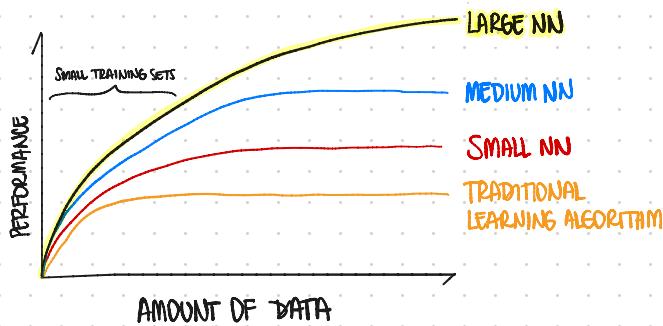


NEURAL NETWORKS AND DEEP LEARNING

COURSE I

WHAT IS A NEURAL NETWORK?

- * Machine learning
- * A computer learns to perform a task by analyzing training samples.
- * Consists of thousands or even millions of simple processing nodes that are interconnected.



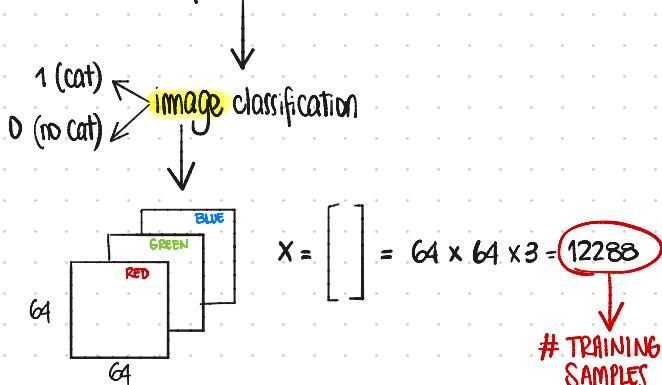
why is deep learning taking off? →

SCALE DRIVES DEEP LEARNING PROGRESS

INPUT (X)	OUTPUT (Y)	APPLICATION
home features	price	real estate
ad, user info	click on ad? (0/1)	online advertising
image	object (1,...,1000)	photo tagging → CNN
audio	text transcript	speech recognition → RNN
image, radar info	position of other cars	autonomous driving → CUSTOM / HYBRID

LOGISTIC REGRESSION AS A NEURAL NETWORK

* logistic regression is an algorithm for binary classification,



#1. given X , we want to find \hat{Y} → PREDICTION (OUTPUT)

$$\hat{Y} = P(Y=1 | X)$$

#2. parameters $w \in \mathbb{R}^{n_x}$
 $b \in \mathbb{R}^{n_x}$ (inter spectrum)

#3. output → result of a Sigmoid function

$$\hat{Y} = \sigma(w^T x + b)$$

NOTATION!

#1 $(x, y), x \in \mathbb{R}^{n_x} \vee y \in \mathbb{R}$

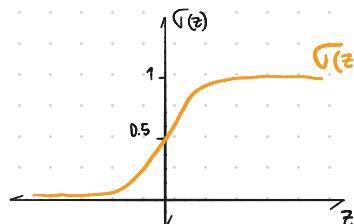
#2 $M \rightarrow$ training samples

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

#3 $X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$

$\uparrow n_x$
 $\leftarrow m \rightarrow$

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



* z LARGE:
 $\sigma(z) \approx \frac{1}{1+0} = 1$

* z LARGE NEGATIVE:
 $\sigma(z) \approx \frac{1}{1+\infty} = 0$

* to train the parameters w and b we need to define a cost function:

 How good is our output \hat{y} ?

LOSS FUNCTION
(ONE DATA SAMPLE)



(APPLIED TO ALL PARAMETERS)
COST FUNCTION



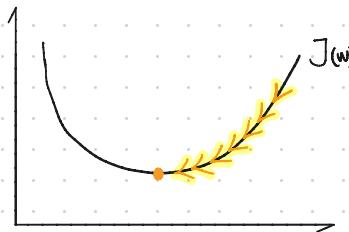
$$Y=1: L(\hat{y}, y) = -\log \hat{y} \Rightarrow \log \hat{y} \text{ small, } \hat{y} \text{ large}$$

$$Y=0: L(\hat{y}, y) = -\log(1-\hat{y}) \Rightarrow \log(1-\hat{y}) \text{ large, } \hat{y} \text{ small}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-\hat{y}^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

how can we find w and b that minimizes $J(w, b)$?

GRADIENT DESCENT! (DIRECTION OF STEEPEST DESCENT)



* you initialize w and b at some random point, and takes steps towards the deepest descent until you get close to the minimum.

REPEAT {

$$\left. \begin{array}{l} w := w - \alpha \frac{dJ(w)}{dw} \end{array} \right\}$$

$$J(w, b) \left\{ \begin{array}{l} w = w - \alpha \frac{dJ(w, b)}{dw} \\ b = b - \alpha \frac{dJ(w, b)}{db} \end{array} \right.$$

this mostly works with logistic regression only

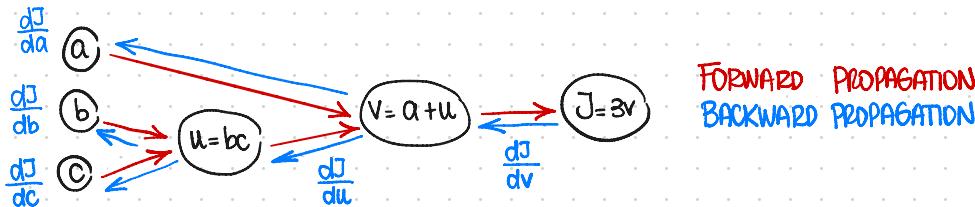
* Computation of neural networks:

- #1 forward propagation → compute output of the neural network
- #2 backward propagation → compute gradients

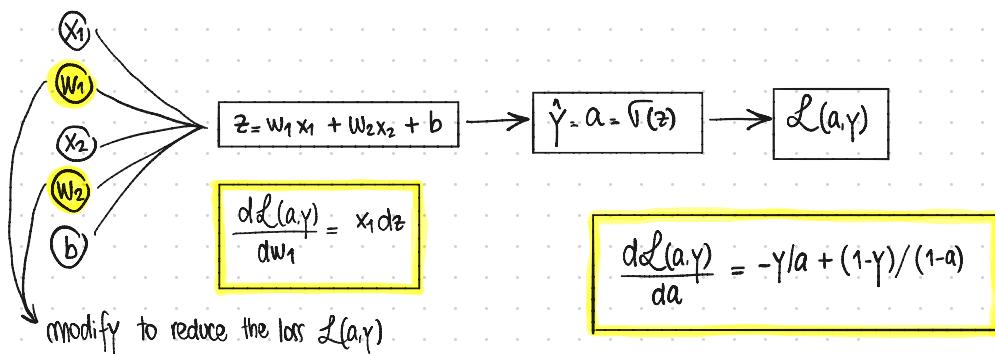
* Example:

$$J(a, b, c) = 3(a + bc)$$

$$u = bc, v = a + u, J = 3v$$



* in logistic regression...



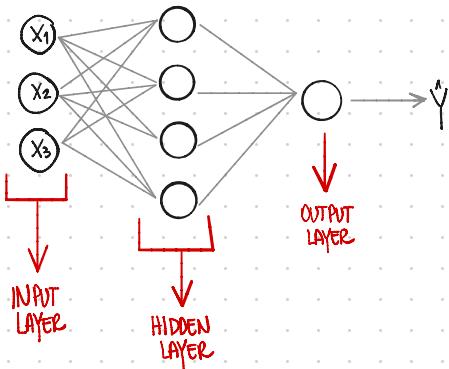
* gradient descent on m examples:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)}, y^{(i)})$$

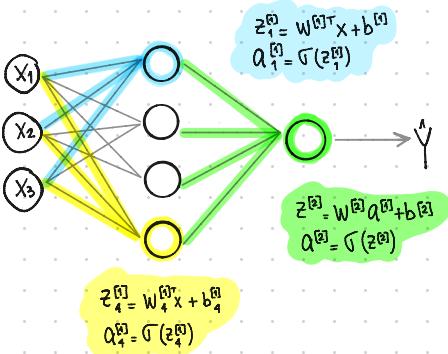
\downarrow

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^\top x^{(i)} + b)$$

SHALLOW NEURAL NETWORKS



* how to compute a NN's output?



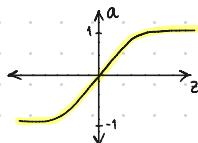
* vectorizing across multiple examples:

$$\begin{aligned} Z^{[1]} &= W^{[1]T} X + b^{[1]} \\ A^{[1]} &= \sigma(Z^{[1]}) \\ Z^{[2]} &= W^{[2]T} A^{[1]} + b^{[2]} \\ A^{[2]} &= \sigma(Z^{[2]}) \end{aligned}$$

$$A^{[1]} = \left[\begin{array}{c} \vdots \\ \text{HIDDEN UNITS} \\ \vdots \end{array} \right] \quad \begin{matrix} \nearrow \\ \text{TRAIN SAMPLES} \end{matrix}$$

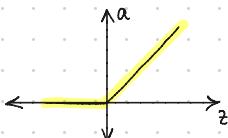
* activation functions:

$$\#1. \quad a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



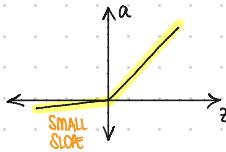
* always works better than $\sigma(z)$ for hidden units because the mean of the activation functions that come out from the hidden layers are close to 0.

$$\#2. \quad \text{ReLU} \quad a = \max(0, z)$$



* the derivative is 1 so long as z is positive, and it's 0 when z is negative.

#3. leaky ReLU $a = \max(0.001z, z)$



#4. linear activation function $g(z) = z$

* it might be useful for the output layer of a linear regression problem.

* derivatives of activation functions:

#1. Sigmoid activation function:

$$\frac{dg(z)}{dz} = \frac{1}{1+e^z} \left(1 - \frac{1}{1+e^z}\right) = g(z)(1-g(z))$$

$$a(1-a) = g'(z)$$

#2. tanh activation function:

$$\frac{dg(z)}{dz} = 1 - [\tanh(z)]^2$$

$$1-a^2 = g'(z)$$

#3. ReLU and leaky ReLU:

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

ReLU

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

leaky ReLU

DOWNSIDE OF $\text{G}(z)$ AND $\tanh(z)$:

* if z is either very large or very small, the gradient (derivative) becomes very small. It may end up being close to 0

↳ this can slow down gradient descent.

GRADIENT DESCENT FOR NEURAL NETWORKS

* parameters \rightarrow

- $w^{[1]}, \dim n^{[1]}, n^{[1]}$
- $b^{[1]}, \dim n^{[1]}, 1$
- $w^{[2]}, \dim n^{[2]}, n^{[2]}$
- $b^{[2]}, \dim n^{[2]}, 1$

* cost function $\rightarrow \mathcal{J}(w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y_i, \hat{y}_i)$

* forward propagation:

$$z^{[1]} = w^{[1]} X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

* back propagation:

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True})$$

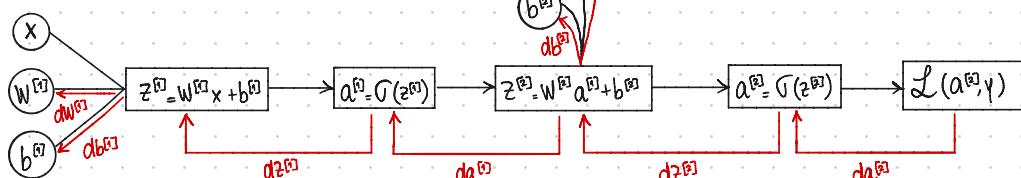
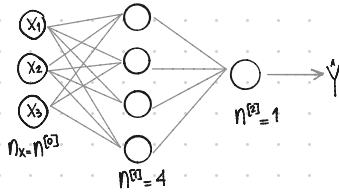
$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]}'(z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$$

BACKPROPAGATION INTUITION

* for a given NN of form:



$$dZ^{[1]} = W^{[2]T} dZ^{[2]} g^{[1]}'(z^{[1]})$$

$$dW^{[2]} = dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = dZ^{[2]}$$

dimensions:

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]}'(z^{[1]})$$

$$(n_h[1], 1) = (n_h[2], n_h[1]) * (n_h[1], 1) * (n_h[1], 1)$$

$$dZ^{[2]} = A^{[2]} - Y$$

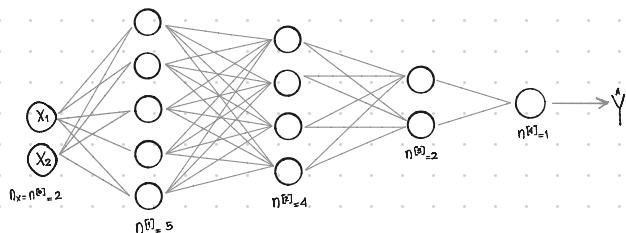
$$dW^{[2]} = dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = dZ^{[2]}$$

RANDOM INITIALIZATION

- * weights initialized with zeros → all hidden units compute the same operation or the same function.
- * initialize weights randomly

DEEP NEURAL NETWORKS



$$\begin{aligned} z^{[1]} &= W^{[1]} \cdot x + b^{[1]} \\ (n^{[1]}, 1) &= (n^{[1]}, n^{[1]}) (n^{[1]}, 1) + (n^{[1]}, 1) \\ (5, 1) &= (5, 2) (2, 1) + (3, 1) \end{aligned}$$

$$\begin{aligned} W^{[1]} &= (n^{[1]}, n^{[2]}) \\ b^{[1]} &= (n^{[1]}, 1) \end{aligned}$$

* vectorization:

$$\begin{aligned} z^{[1]} &= W^{[1]} X + b^{[1]} \\ (n^{[1]}, m) &= (n^{[1]}, n^{[1]}) (n^{[1]}, m) + (n^{[1]}, 1) \end{aligned}$$

- * building blocks of deep neural networks:

