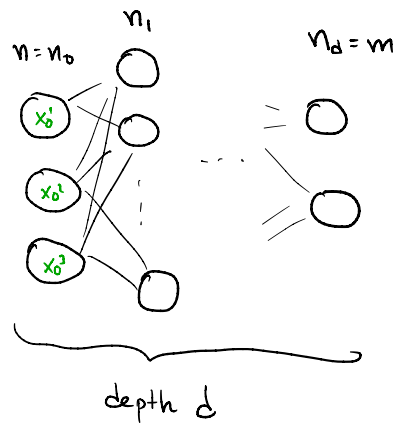
defined recursively $x_0 = x$

$$x_1 = A(x_0 \cdot w_1 + b_1) \in \mathbb{R}^{n_1}$$

$$x_2 = A(x_1 \cdot w_2 + b_2) \in \mathbb{R}^{n_2}$$

\vdots

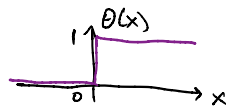
$$f(x) = x_d = A(x_{d-1} \cdot w_d + b_d) \in \mathbb{R}^{n_d} = \mathbb{R}^m$$



Observe: if A is trivial ($A(x) = x$), then f collapses to a single linear function
 $f(x) = x \cdot \tilde{w} + \tilde{b}$ (for suitable \tilde{w}, \tilde{b}).

Boring! Linear fits.

A more interesting limit: $A(x) = \Theta(x)$ step function



\leadsto "perceptron" networks from 50's & 60's

model for (iterative, inter-dependent) decision making. (See supplementary reading.)

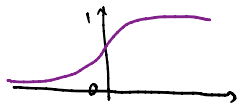
Problem w/ $A(x) = \Theta(x)$: not differentiable (not even cts) \rightarrow difficult to learn

12

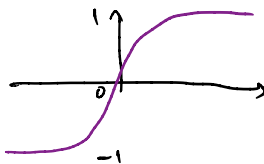
Learning: how to alter parameters w, b to improve the outcome of the model.

Rob: in much of modern ML, $A(x)$ is taken to be a differentiable function, w/ a nice derivative

e.g. $A(x) = \sigma(x) = \frac{1}{1+e^{-x}}$ or $A(x) = \tanh(x)$



$$A'(x) = A(x)(1-A(x))$$



$$A'(x) = 1 - A(x)^2$$

and "outcome" is measured by a loss function $\ell: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$

$$\ell(x, y) = L(f(x), y) \stackrel{\text{e.g.}}{=} \frac{|f(x) - y|^2}{2}$$

Aside: given a metric space W and a function $l: W \rightarrow \mathbb{R}$
 one can define gradient flow on W .

In local coords w^i , metric $g_{ij} dw^i dw^j$ (or $g_{ij} = g(\frac{\partial}{\partial w^i}, \frac{\partial}{\partial w^j})$)
 the gradient flow ODE is $\frac{\partial w^i}{\partial t} = - g^{ij} \frac{\partial l}{\partial w^j}$.

With discrete steps, this becomes $\Delta w^i = - c \underset{\substack{\uparrow \\ \text{learning rate}}}{g^{ij}} \frac{\partial l}{\partial w^j}$.

For a DNN, really $l: \mathbb{R}^n \times \underbrace{\mathbb{R}^{n \times n_1 + n_1 \times n_2 + n_2 \times n_3 + \dots + n_{d-1} \times n_d + n_d}}_{\text{space of } w, b \text{ params, assumed vec space w/ Euclidean metric}} \times \mathbb{R}^m \rightarrow \mathbb{R}$ $l(x; w, b; y) = \frac{|f_{w,b}(x) - y|^2}{2}$

$$l(x, y) = L(f(x), y)$$

x, y are fixed, only care about w, b dependence

$$\boxed{g^{ij} = \delta^{ij}}$$

with appropriate index contractions, transposes, etc.

$$\begin{aligned} dl &= \frac{\partial L}{\partial f(x)} df(x) = \frac{\partial L}{\partial f(x)} A'(x_{d-1}, w_d + b_d) (x_{d-1} \cdot dw_d + db_d + \underbrace{dx_{d-1} \cdot w_d}_{\substack{\frac{\partial L}{\partial w_{d-1}} \quad \frac{\partial L}{\partial b_{d-1}} \\ \text{keep going}}}) \\ &= \frac{\partial L}{\partial f} A'(x_{d-1}, w_d + b_d) (x_{d-1} dw_d + db_d + A'(x_{d-2}, w_{d-1} + b_{d-1}) (x_{d-2} dw_{d-1} + db_{d-1} + \underbrace{dx_{d-2} \cdot w_{d-1}}_{\text{keep going}}) w_d) \end{aligned}$$

\leadsto "back-propagation" to compute components of the gradient.