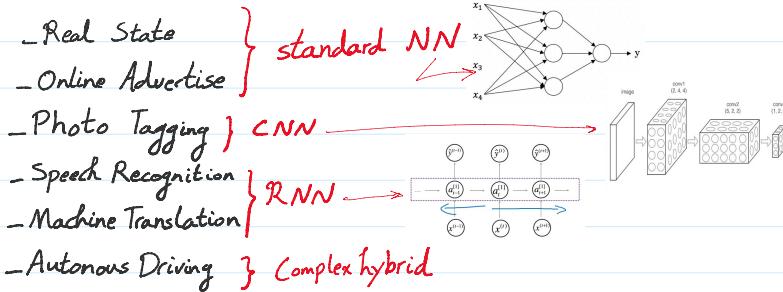


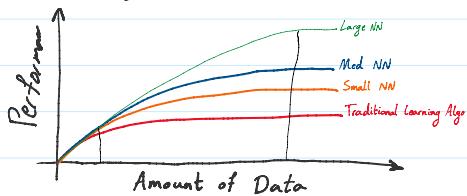
Application:



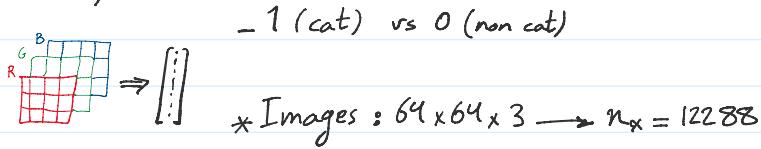
Supervised Learning:

- Structured Data →
- Unstructured Data → * Audio
* Image
* Text

Deep Learning Take-off:



Binary Classification:



$$(x, y) : x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

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

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}_{n_x \times m} \rightarrow X \in \mathbb{R}^{n_x \times m} \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \rightarrow Y \in \mathbb{R}^{1 \times m}$$

Logistic Regression:

$$\text{given } x, \hat{y} = P(y=1 | x)$$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

$$\text{Output: } \hat{y} = \sigma(w^T x + b)$$

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



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



Cost Function:

Loss (error) function: $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

$$\rightarrow L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

$$\begin{aligned} \text{Cost } J(w, b) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})) \end{aligned}$$

Gradient Descent:

- find w, b that minimize $J(w, b)$

$$\begin{aligned} \rightarrow w &:= w - \alpha \frac{\partial J(w, b)}{\partial w} \quad dw \\ \rightarrow b &:= b - \alpha \frac{\partial J(w, b)}{\partial b} \quad db \\ \text{learning rate } \alpha &\quad \swarrow \quad \searrow \end{aligned}$$

LR GD:

$$\begin{cases} z = w^T x + b \\ \hat{y} = a = \sigma(z) \\ L(a, y) = -(y \log a + (1-y) \log (1-a)) \end{cases}$$

$$\Rightarrow \begin{cases} dz = \frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} = a - y \\ da = \frac{\partial L}{\partial a} = \frac{-y}{a} + \frac{1-y}{1-a} \\ dw = \frac{\partial L}{\partial w} = \alpha dz \\ db = dz \end{cases}$$

On m examples $\rightarrow J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) \quad * a^{(i)} = \sigma^{(i)}(w^T x^{(i)} + b)$

$$\left. \begin{aligned} dZ &= [dz^{(1)} \dots dz^{(m)}] \\ A &= [a^{(1)} \dots a^{(m)}] \quad Y = [y^{(1)} \dots y^{(m)}] \end{aligned} \right\} dZ = A - Y$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$dw = \frac{1}{m} X dZ^T = \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} = \frac{1}{m} [x^{(1)} dz^{(1)} \dots x^{(m)} dz^{(m)}]$$

$$w_- = \alpha dw \quad b_- = \alpha db$$

Python Broadcasting :

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 10 & 10 & 10 \\ 10 & 10 & 10 \end{bmatrix} = \begin{bmatrix} 11 & 12 & 13 \\ 14 & 15 & 16 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 10 \\ 20 \\ 20 \end{bmatrix} = \begin{bmatrix} 11 & 12 & 13 \\ 24 & 25 & 26 \end{bmatrix}$$

Notations:

$(i) \rightarrow i^{\text{th}}$ example
 $(l) \rightarrow l^{\text{th}}$ layer

m : Num of examples

$n_h^{[l]}$: Num of hidden layers of l^{th}

n_x : Input size

n_y : Num of classes

$n_x = n_h^{[0]}$

$n_y = n_h^{[L+1]}$

$$X \in \mathbb{R}^{n_x \times m}$$

$x^{(i)} \in \mathbb{R}^{n_x}$ column

$$Y \in \mathbb{R}^{n_y \times m}$$

label matrix
 $y^{(i)} \in \mathbb{R}^{n_y}$ output label of i^{th}

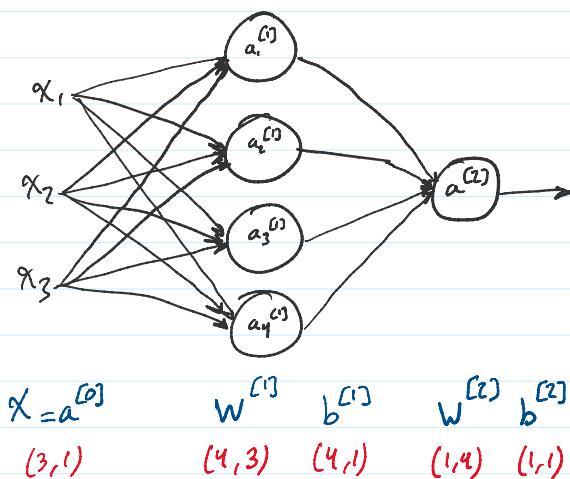
$$W^{[l]} \in \mathbb{R}^{n_{\text{prev}} \times n_{\text{next}}}$$

$$b^{[l]} \in \mathbb{R}^{n_{\text{next}}}$$

weights
bias

$$a = g^{[l]}(W_x x^{(i)} + b_l) = g^{[l]}(z_i) \quad \hat{y}^{(i)} = \text{softmax}(W_h h + b_2)$$

Neural Network:



$$z^{[1]} = W^{[1]} X + b^{[1]} \quad (4, 1)$$

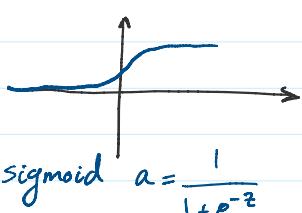
$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \quad (1, 1)$$

$$a^{[2]} = \sigma(z^{[2]})$$

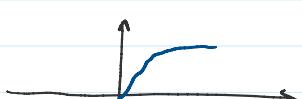
$$\left. \begin{array}{l} Z^{[1]} = W^{[1]} X + b^{[1]} \\ A^{[1]} = \sigma(Z^{[1]}) \\ Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} = \sigma(Z^{[2]}) \end{array} \right\} m \text{ Examples}$$

Activation Functions:



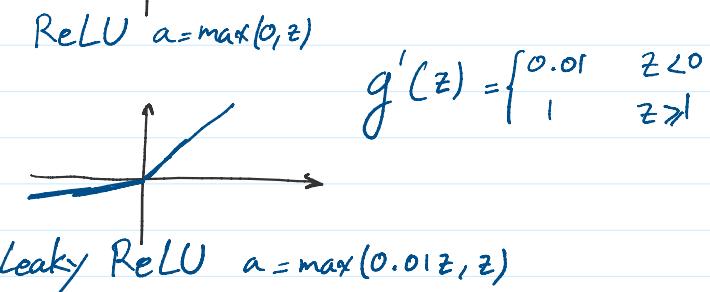
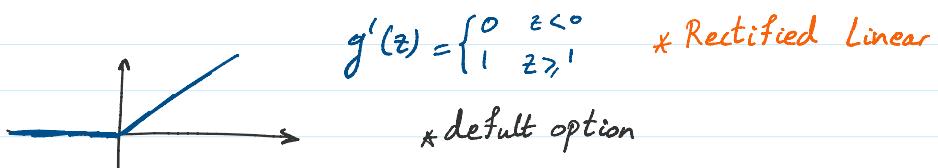
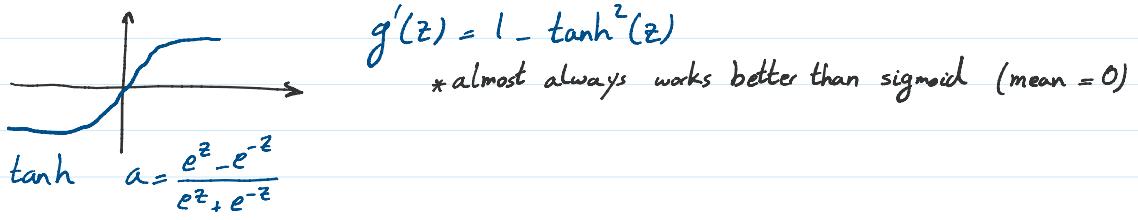
$$g'(z) = a(1-a)$$

* output layer when binary classification.



$$g'(z) = 1 - \tanh^2(z)$$

* almost always works better than sigmoid (mean = 0)



Gradient Descent for NN:

forward:

$$\begin{matrix} Z^{[1]} \\ Z^{[2]} \end{matrix} \rightarrow \begin{matrix} A^{[1]} \\ A^{[2]} \end{matrix}$$

backpropagation:

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

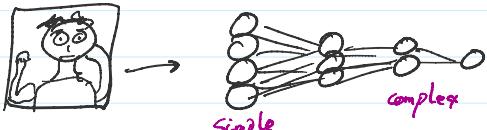
$$\begin{aligned} dW^{[2]} &= \frac{1}{m} dZ^{[2]} A^{[1]T} \\ db^{[2]} &= \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims}=T) \end{aligned}$$

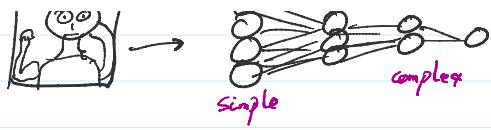
$$dZ^{[1]} = W^{[2]T} dZ^{[2]} \times \underbrace{g'(Z^{[1]})}_{\text{element wise}} \quad (n^{[1]}, m) \quad (n^{[2]}, m)$$

$$\begin{aligned} dW^{[1]} &= \frac{1}{m} dZ^{[1]} X^T \\ db^{[1]} &= \frac{1}{m} \text{np.sum}(\dots) \end{aligned}$$

1 layer (Shallow) \longrightarrow 5 layers (Deep)

$$\begin{aligned} W^{[L]} & (n^{[L]}, n^{[L-1]}) \\ b^{[L]} & (n^{[L]}, 1) \\ dW^{[L]} & (n^{[L]}, n^{[L-1]}) \\ db^{[L]} & (n^{[L]}, 1) \end{aligned}$$





* shallower networks require exponentially more hidden units.