

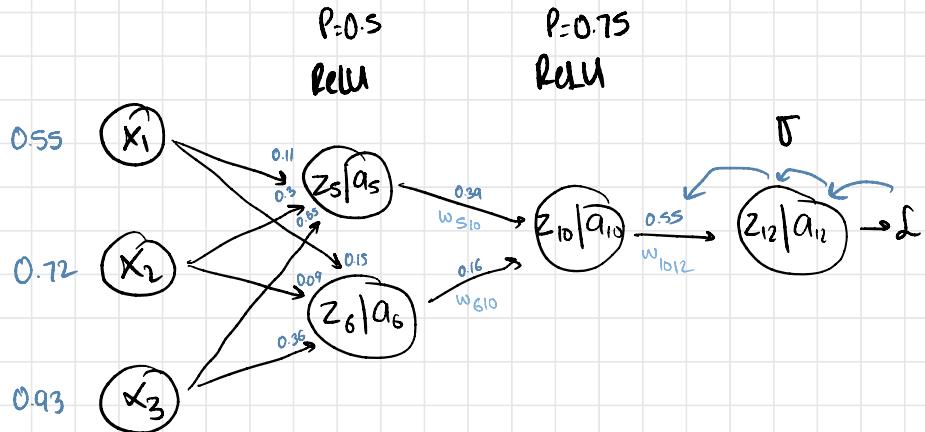
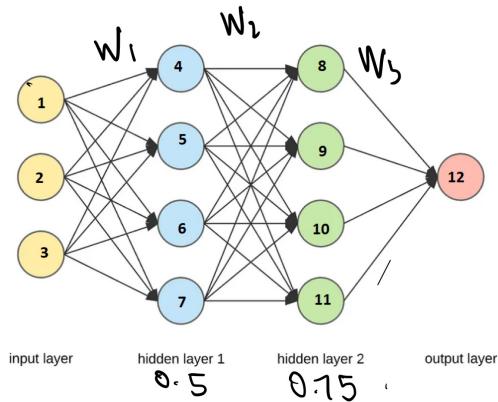
Use the dropout probability for 0.5 for hidden layer 1 and 0.75 for hidden layer 2. All forward and backward propagation formulas should be shown clearly.

i. [1 Point] Compute forward propagation on the thinned network.

ii. [1 Point] Compute the error using the squared error function if $y=0.98$.

iii. [3 Points] Perform backpropagation on a thinned network and update the weights.

iv. [1 Point] Test the updated network on new test data if $X=[[0.33], [0.45], [0.81]]$. Compute the output.



learning rate = 0.7

$$L = \frac{1}{2} (y - a)^2$$

$$\frac{\delta L}{\delta a} (a-y)$$

$$y = 0.98$$

4 FF

$$z_5 = (0.55)(0.11) + (0.72)(0.3) + (0.93)(0.65) = 0.881$$

$$a_5 = 0.881$$

$$z_6 = (0.55)(0.15) + (0.72)(0.09) + (0.93)(0.36) = 0.4821$$

$$a_6 = 0.4821$$

$$z_{10} = 0.881(0.39) + 0.4821(0.16) = 0.420726$$

$$a_{10} = 0.420726$$

$$z_{12} = 0.420726(0.55) = 0.2313905$$

$$a_{12} = \frac{1}{1+e^{-z_{12}}} = 0.55559$$

$a_{10} \cdot w_{1012}$

b Compute error using sq error function if $y = 0.98$

$$\text{LOSS} = \frac{1}{2} (0.98 - 0.55559)^2$$

$$= 0.089213$$

5 BP

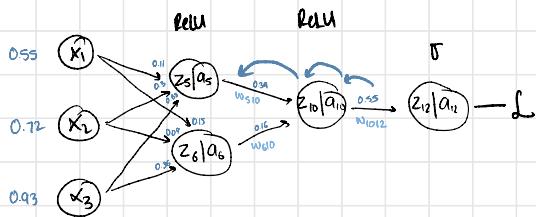
$$\frac{\delta L}{\delta a_{12}} = a_{12} - y \quad \frac{\delta a_{12}}{\delta z_{12}} = a_{12}(1-a_{12}) \quad \frac{\delta z_{12}}{\delta w_{1012}} = a_{10}$$

$$\frac{\delta L}{\delta w_{1012}} = \frac{\delta L}{\delta a_{12}} \cdot \frac{\delta a_{12}}{\delta z_{12}} \cdot \frac{\delta z_{12}}{\delta w_{1012}} = (a_{12} - y) \cdot (a_{12}(1-a_{12})) \cdot a_{10} = -0.04$$

$$w_{1012}^+ = 0.55 - 0.7(-0.04) = 0.519$$

$$\frac{\delta L}{\delta b_{12}} = (a_{12} - y) \cdot (a_{12}(1-a_{12})) = -0.10419$$

$$b_{12}^+ = 0 - 0.7(-0.10419) = 0.07294$$



$$\frac{\delta L}{\delta w_{510}} = \frac{\delta L}{\delta a_{10}} \cdot \frac{\delta a_{10}}{\delta z_{10}} \cdot \frac{\delta z_{10}}{\delta w_{510}}$$

$$= (\alpha_{12}-y) \cdot (\alpha_{12}(1-\alpha_{12})) \cdot W_{1012} \cdot (0,1) \cdot a_5$$

$$= -0.0505$$

$$\frac{\delta L}{\delta a_{10}} = \frac{\delta L}{\delta a_{12}} \cdot \frac{\delta a_{12}}{\delta z_{12}} \cdot \frac{\delta z_{12}}{\delta a_{10}}$$

$$= (\alpha_{12}-y) \cdot (\alpha_{12}(1-\alpha_{12})) \cdot W_{1012}$$

$$W_{510}^t = 0.39 - (0.7)(-0.0505)$$

$$= 0.42535$$

$$\frac{\delta L}{\delta w_{610}} = (\alpha_{12}-y) \cdot (\alpha_{12}(1-\alpha_{12})) \cdot W_{1012} \cdot (0,1) \cdot a_6$$

$$= 0.0276$$

$$W_{610}^t = 0.16 - (0.7)(0.0276)$$

$$= 0.17932$$

$$\frac{\delta L}{\delta b_{10}} = (\alpha_{12}-y) \cdot (\alpha_{12}(1-\alpha_{12})) \cdot W_{1012} \cdot (0,1)$$

$$= -0.0573$$

$$b_{10}^t = 0 - (0.7)(-0.0573)$$

$$= 0.0401$$

↳ Repeat of last left

$$w_{15}^t = 0.118603$$

$$w_{16}^t = 0.153528$$

$$w_{25}^t = 0.31127$$

$$w_{26}^t = 0.09462$$

$$w_{35}^t = 0.664546$$

$$w_{36}^t = 0.36596$$

$$b_5^t = 0.0156453$$

$$b_6^t = 0.006412$$

↳ Test updated Network

$$X_1 = 0.33$$

$$X_2 = 0.45$$

$$X_3 = 0.81$$

Updated weights, multiply by P

$$Z_4 = 0.33(0.21) + 0.45(0.36) + 0.81(0.105) \\ = 0.31635$$

Updated weights, multiplying

$$W_1 = \left[[0.21, 0.36, 0.105], [0.0593, 0.1556, 0.3272] \right. \\ \left. [0.076764, 0.047317, 0.1827], [0.16, 0.755, 0.06] \right]$$

$$W_2 = \left[[0.1275, 0.2175, 0.2625, 0.3825], \\ [0.6825, 0.06, 0.3825, 0.2775], \\ [0.4575, 0.3190125, 0.13449, 0.0075], \\ [0.205, 0.1575, 0.615, 0.105] \right]$$

$$W_3 = [0.44, 0.03, 0.581, 0.81]$$

$$Z_4 = 0.31635$$

$$a_4 = 0.31635$$

$$Z_5 = 0.3747528$$

$$a_5 = 0.3747528$$

$$Z_6 = 0.20124742$$

$$a_6 = 0.20124742$$

$$Z_7 = 0.21615$$

$$a_7 = 0.21615$$

$$Z_8 = 0.25734$$

$$a_8 = 0.25734$$

$$Z_9 = 0.37535$$

$$a_9 = 0.37535$$

$$Z_{10} = 0.293$$

$$a_{10} = 0.293$$

$$Z_{11} = 0.295647$$

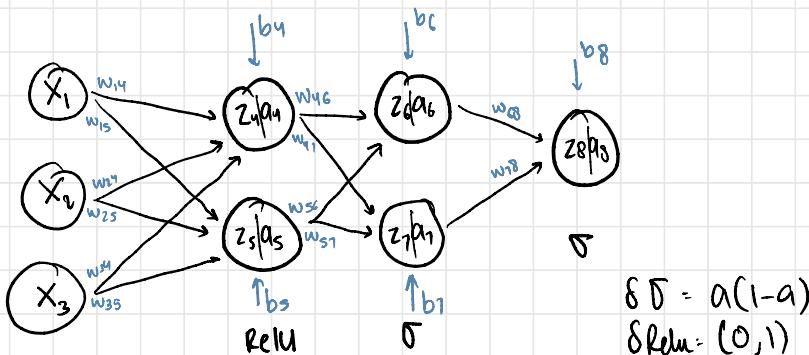
$$a_{11} = 0.295647$$

$$Z_{12} = 0.60713$$

$$\boxed{a_{12} = 0.64728}$$

21KU688

1. [20 Points] Draw 3-layer Neural Networks (Input layer is not counted here) with 2 nodes in each hidden layer and 1 output node. Input size is 3. Activation function is relu for 1st hidden layer and sigmoid for 2nd hidden and last layer. Write mathematical equations of all hidden and output layer nodes for a forward pass. Write down mathematical equations of the partial derivative of Loss w.r.t. all weights and bias for a backward pass. Simplify all these equations. Use (i) Squared Loss (ii) Cross Entropy Loss



Feed Forward

$$Z_4 = u_1 \cdot w_{14} + u_2 \cdot w_{24} + u_3 \cdot w_{34} + b_4$$

$$Z_5 = u_1 \cdot w_{15} + u_2 \cdot w_{25} + u_3 \cdot w_{35} + b_5$$

$$a_4 = \text{ReLU}(z_4) = \max(0, z_4)$$

$$a_5 = \text{ReLU}(z_5) = \max(0, z_5)$$

$$Z_6 = a_4 \cdot w_{46} + a_5 \cdot w_{56} + b_6$$

$$a_6 = \sigma(z_6) = \frac{1}{1+e^{-z_6}}$$

$$Z_7 = a_4 \cdot w_{47} + a_5 \cdot w_{57} + b_7$$

$$a_7 = \sigma(z_7) = \frac{1}{1+e^{-z_7}}$$

$$Z_8 = a_6 \cdot w_{68} + a_7 \cdot w_{78} + b_8$$

$$a_8 = \sigma(z_8) = \frac{1}{1+e^{-z_8}}$$

Backward Propagation

1) Squared Loss

$$L = \frac{1}{2} (y - a_8)^2$$

$$L = \frac{1}{2} (y - a_8)^2$$

$$\frac{\delta L}{\delta a_8} = \frac{1}{2} (2)(-1)(y - a_8)$$

$$\frac{\delta L}{\delta w_{68}} = \frac{1}{2}$$

$$\frac{\delta L}{\delta a_8} = a_8 - y$$

$$\frac{\delta a_8}{\delta z_8} = a_8(1-a_8)$$

$$\frac{\delta z_8}{\delta w_{68}} = a_6$$

$$\frac{\delta z_8}{\delta w_{78}} = a_7$$

$$\frac{\delta L}{\delta w_{68}} = \frac{\delta L}{\delta a_8} \cdot \frac{\delta a_8}{\delta z_8} \cdot \frac{\delta z_8}{\delta w_{68}}$$

$$\rightarrow (a_8 - y) \cdot a_8(1-a_8) \cdot a_6$$

$$w_{68}^+ = w_{68} - \alpha \frac{\delta L}{\delta w_{68}}$$

$$\rightarrow \frac{\delta L}{\delta w_{78}} = (a_8 - y) \cdot a_8(1-a_8) \cdot a_7$$

$$w_{78}^+ = w_{78} - \alpha \frac{\delta L}{\delta w_{78}}$$

$$\rightarrow \frac{\delta L}{\delta b_8} = (a_8 - y) \cdot a_8(1-a_8)$$

H1

$$\rightarrow \frac{\delta L}{\delta w_{46}} = \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta w_{46}}$$

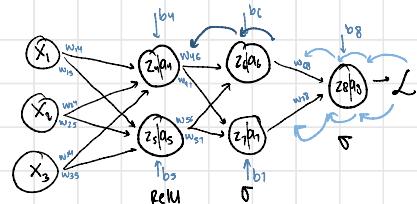
$$= (a_8 - y) (a_8(1-a_8)) \cdot w_{68} (a_6(1-a_6)) \underline{a_6}$$

$$\frac{\delta L}{\delta a_6} = \frac{\delta L}{\delta a_8} \times \frac{\delta a_8}{\delta z_8} \times \frac{\delta z_8}{\delta a_6}$$

$$(a_8 - y) (a_8(1-a_8)) \cdot w_{68}$$

$$\rightarrow \frac{\delta L}{\delta w_{56}} = (a_8 - y) (a_8(1-a_8)) \cdot w_{68} (a_6(1-a_6)) a_5$$

$$\frac{\delta L}{\delta b_6} = (a_8 - y) (a_8(1-a_8)) \cdot w_{68} (a_6(1-a_6))$$



$$\frac{\delta L}{\delta w_{47}} = \frac{\delta L}{\delta a_7} \cdot \frac{\delta a_7}{\delta z_7} \cdot \frac{\delta z_7}{\delta w_{47}}$$

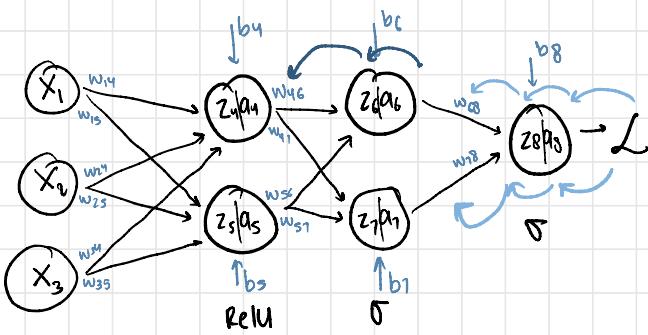
$$= (a_8 - y) \cdot a_8(1-a_8) \cdot w_{78} \cdot a_7(1-a_7) \cdot a_6$$

$$\frac{\delta L}{\delta a_7} = \frac{\delta L}{\delta a_8} \cdot \frac{\delta a_8}{\delta z_8} \cdot \frac{\delta z_8}{\delta a_7}$$

$$= (a_8 - y) \cdot a_8(1-a_8) \cdot w_{78}$$

$$\frac{\delta L}{\delta w_{51}} = (a_8 - y) \cdot a_8(1-a_8) \cdot w_{78} \cdot a_7(1-a_7) \cdot a_5$$

$$\frac{\delta L}{\delta b_7} = (a_8 - y) \cdot a_8(1-a_8) \cdot w_{78} \cdot a_7(1-a_7)$$



H2

$$\frac{\delta L}{\delta w_{14}} \cdot \frac{\delta L}{\delta a_4} \cdot \frac{\delta a_4}{\delta z_4} \cdot \frac{\delta z_4}{\delta w_{14}}$$

$$\begin{aligned}\frac{\delta L}{\delta a_4} &= \frac{\delta L}{a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{a_4} + \frac{\delta L}{a_7} \cdot \frac{\delta a_7}{\delta z_7} \cdot \frac{\delta z_7}{a_4} \\ &= (\alpha_8 - y)(\alpha_8(1-\alpha_8)) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} + \\ &\quad (\alpha_7 - y)(\alpha_7(1-\alpha_7)) \cdot w_{78} \cdot a_7(1-a_7) \cdot w_{47}\end{aligned}$$

$$\frac{\delta}{\delta w_{14}} \left[(\alpha_8 - y)(\alpha_8(1-\alpha_8)) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} + (\alpha_7 - y)(\alpha_7(1-\alpha_7)) \cdot w_{78} \cdot a_7(1-a_7) \cdot w_{47} \right].$$

$$\frac{\delta}{\delta w_{24}} \left[(\alpha_8 - y)(\alpha_8(1-\alpha_8)) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} + (\alpha_7 - y)(\alpha_7(1-\alpha_7)) \cdot w_{78} \cdot a_7(1-a_7) \cdot w_{47} \right].$$

$$\frac{\delta}{\delta w_{34}} \left[(\alpha_8 - y)(\alpha_8(1-\alpha_8)) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} + (\alpha_7 - y)(\alpha_7(1-\alpha_7)) \cdot w_{78} \cdot a_7(1-a_7) \cdot w_{47} \right].$$

$$\frac{\delta}{\delta w_{15}} \left[(\alpha_8 - y)(\alpha_8(1-\alpha_8)) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{56} + (\alpha_7 - y)(\alpha_7(1-\alpha_7)) \cdot w_{78} \cdot a_7(1-a_7) \cdot w_{57} \right].$$

$$\frac{\delta}{\delta w_{25}} \left[(\alpha_8 - y)(\alpha_8(1-\alpha_8)) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{56} + (\alpha_7 - y)(\alpha_7(1-\alpha_7)) \cdot w_{78} \cdot a_7(1-a_7) \cdot w_{57} \right].$$

$$\frac{\delta}{\delta w_{35}} \left[(\alpha_8 - y)(\alpha_8(1-\alpha_8)) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{56} + (\alpha_7 - y)(\alpha_7(1-\alpha_7)) \cdot w_{78} \cdot a_7(1-a_7) \cdot w_{57} \right].$$

ii) Cross entropy

$$L = -[y \log a_8 + (1-y) \log (1-a_8)]$$

$\frac{\delta L}{\delta w_{68}} = \frac{\delta L}{\delta a_8} \cdot \frac{\delta a_8}{\delta z_8} \cdot \frac{\delta z_8}{\delta w_{68}}$
 $\frac{\delta L}{\delta w_{68}} = \left(\frac{y}{a_8} + \frac{1-y}{1-a_8} \right) \cdot a_8(1-a_8) a_6$

$\frac{\delta L}{\delta b_8} = (a_8 - y)$
 $\frac{\delta L}{\delta z_8} = (a_8 - y) a_7$
 $\frac{\delta L}{\delta a_8} = (a_8 - y) a_6$

H2

$$\frac{\delta L}{\delta w_{46}} = \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta w_{46}}$$

$$\frac{\delta L}{\delta w_{46}} = (a_8 - y) \cdot w_{68} \cdot a_6 (1-a_6) a_4$$

$$\frac{\delta L}{\delta w_{56}} = (a_8 - y) \cdot w_{68} \cdot a_6 (1-a_6) a_5$$

$$\frac{\delta L}{\delta b_6} = (a_8 - y) \cdot w_{68} \cdot a_6 (1-a_6)$$

$$\frac{\delta L}{\delta w_{41}} = (a_8 - y) \cdot w_{78} \cdot a_7 (1-a_7) \cdot a_4$$

$$\frac{\delta L}{\delta w_{51}} = (a_8 - y) \cdot w_{78} \cdot a_7 (1-a_7) \cdot a_5$$

$$\frac{\delta L}{\delta b_7} = (a_8 - y) \cdot w_{78} \cdot a_7 (1-a_7)$$

H2

$$\frac{\delta L}{\delta w_{14}} = \frac{\delta L}{\delta a_9} \cdot \frac{\delta a_9}{\delta z_4} \cdot \frac{\delta z_4}{\delta w_4}$$

$$\frac{\delta L}{\delta w_4} = \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta a_9} + \frac{\delta L}{\delta a_7} \cdot \frac{\delta a_7}{\delta z_7} \cdot \frac{\delta z_7}{\delta a_1}$$

$$= (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} + [a_8 - y] \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{47}$$

$$\frac{\delta L}{\delta w_{14}} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{47} \end{array} \right] \times (0, 1) \cdot v_1$$

$$\frac{\delta L}{\delta w_{21}} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{47} \end{array} \right] \times (0, 1) \cdot v_2$$

$$\frac{\delta L}{\delta w_{34}} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{47} \end{array} \right] \times (0, 1) \cdot v_3$$

$$\frac{\delta L}{\delta p_4} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{46} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{47} \end{array} \right] \times (0, 1)$$

$$\frac{\delta L}{\delta a_8} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{56} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{57} \end{array} \right] \times (0, 1) \cdot v_1$$

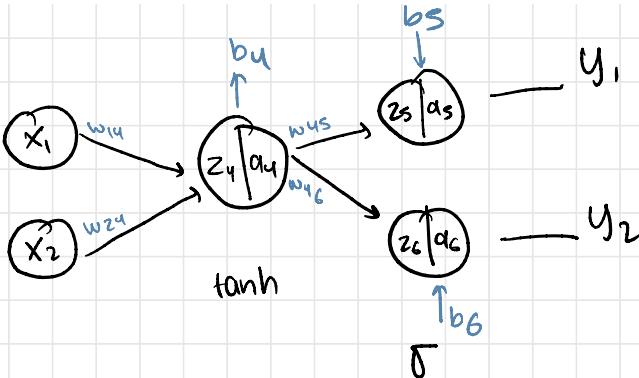
$$\frac{\delta L}{\delta z_5} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{56} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{57} \end{array} \right] \times (0, 1) \cdot v_1$$

$$\frac{\delta L}{\delta z_5} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{56} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{57} \end{array} \right] \times (0, 1) \cdot v_2$$

$$\frac{\delta L}{\delta z_5} = \left[\begin{array}{l} (a_8 - y) \cdot w_{68} \cdot a_6(1-a_6) \cdot w_{56} \\ (a_8 - y) \cdot w_{78} \cdot (a_7(1-a_7)) \cdot w_{57} \end{array} \right] \times (0, 1) \cdot v_3$$

H 0

2. [10 Points] Draw 2-layer Neural Networks with 1 node in the hidden layer and 2 output nodes. Input size is 2. The activation function in the hidden layer is tanh while in the output layer, sigmoid is used. Write mathematical equations of all hidden and output layer nodes for a forward pass. Write down mathematical equations of the partial derivative of Loss w.r.t. all weights and bias for a backward pass. Simplify all these equations. Use Cross Entropy Loss.



Feed Forward

$$z_4 = u_1 w_{14} + u_2 w_{24} + b_4 \rightarrow a_4 = \tanh(z_4)$$

$$z_5 = a_4 w_{45} + b_5 \rightarrow a_5 = \sigma(z_5) = \frac{1}{1 + e^{-z_5}}$$

$$z_6 = a_4 w_{46} + b_6 \rightarrow a_6 = \sigma(z_6) = \frac{1}{1 + e^{-z_6}}$$

Backward

↳ Cross Entropy loss

$$-y \log \hat{y} + (1-y) \log(1-\hat{y})$$

$$L = -[y \log a_5 + (1-y) \log(1-a_5)] - [y \log a_6 + (1-y) \log(1-a_6)]$$

$$\frac{\delta L}{a_5} = \frac{-y}{a_5} + \frac{1-y}{1-a_5} \quad \frac{\delta L}{a_6} = \frac{-y}{a_6} + \frac{1-y}{1-a_6}$$

$$\frac{\delta L}{w_{45}} = \frac{\delta L}{a_5} \cdot \frac{\delta a_5}{\delta z_5} \cdot \frac{\delta z_5}{\delta w_{45}}$$

$$= \left[-\frac{y}{a_5} + \frac{1-y}{1-a_5} \right] \cdot [a_5(1-a_5)] \cdot a_4$$

$$= -\frac{(y)(1-a_5) + a_5(1-y)}{a_5(1-a_5)} \cdot [a_5(1-a_5)] \cdot a_4$$

$$= -y + a_5 y + a_5 - a_5 y$$

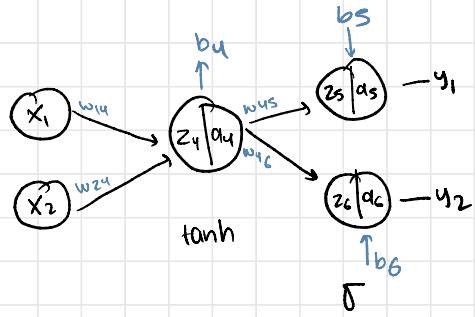
$$= (a_5 - y) a_4$$

$$\frac{\delta L}{\delta w_{45}} = (a_5 - y_1) a_4$$

$$\frac{\delta L}{\delta b_5} = a_5 - y_1$$

$$\frac{\delta L}{\delta w_{46}} = (a_6 - y_2) a_4$$

$$\frac{\delta L}{\delta b_6} = a_6 - y_2$$



H1

$$\frac{\delta L}{\delta a_4} = \frac{\delta L}{\delta a_5} \cdot \frac{\delta a_5}{\delta z_5} \cdot \frac{\delta z_5}{\delta a_4} + \frac{\delta L}{\delta a_6} \cdot \frac{\delta a_6}{\delta z_6} \cdot \frac{\delta z_6}{\delta a_4}$$

$$\begin{aligned} \frac{\delta L}{\delta w_{45}} &= \frac{\delta L}{\delta a_4} \cdot \frac{\delta a_4}{\delta z_5} \cdot \frac{\delta z_5}{\delta w_{45}} \\ &= [(a_5 - y_1)w_{45} + (a_6 - y_2)w_{46}] \cdot (1 - a_4^2) \cdot n_1 \end{aligned}$$

$$\frac{\delta L}{\delta w_{24}} = [(a_5 - y_1)w_{45} + (a_6 - y_2)w_{46}] \cdot (1 - a_4^2) \cdot n_2$$

$$\frac{\delta L}{\delta b_4} = [(a_5 - y_1)w_{45} + (a_6 - y_2)w_{46}] \cdot (1 - a_4^2)$$

3. [20 Points] Complete the following exercise using formulas computed in Question #1 for

weights of L1

$W_1 = [[0.42, 0.72, 0.21], [0.11, 0.3, 0.65]]$

$B_1 = 0.01$ (for all neurons)

weights of L2

$W_2 = [[0.17, 0.29], [0.91, 0.08]]$

$B_2 = 0.05$ (for all neurons)

weights of L3

$W_3 = [[0.61, 0.39]]$ $B_3 = 0.08$

Learning rate = 0.7

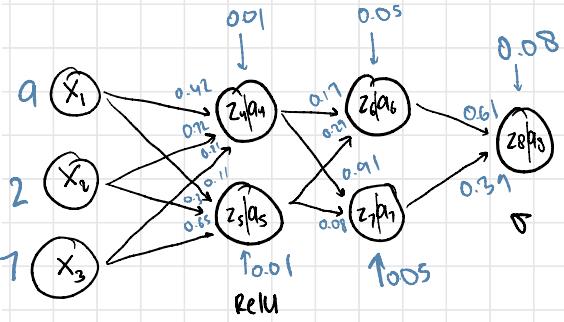
(i) Squared Loss Function

(ii) Cross Entropy Loss Function

F1	F2	F3	Class
4 th digit of student ID + 1	2	2 nd digit of student ID + 1	A
5	8	3 rd digit of student ID + 1	B

Be like a neural network. Learn from your mistakes :)

21K-4688
921A



$$\alpha = 0.7$$

Q3) 1. Squared Loss Function

$$Z_4 = 0.42(a) + 2(0.72) + (7)(0.21) = 6.7$$

$$Z_5 = 6.15$$

$$Z_6 = 2.9725$$

$$Z_7 = 6.639$$

$$Z_8 = 1.0498$$

$$a_4 = 6.7$$

$$a_5 = 6.15$$

$$a_6 = 0.9513$$

$$a_7 = 0.9987$$

$$a_8 = 0.7107$$

$$\delta w_{14} = -0.0123$$

$$w_{14} = 0.4226$$

$$\delta b_4 = -0.0014$$

$$b_4 = 0.011$$

$$\delta w_{24} = -0.0027$$

$$w_{24} = 0.7219$$

$$\delta b_5 = -0.0021$$

$$b_5 = 0.0115$$

$$\delta w_{34} = -0.0096$$

$$w_{34} = 0.2161$$

$$\delta b_6 = -0.0073$$

$$b_6 = 0.0551$$

$$\delta w_{15} = -0.0197$$

$$w_{15} = 0.1234$$

$$\delta b_7 = -0.0001$$

$$b_7 = 0.0561$$

$$\delta w_{25} = -0.0043$$

$$w_{25} = 0.303$$

$$\delta b_8 = -0.2593$$

$$b_8 = 0.2615$$

$$\delta w_{35} = -0.0199$$

$$w_{35} = 0.6605$$

$$\delta w_{46} = -0.0041$$

$$w_{46} = 0.2041$$

$$\delta w_{56} = -0.0045$$

$$w_{56} = 0.3215$$

$$\delta w_{47} = -0.0009$$

$$w_{47} = 0.0102$$

$$\delta w_{57} = -0.0002$$

$$w_{57} = 0.0806$$

$$\delta w_{68} = -0.2466$$

$$w_{68} = 0.7926$$

$$\delta w_{78} = -0.2589$$

$$w_{78} = 0.5712$$

2) Cross Entropy

chart

$$Z_4 = 0.42(a) + 2(0.72) + (7)(0.21) = 6.7$$

$$a_4 = 6.7$$

$$Z_5 = 6.15$$

$$a_5 = 6.15$$

$$Z_6 = 2.9725$$

$$a_6 = 0.9513$$

$$Z_7 = 6.639$$

$$a_7 = 0.9987$$

$$Z_8 = 0.8336$$

$$a_8 = 0.7107$$

$$\delta w_{14} = -0.0024$$

$$w_{14} = 0.4217$$

$$\delta w_{24} = -0.0005$$

$$w_{24} = 0.7204$$

$$\delta w_{34} = -0.0013$$

$$w_{34} = 0.2113$$

$$\delta w_{15} = -0.0057$$

$$w_{15} = 0.1126$$

$$\delta b_4 = -0.0003$$

$$b_4 = 0.0102$$

$$\delta w_{25} = -0.008$$

$$w_{25} = 0.3006$$

$$\delta b_5 = -0.0004$$

$$b_5 = 0.0103$$

$$\delta w_{35} = -0.0021$$

$$w_{35} = 0.652$$

$$\delta b_6 = -1.4067$$

$$b_6 = 0.051$$

$$\delta w_{46} = -0.0094$$

$$w_{46} = 0.1766$$

$$\delta w_{56} = -0.0087$$

$$w_{56} = 0.2961$$

$$\delta w_{47} = -0.0062$$

$$w_{47} = 0.9101$$

$$\delta w_{57} = -0.0062$$

$$w_{57} = 0.0801$$

$$\delta w_{68} = -0.074$$

$$w_{68} = 0.6432$$

$$\delta w_{78} = -0.0491$$

$$w_{78} = 0.4248$$

CONVOLUTION NEURAL NETWORKS 6

CNN

An image is nothing but a matrix of pixel values, right? So why not just flatten the image (e.g. 3x3 image matrix into a 9x1 vector) and feed it to a Multi-Level Perceptron for classification purposes? Uh.. not really.

To capture Spatial/temporal dependencies.



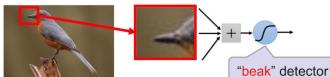
1
1
0
4
2
1
0
2
1

- Flattening ignores spatial dependencies.
- CNNs preserve spatial and temporal relationships by using convolutions to detect local patterns.

NN with images

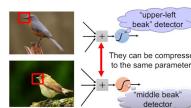
Some patterns are much smaller than the whole image.

Can represent a small region with fewer parameters



Same pattern appears in different places: They can be compressed!

What about training a lot of such "small" detectors and each detector must "move around".



Edge Detection

Edge is Where Change Occurs

Change is measured by derivative in 1D

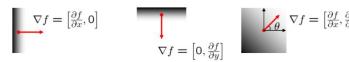
Biggest change, derivative has maximum magnitude
Or 2nd derivative is zero.

image gradient

The gradient of an image:

$$\nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$$

The gradient points in the direction of most rapid change in intensity



The gradient direction is given by:

$$\theta = \tan^{-1} \left(\frac{\frac{\partial f}{\partial y}}{\frac{\partial f}{\partial x}} \right)$$

- how does this relate to the direction of the edge?

The edge strength is given by the gradient magnitude.

$$\|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$

Gradient operators

$$\begin{array}{cc} \Delta_1 & \Delta_2 \\ \begin{matrix} 0 & 1 \\ -1 & 0 \end{matrix} & \begin{matrix} 1 & 0 \\ 0 & -1 \end{matrix} \end{array} \quad \begin{array}{cc} \Delta_1 & \Delta_2 \\ \begin{matrix} -1 & 0 & 1 \\ 1 & 1 & 1 \\ -1 & 0 & 1 \end{matrix} & \begin{matrix} 1 & 1 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & -1 \end{matrix} \end{array}$$

(a) (b)

$$\begin{array}{cc} \Delta_1 & \Delta_2 \\ \begin{matrix} -1 & 0 & 1 \\ 2 & 0 & 2 \\ -1 & 0 & 1 \end{matrix} & \begin{matrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{matrix} \end{array} \quad \begin{array}{cc} \Delta_1 & \Delta_2 \\ \begin{matrix} -3 & -1 & 1 & 3 \\ -3 & -1 & 1 & 3 \\ 1 & 1 & 1 & 1 \\ -3 & -1 & 1 & 3 \end{matrix} & \begin{matrix} 3 & 3 & 3 & 3 \\ -3 & -3 & -3 & -3 \\ 1 & 1 & 1 & 1 \\ -3 & -3 & -3 & -3 \end{matrix} \end{array}$$

(c) (d)

(a): Roberts' cross operator (b): 3x3 Prewitt operator

(c): Sobel operator (d) 4x4 Prewitt operator

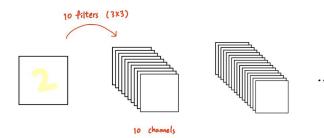
CONVOLUTION

Convolution mathematically combines two signals (an input image and a filter/kernel) to produce a new signal (feature map).

CNNs

A CNN is a neural network with some convolutional layers (and some other layers). A convolutional layer has a number of filters that does convolutional operation.

Result array from a convolution operation using a filter is called a feature map. If you use 10 filters, then there will be 10 output channels.

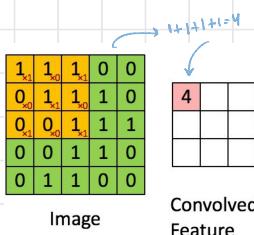
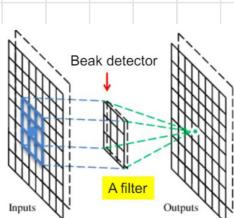


How Convolution Works (Diagram explained):

- A filter (small matrix) slides across the input image (larger matrix), performing element-wise multiplication followed by summation.
- The resulting smaller matrix is a feature map, capturing local features such as edges and textures.

Filters/Kernels:

- Filters identify various features (edges, textures, patterns).
- Initially random, filters are trained through backpropagation to recognize specific meaningful patterns.



1	0	0	0	0	0	1
0	1	0	0	1	0	0
0	0	1	1	0	0	0
1	0	0	0	1	0	0
0	1	0	0	1	0	0
0	0	1	0	1	0	0

6 x 6 image

1	-1	-1
-1	1	1
-1	-1	1

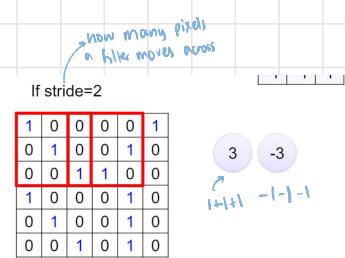
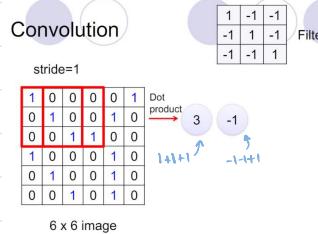
Filter 1

-1	1	-1
1	-1	1
-1	1	-1

Filter 2

⋮ :

Each filter detects a small pattern (3 x 3).



Convolution

stride=1

0	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	1
-1	-1	1

Filter 1

3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

6 x 6 image

Convolution

stride=1

1	0	0	0	0	0
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

Filter 2

-1	1	-1
1	-1	1
-1	1	-1

1	-1	-1
-1	1	1
-1	-1	1

Repeat this for each filter

1	0	0	0	0	0
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

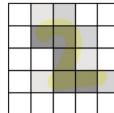
6 x 6 image

Two 4 x 4 Images
Forming 2 x 4 x 4 matrix

FOR GREY IMAGES

→ 1 filter

2



1	0.6	0.3	1	1
1	0	0	0.3	1
1	1	0	0.6	1
1	0.6	0	0	0.3
1	1	1	1	1

1	0.6	0.3	1	1
1	0	0	0.3	1
1	1	0	0.6	1
1	0.6	0	0	0.3
1	1	1	1	1

0	1	0
1	2	1
0	1	0

multiply numbers in same location

1	0	0.6	0.3	1	1
1	0	0	0.3	1	1
1	0	0	0.6	1	1
1	0.6	0	0	0.3	1
1	1	1	1	1	1

sum all

next cell	→			
1	0.6	0.3	1	1
1	0	0	0.3	1
1	1	0	0.6	1
1	0.6	0	0	0.3
1	1	1	1	1

0	0.3	0
0	0	0.3
0	0	0

2.6 0.6

1	0.6	0.3	1	1
1	0	0	0.3	1
1	1	0	0.6	1
1	0.6	0	0	0.3
1	1	1	1	1

This sliding operation
is called convolution

1	0.6	0.3	1	1
1	0	0	0.3	1
1	1	0	0.6	1
1	0.6	0	0	0.3
1	1	1	1	1

2.6 0.6 3.2
3.6

filter

1	0.6	0.3	1	1
1	0	0	0.3	1
1	1	0	0.6	1
1	0.6	0	0	0.3
1	1	1	1	1

2.6	0.6	3.2
3.6	1.6	2.5
4.2	1.6	1.4

data

result

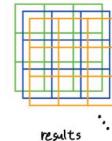
FOR GREY IMAGES

filters



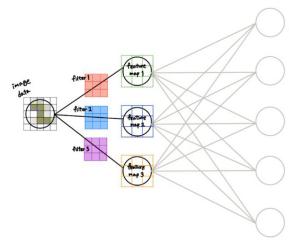
1	0.6	0.3	1	1
1	0	0	0.3	1
1	1	0	0.6	1
1	0.6	0	0	0.3
1	1	1	1	1

data



results

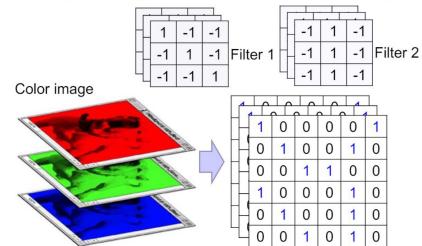
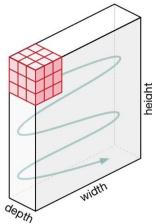
FOR GREY IMAGES



For RGB: color images

→ 3 filters
↓ ↓ ↓
R G B

- The filter moves to the right with a certain Stride Value till it parses the complete width.
- Moving on, it hops down to the beginning (left) of the image with the same Stride Value and repeats the process until the entire image is traversed.



0	0	0	0	0	0	0	...
0	156	155	156	158	158	...	
0	153	154	157	159	159	...	
0	149	151	155	158	159	...	
0	146	146	149	153	158	...	
0	145	143	143	148	158	...	
...	

Input Channel #1 (Red)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1

308

$$156 - 155 + \\ 153 + 154$$

0	0	0	0	0	0	0	...
0	167	166	167	169	169	...	
0	164	165	168	170	170	...	
0	160	162	166	169	170	...	
0	156	156	159	163	168	...	
0	155	153	153	158	168	...	
...	

Input Channel #2 (Green)

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2

-498

$$-167 - 166 + \\ -165$$

0	0	0	0	0	0	0	...
0	163	162	163	165	165	...	
0	160	161	164	166	166	...	
0	156	158	162	165	166	...	
0	155	155	158	162	167	...	
0	154	152	152	157	167	...	
...	

Input Channel #3 (Blue)

0	1	1
0	1	0
1	-1	1

Kernel Channel #3

164

$$163 - 160 + 161$$

Bias = 1

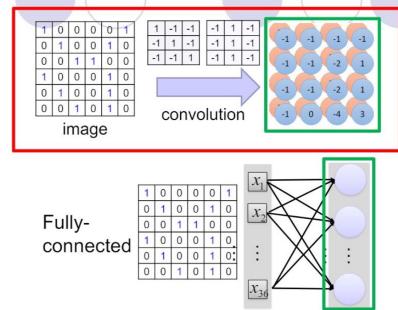
-25			...
			...
			...
			...
...

$$\text{Input} \times \text{Input} \times \text{Input} \times \text{Bias} \downarrow \\ \text{Filter R} + \text{Filter G} + \text{Filter B} + 1 = \text{Output}$$

Convolution v.s. Fully Connected

Fully Connected vs. CNN:

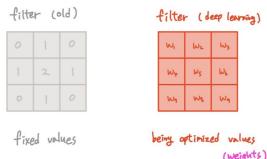
- Fully Connected Networks:** Every neuron is connected to every other neuron in the next layer. This requires an immense number of parameters for images, leading to complexity and inefficient learning.
- CNNs:** Introduce layers that recognize local patterns, significantly reducing complexity by focusing only on small areas of the input at a time.



SO WHAT KIND OF VALUES SHOULD BE IN THE FILTERS?

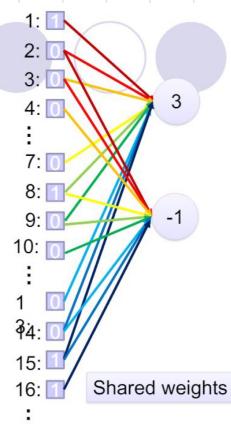
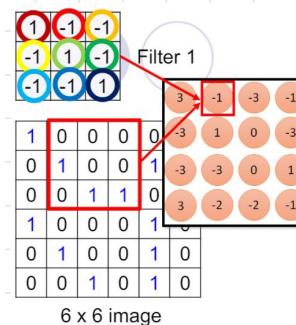
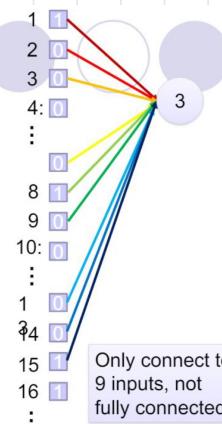
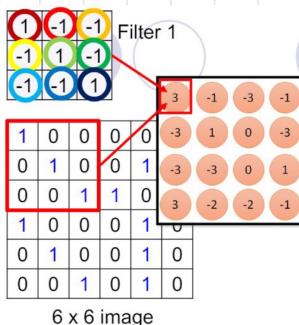
Before the era of deep learning, the values in the filter were fixed. Have you ever used "sharpen filter" in a photo editor?

The values in a convolution filter are the weights of a Convolutional Neural Network.



CNN PROS

↳ reduces complexity ↳ few parameters
 ↳ shared weights



↳ Shared weights

CNN

Filters: Values and Training

- Initially, filters have random values.
- Through training (gradient descent), these values become optimized to detect meaningful patterns (e.g., edges, textures, shapes).

CNN Layers and Feature Extraction (Diagram explained):

- First Layer: Captures basic features (edges, colors).
- Middle Layers: Identify intermediate features (shapes, complex patterns).
- Final Layers: Capture high-level semantic features (objects, scenes).

Convolution layer takes an input feature map of dimension $W \times H \times N$ and produces an output feature map of dimension $\tilde{W} \times \tilde{H} \times M$

Each layer is defined using following parameters:

- # Input channels (N)
- # Output channels (M)
- Kernel size
- Padding
- Stride → how many pixels a filter moves across

↳ Stride 1: Large feature map
↳ more detailed

↳ Stride 2: Smaller feature map
↳ more faster → as it jumps 2 pixels
↳ reduces complexity
↳ lower dim of data

Problem: In AlexNet, the input image is of size 227x227x3. The first convolutional layer has 96 kernels of size 11x11x3. The stride is 4 and padding is 0. Therefore, the size of the output image right after the first bank of convolutional layers is

$$O = \frac{227 - 11 + 2 \times 0}{4} + 1 = 55$$

So, the output image is of size 55x55x96 (one channel for each kernel)

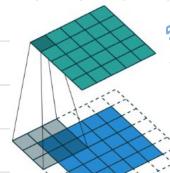


Figure: In this example, 5x5 input is convolved with 3x3 kernel with stride=padding=1 to produce an output of size 5x5.

$$O: \frac{5 + 2(1) - 3}{1} + 1$$

$$= 5 \times 5$$

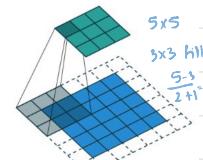


Figure: In this example, 5x5 input is convolved with 5x5 kernel with stride=2 and padding=1 to produce an output of size 3x3.

$$O: \frac{5 + 2(1) - 5}{2} + 1$$

$$\cdot \frac{5-2}{2} + 1 \\ = 3 \times 3$$

dim of $\text{OUTPUT} = \frac{\text{input size}}{\text{stride}} + \text{padding} + 1$ feature map

$$Q) S: 1, P: 0, K: 3, \text{input: } 28$$

$$\text{output: } \frac{28 + 2(0) - 3}{1} + 1$$

$$= 26$$

$$= 26 \times 26$$

$$O: \frac{227 + 2(0) - 11}{4} + 1$$

$$= 55$$

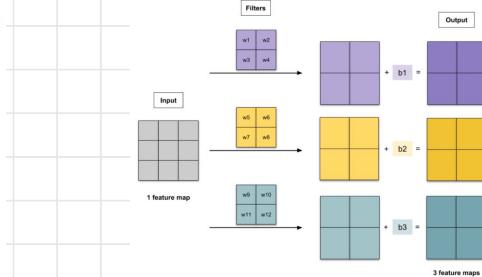
Parameters

num. Parameters:	Weights $i \times f \times f \times o$	+ biases $+ o$
	no of input maps / channels only 1 map	no of output maps 1 filter is applied to each output map

grey scale image - 1 channel
RGB image - 3 channels

Greyscale image with 2×2 filter, output 3 channels

- $i = 1$ (greyscale has only 1 channel)
- $f = 2$
- $o = 3$
- num_params
 $= [i \times (f \times f) \times o] + o$
 $= [1 \times (2 \times 2) \times 3] + 3$
 $= 15$



RGB image with 2×2 filter, output of 1 channel

- $i = 3$ (3 channels)
- $f = 2$
- $o = 1$
- num_params
 $= [i \times (f \times f) \times o] + o$
 $= [3 \times (2 \times 2) \times 1] + 1$
 $= 13$

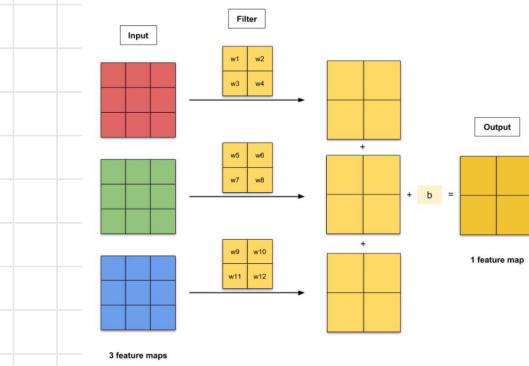
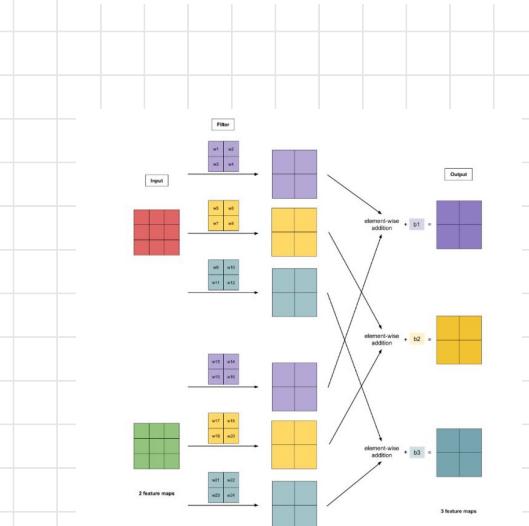


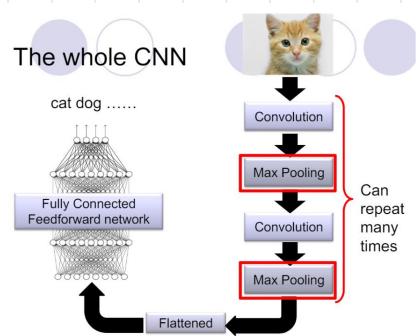
Image with 2 channels, with 2×2 filter, and output of 3 channels

- $i = 2$ (2 channels)
- $f = 2$
- $o = 3$
- num_params
 $= [i \times (f \times f) \times o] + 3$
 $= [2 \times (2 \times 2) \times 3] + 3$
 $= 27$



POOLING LAYERS

- Similar to the Convolutional Layer, the Pooling layer is responsible for
- Reduce dimensionality (height and width) of feature maps.
 - Lower computational complexity and control overfitting.



TYPES OF POOLING

- Max Pooling: Retains only the maximum value in a pooling window, suppressing noise.
- Min Pooling: Retains the minimum value (rarely used).
- Average Pooling: Retains the average value, smoothing feature maps.

MAX POOL

Max Pooling also performs as a **Noise Suppressor**.

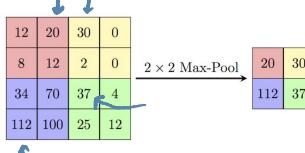
It discards the noisy activations altogether and also performs de-noising along with dimensionality reduction.

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

- A pooling window (e.g., 2x2) slides across the feature map, taking the maximum value in each region.
- Result: reduced spatial dimension, preserving the strongest activations.

Max Pool



Max Pooling

Features		Pool size	Stride
-4	0	-2	4
3	1	0	2
1	0	1	1
4	6	5	1
-1	2	0	0

3	4
6	5

-4	-2
-1	0

0	1
2	1

Why Pooling

- Subsampling pixels will not change the object



We can subsample the pixels to make image smaller → fewer parameters to characterize the image

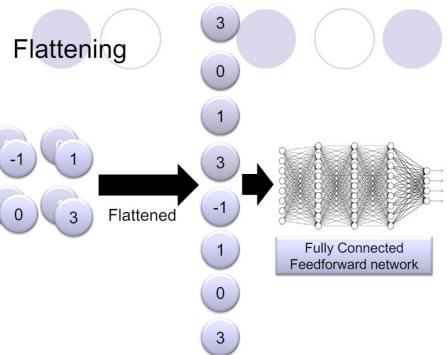
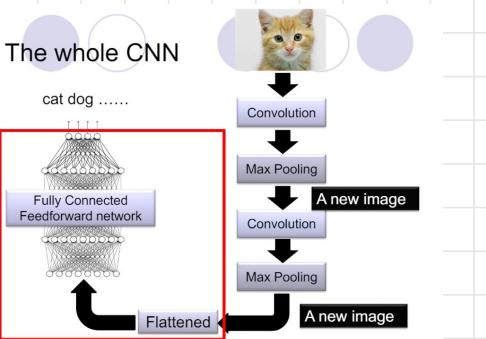
CNN

A CNN compresses a fully connected network in following ways:

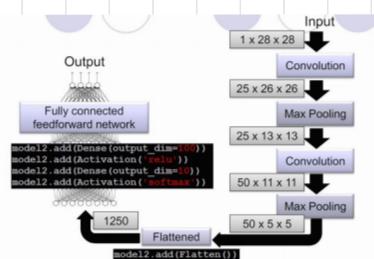
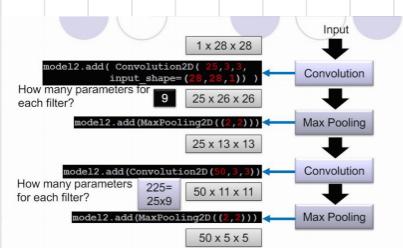
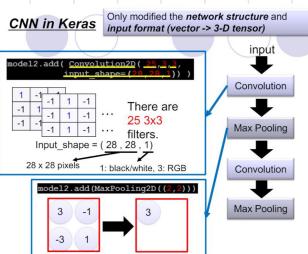
- Reducing number of connections
- Shared weights on the edges
- Max pooling further reduces the complexity

Three Fundamental CNN Concepts:

- Local receptive fields: Individual neurons respond to specific areas, capturing local dependencies.
- Shared weights: Reduces the number of parameters by reusing filter weights across the image.
- Spatial/Temporal sub-sampling: Reduces the feature map size through pooling, maintaining important information while reducing complexity.



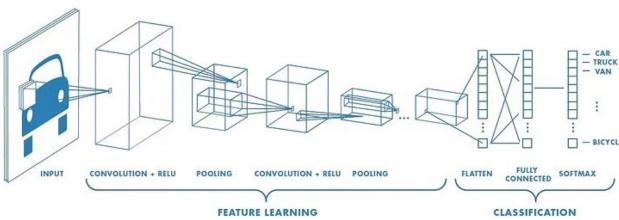
CNN in Keras



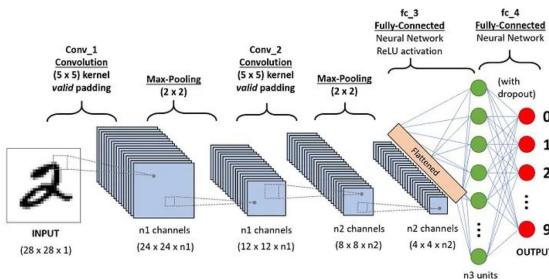
Practical Summary (What Happens in CNN training):

- Initialize filters randomly.
- Apply convolution to create initial feature maps.
- Apply activation functions (e.g., ReLU).
- Use pooling layers to reduce feature dimensions.
- Connect resulting features to fully connected layers for final classification.
- Optimize filter weights through backpropagation to improve feature detection capabilities and predictions.

CNN



CNN



Conclusion:

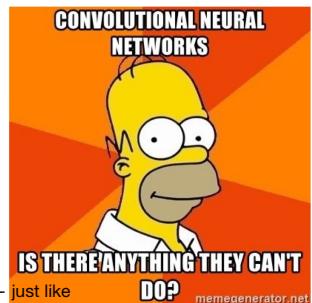
- CNNs are a powerful type of neural network optimized for handling images by exploiting spatial relationships, significantly outperforming traditional fully connected neural networks in visual recognition tasks.

what are spacial features?

to the way pixels are arranged and how they relate to one another in terms of position.
e.g

- An edge is a transition between light and dark pixels in close proximity.
- A face in a photo has eyes above the nose, and the nose above the mouth — that's a spatial relationship between facial features.
- In a photo of a cat, the ears are usually at the top of the head — their position relative to other features matters.

CNNs use this spatial information to understand shapes, edges, textures, and ultimately objects — just like our brains do.



CNN ARCHITECTURES

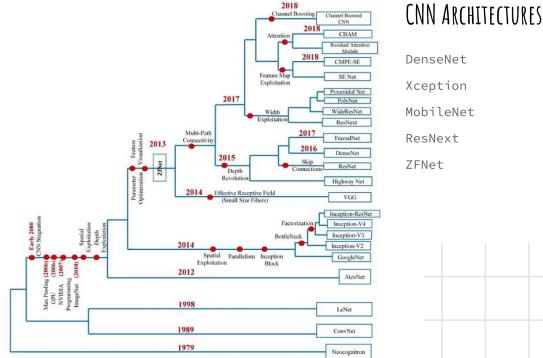
7

CNN Architectures

CNN Architectures refer to various structured designs of convolutional neural networks tailored to perform specific tasks efficiently (classification, detection, segmentation).

Key architectures discussed in the slides are:

CNN ARCHITECTURES



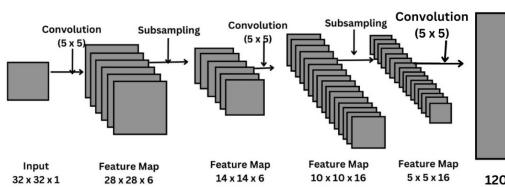
RESEARCH PAPERS

RESEARCH PAPERS EVERYWHERE

1. LeNet

- Early CNN architecture developed by Yann LeCun.
- Initially created for recognizing handwritten digits.
- Composed of convolutional layers, pooling layers, and fully connected layers.

LENET

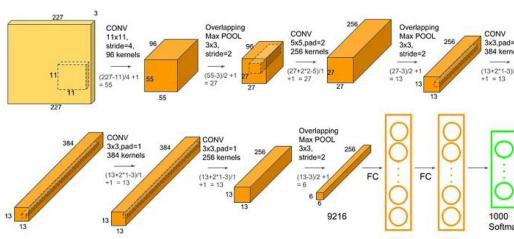


2. AlexNet (ILSVRC 2012 Winner)

- Significantly deeper and larger CNN than LeNet.
- Introduced ReLU activation function to tackle vanishing gradients.
- First CNN to leverage GPU acceleration.
- Achieved breakthrough performance on ImageNet classification.

ALEX NET

ILSVRC 2014 WINNER



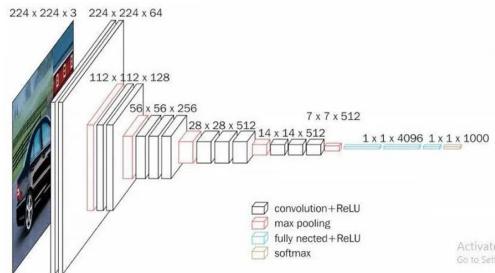
3. VGG-16 (ILSVRC 2014 Runner-up)

- Characterized by its simplicity and uniform structure.
- Uses small (3x3) convolution filters stacked to increase depth.
- 16 layers (13 convolutional, 3 fully connected).

VGG -16

ILSVRC 2014 Runner Up

Visual Geometry Group



Input



Activate
Go to Setti

4. GoogleNet (ILSVRC 2014 Winner)

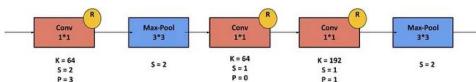
- Introduced the Inception module, significantly improving efficiency.
- Features parallel convolutions of different sizes (1×1 , 3×3 , 5×5), concatenating outputs.
- Reduces dimensionality with Stem Block (initial layer sequence reducing size and complexity) and multiple stacked Inception Blocks.

GOOGLENET

ILSVRC 2014 Winner



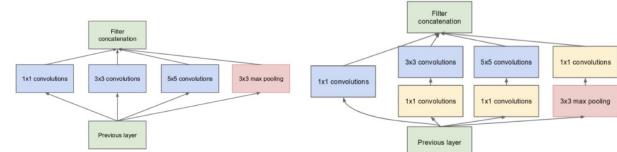
GOOGLENET - STEM BLOCK



Inception Block Diagram:

- Parallel branches process input differently (1×1 convolutions, pooling).
- Outputs concatenated into a deeper, more powerful representation.

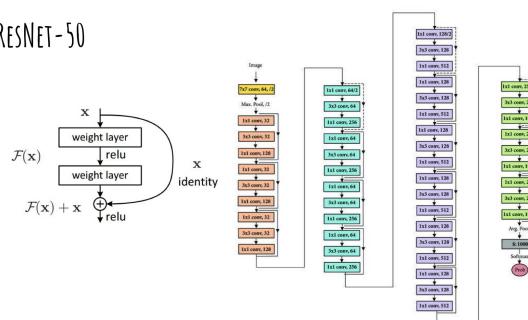
GOOGLENET - INCEPTION BLOCK



5. ResNet-50

- Introduced the concept of Residual Connections (skip connections).
- Addresses the vanishing gradient problem, enabling extremely deep networks (50+ layers).
- Skip connections allow layers to learn residual functions, significantly improving accuracy.

RESNET-50



6. Inception-ResNet

- Hybrid architecture combining GoogleNet's inception modules with ResNet's skip connections.
- Exploits both width (parallel convolutions) and depth (residual connections).
- Enhanced feature extraction capabilities.

INCEPTION - RESNET

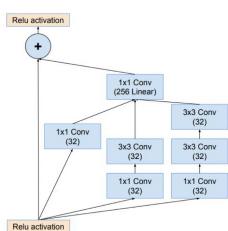


Figure 10. The schema for 35×35 grid (Inception-ResNet-A) module of Inception-ResNet-v1 network.

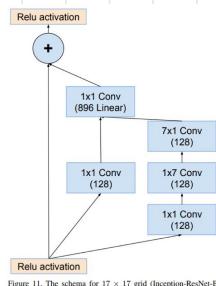


Figure 11. The schema for 17×17 grid (Inception-ResNet-B) module.

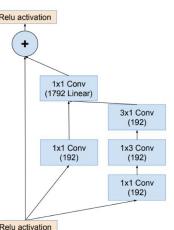


Figure 13. The schema for 8×8 grid (Inception-ResNet-C) module of Inception-ResNet-v1 network.

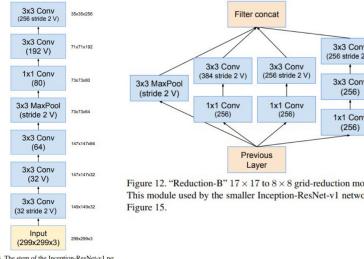


Figure 12. "Reduction-B" $17 \times 17 \times 8$ grid-reduction module. This module used by the smaller Inception-ResNet-v1 network in Figure 15.

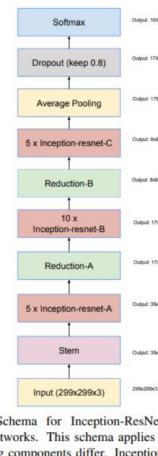


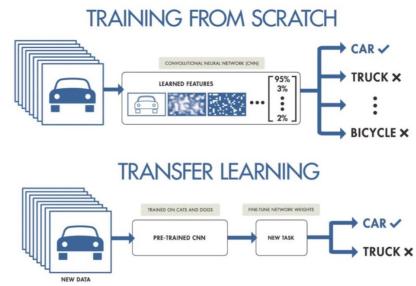
Figure 15. Scheme for Inception-ResNet-v1 and Inception-ResNet-v2 networks. This scheme applies to both networks but the underlying components differ. Inception-ResNet-v1 uses the blocks as described in Figures 14, 10, 7, 11, 12 and 13. Inception-ResNet-v2 uses the blocks as described in Figures 3, 16, 7, 17, 18 and 19. The output sizes in the diagram refer to the activation vector tensor shapes of Inception-ResNet-v1.

Transfer Learning

Applying knowledge learned from solving one task to a different but related task.

Examples:

Using knowledge of riding a bicycle to ride a motorcycle.



Approaches in Transfer Learning

Feature Extraction:

Freeze convolutional layers of a pretrained model.
Retrain only the fully connected layers on new dataset.
Useful when datasets are similar.

Fine Tuning:

Unfreeze the last convolutional layers and fully connected layer + output layer.
Retrain these layers alongside fully connected layers.
Useful when datasets or labels are significantly different.

For deep learning models, the data needs to be labelled, and if the dataset is not publicly available, then we have to manually label the data which is time consuming.

Deep learning models takes a lot of time to get trained.

Pre-trained models: VGG16, ResNet, MobileNet, EfficientNet

Object Detection

Detecting and locating objects within images or videos.

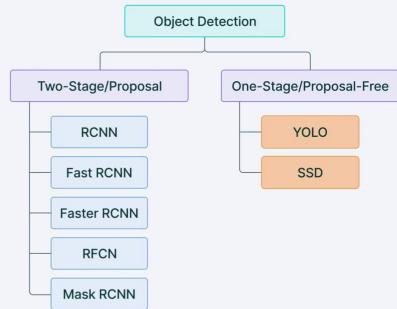
Applications:

Surveillance, self-driving cars, robotics

Types of Object Detection

1. Single shot
2. Two Shot Detection

One and two stage detectors



1. Single-Shot Detection:

- Processes images in one pass.
- Faster and efficient, suitable for real-time applications.
- Example: YOLO (You Only Look Once).

However, single-shot object detection is generally less accurate than other methods, and it's less effective in detecting small objects. Such algorithms can be used to detect objects in real time in resource-constrained environments.

2. Two-Shot Detection:

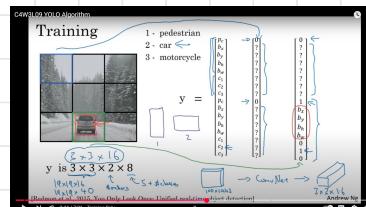
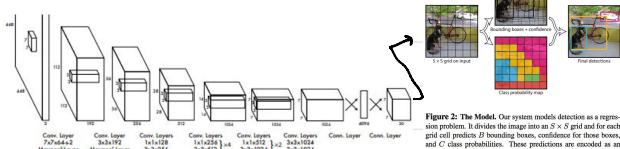
- Two-step process:
 - ↳ First pass: Generate candidate regions (proposals).
 - ↳ Second pass: Classify and refine these proposals.
- More accurate but computationally heavier.
- Example: R-CNN, Faster R-CNN.

→ Potential object locations

YOLO Diagram Explained:

Input image passed once through the network.
Outputs bounding boxes with class probabilities.

YOLO



IoU (Intersection over Union):

- Metric measuring overlap between predicted bounding box and ground-truth.
- Ranges between
 - ↳ 0 (no overlap) and
 - ↳ 1 (perfect overlap).
- Higher IoU = better detection accuracy.

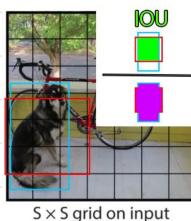
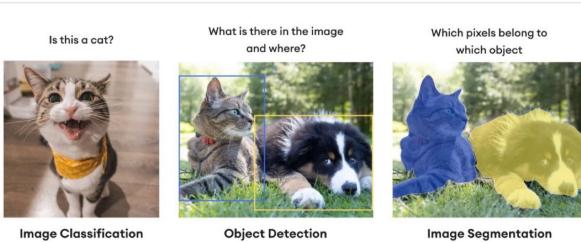


Image Segmentation

- Assigning labels to every pixel of an image.
- Pixels with the same label share certain characteristics.



Types of image segmentation

Semantic Segmentation:

Labels each pixel according to its class
(all birds belong to the same class).

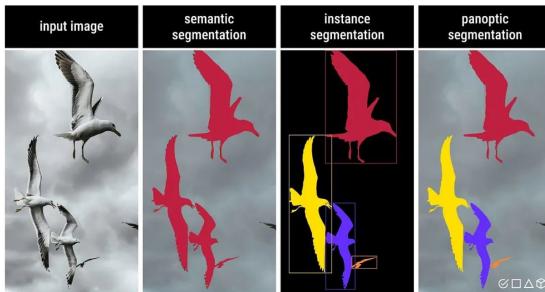
Instance Segmentation:

Assigns a unique label to each instance of an object
(each bird belongs to a different class).

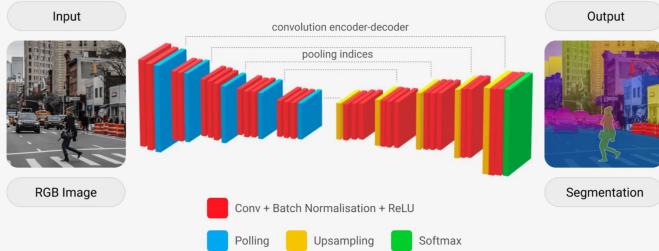
Panoptic Segmentation:

Combines semantic and instance segmentation.
Each instance gets a unique label, grouped under semantic categories.
(each bird has its own class, but they're all identified as a bird)

IMAGE SEGMENTATION



A CNN ARCHITECTURE FOR IMAGE SEGMENTATION.

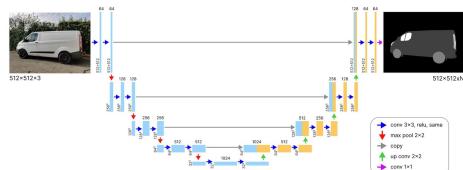


Key Segmentation Architectures

U-Net:

- Primarily used in biomedical image segmentation.
- Encoder-decoder structure.
- Symmetrical, with skip connections copying feature maps directly from encoder to decoder.

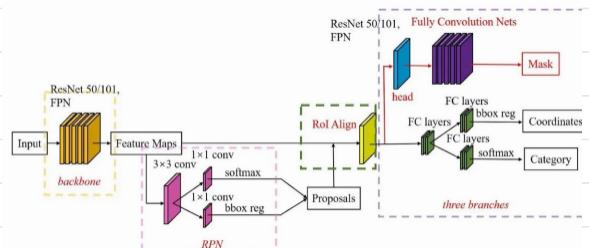
U-Net



MASK R-CNN

Mask R-CNN:

- Extends Faster R-CNN for instance segmentation.
- Adds a parallel branch predicting segmentation masks for each detected instance.
- Highly accurate but computationally demanding.



RNNs

Recurrent Neural Networks (RNNs)

Recurrent Neural Networks are designed specifically for handling sequential data, such as text, speech, or time series data.

Unlike standard neural networks, RNNs have loops that allow information to persist across sequential inputs.

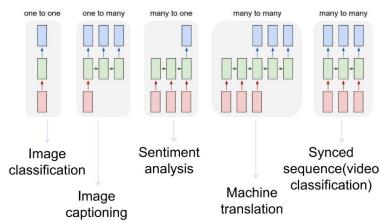
Common applications of RNNs include:

- Generating Text: Predict next words or characters in a sentence.
- Machine Translation: Translating text from one language to another.
- Speech Recognition: Transcribing spoken language into text.
- Generating Image Descriptions: Describing images with meaningful captions.
- Chatbots: Conversational AI capable of generating relevant responses.

Examples of sequence data

Speech recognition		"The quick brown fox jumped over the lazy dog."
Music generation		♪
Sentiment classification		"There is nothing to like in this movie."
DNA sequence analysis		AGCCCTGTGAGGAACTAG
Machine translation		Voulez-vous chanter avec moi?
Video activity recognition		Do you want to sing with me? Running
Name entity recognition		Yesterday, Harry Potter met Hermione Granger.

Motivation

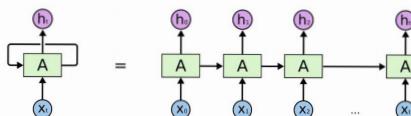


Core Idea of RNNs

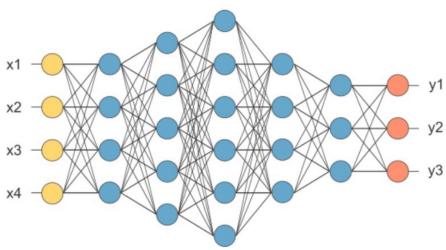
- Standard feed-forward neural networks process inputs independently.
- RNNs maintain hidden states that carry forward information from previous inputs, allowing them to remember context.

Key Components:

- Input Layer: Receives input at each time step.
- Hidden Layer(s): Maintain states that capture temporal dependencies.
- Output Layer: Produces output at each step or at the end of the sequence.

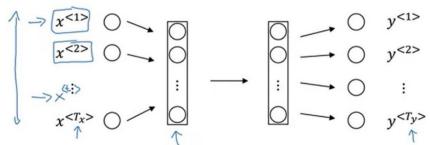


Fully Connected Network



This is our fully connected network. If x_1, \dots, x_n , n is very large and growing, this network would become too large. We now will input **one x_i at a time**, and **re-use the same edge weights**.

Why not a standard network?



Problems:

- - Inputs, outputs can be different lengths in different examples.
- - Doesn't share features learned across different positions of text.

Feedforward Neural Network (FNN):

Processes inputs independently without any memory.

Each input flows one way: input \rightarrow hidden layer(s) \rightarrow output.

No awareness of sequence or context.

Example use: Image classification.

Recurrent Neural Network (RNN):

Designed for sequential data.

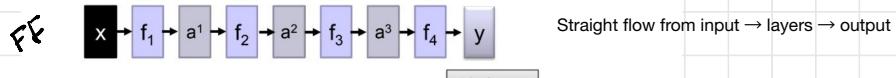
Has loops—each output depends not just on the current input but also on past hidden states.

Retains memory across time steps via hidden states (h_t depends on h_{t-1}).

Example use: Language modeling, time series prediction.

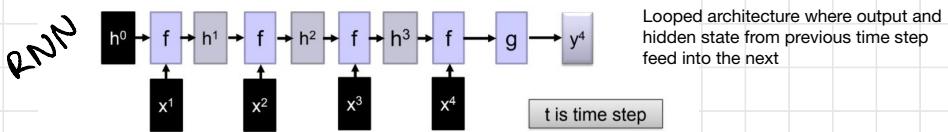
Feed-forward vs Recurrent Network

1. Feedforward network does not have input at each step
2. Feedforward network has different parameters for each layer



$$a^l = f_l(a^{l-1}) = \sigma(W^l a^{l-1} + b^l)$$

Straight flow from input \rightarrow layers \rightarrow output



$$h^l = f(h^{l-1}, x^l) = \sigma(W^l h^{l-1} + W x^l + b^l)$$

Looped architecture where output and hidden state from previous time step feed into the next

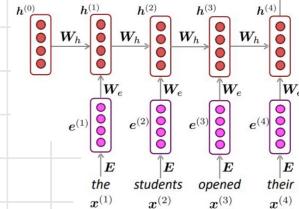
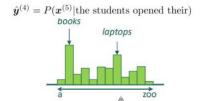
Andrew

RMN Advantages:

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Model size doesn't increase for longer input context
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back



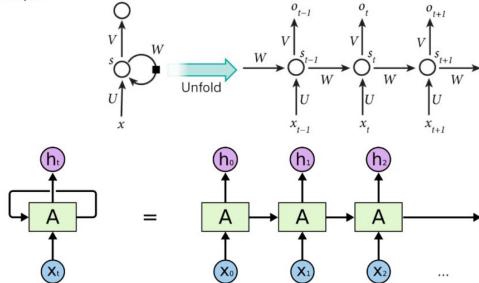
Recurrent Neural Network

Input \rightarrow Hidden \rightarrow Output

it becomes

Input + Previous Hidden \rightarrow Hidden \rightarrow

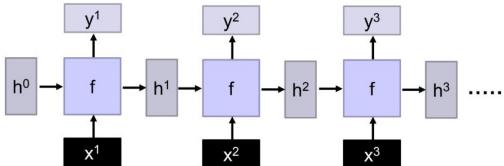
Output



The recurrent loop is formed because hidden state from one step (h_{t-1}) is passed to the next step (h_t).

How does RNN reduce complexity?

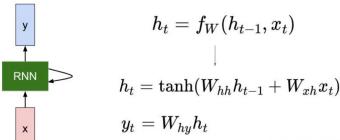
- Given function f : $h', y = f(h, x)$
h and h' are vectors with the same dimension



No matter how long the input/output sequence is, we only need one function f . If f 's are different, then it becomes a feedforward NN. This may be treated as another compression from fully connected network.

Simple RNN

The state consists of a single "hidden" vector h_t :



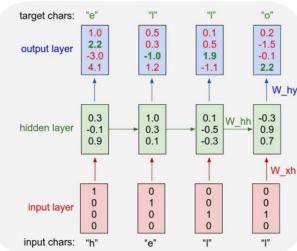
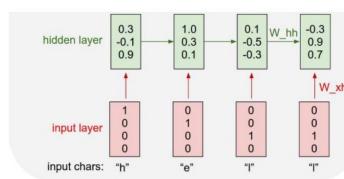
Sometimes called a "Vanilla RNN" or an "Elman RNN" after Prof. Jeffrey Elman

Example: Character-level Language Model

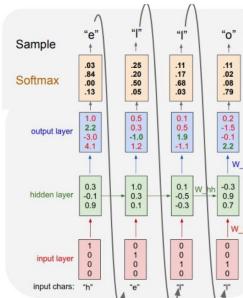
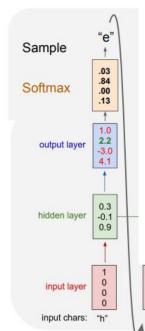
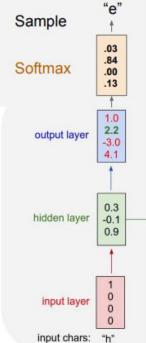
Vocabulary:
[h,e,l,o]

Example training
sequence:
"hello"

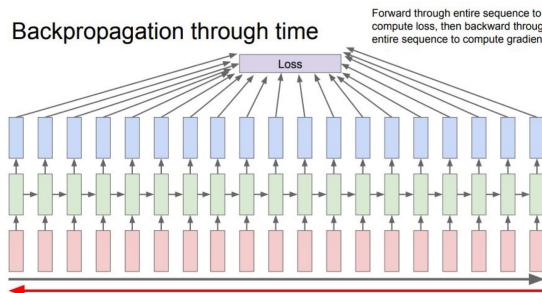
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



At test-time sample
characters one at a time,
feed back to model

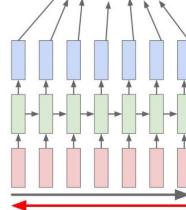


Backpropagation through time

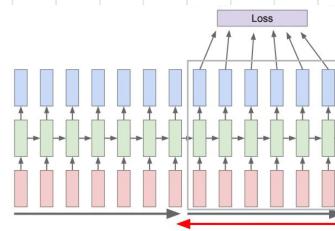


Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient

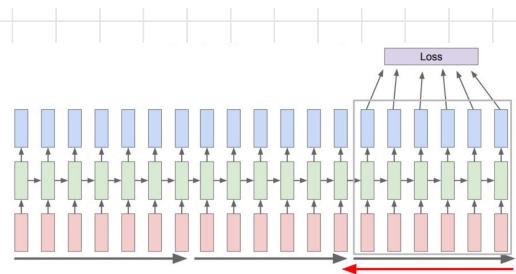
Truncated Backpropagation through time



Run forward and backward through chunks of the sequence instead of whole sequence



Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps



Limitations of Simple RNNs

Vanishing Gradients:

Difficulty in learning long-range dependencies due to gradients becoming very small (close to zero) over many time steps.

Exploding Gradients:

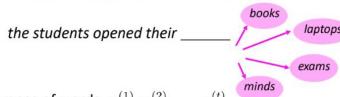
Gradients can become excessively large, causing unstable learning.

SLOW training

very
difficulty in handling long seq,

Language Modeling

- Language Modeling is the task of predicting what word comes next



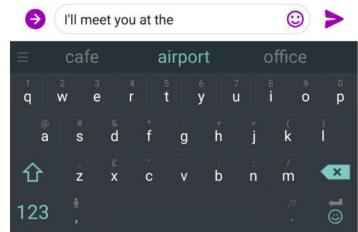
- More formally: given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$$

where $x^{(t+1)}$ can be any word in the vocabulary $V = \{w_1, \dots, w_{|V|}\}$

- A system that does this is called a Language Model

You use Language Models every day!



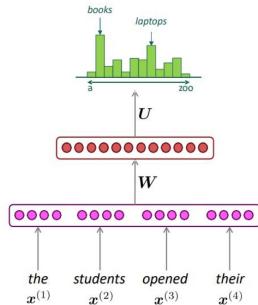
A fixed-window neural Language Model

output distribution
 $\hat{y} = \text{softmax}(U\mathbf{h} + b_2) \in \mathbb{R}^{|V|}$

hidden layer
 $\mathbf{h} = f(W\mathbf{e} + b_1)$

concatenated word embeddings
 $\mathbf{e} = [e^{(1)}, e^{(2)}; e^{(3)}; e^{(4)}]$

words / one-hot vectors
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$



Approximately: Y. Bengio, et al. (2000/2003): A Neural Probabilistic Language Model

Improvements over n-gram LM:

- No sparsity problem
- Don't need to store all observed n-grams

Remaining problems:

- Fixed window is too small
- Enlarging window enlarges W
- Window can never be large enough!
- $x^{(1)}$ and $x^{(4)}$ are multiplied by completely different weights in W .

No symmetry in how the inputs are processed
 a neural architecture
 that can process any length input

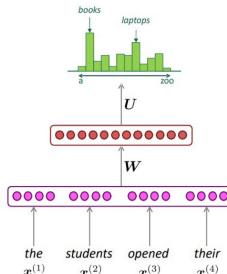
what is the |

what is the weather
 what is the meaning of life
 what is the dark web
 what is the xfl
 what is the doomsday clock
 what is the weather today
 what is the keto diet
 what is the american dream
 what is the speed of light
 what is the bill of rights

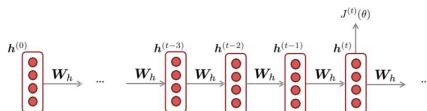
Google Search

I'm Feeling Lucky

Google



Backpropagation for RNNs



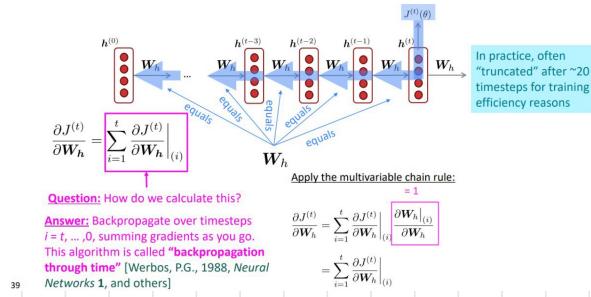
Question: What's the derivative of $J^{(t)}(\theta)$ w.r.t. the **repeated** weight matrix W_h ?

$$\text{Answer: } \frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \frac{\partial J^{(t)}}{\partial W_h} \Big|_{(i)}$$

"The gradient w.r.t. a repeated weight is the sum of the gradient w.r.t. each time it appears"

Why?

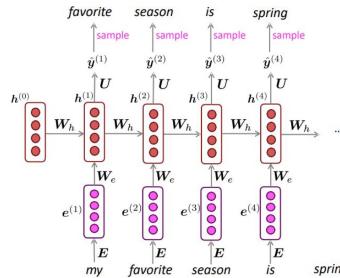
Training the parameters of RNNs: Backpropagation for RNNs



39

Generating text with a RNN Language Model

Just like a n-gram Language Model, you can use a RNN Language Model to generate text by **repeated sampling**. Sampled output becomes next step's input.

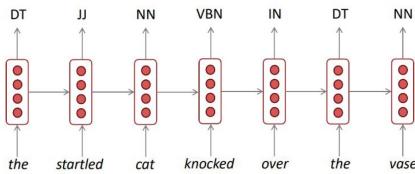


Recap

- **Language Model:** A system that **predicts the next word**
- **Recurrent Neural Network:** A family of neural networks that:
 - Take **sequential input of any length**
 - Apply the **same weights** on each step
 - Can optionally produce output on each step
- **Recurrent Neural Network \neq Language Model**
- We've shown that RNNs are a great way to build a LM.
- But RNNs are useful for much more!

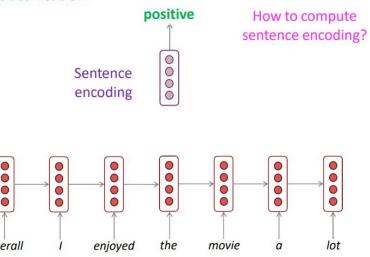
RNNs can be used for tagging

e.g., part-of-speech tagging, named entity recognition



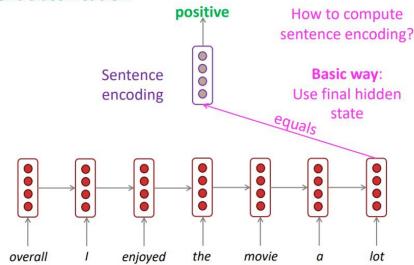
RNNs can be used for sentence classification

e.g., sentiment classification



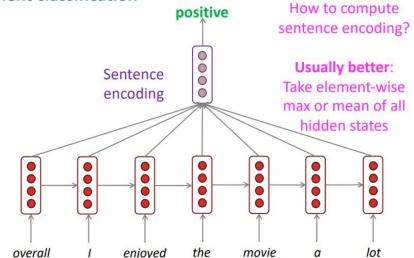
RNNs can be used for sentence classification

e.g., sentiment classification



RNNs can be used for sentence classification

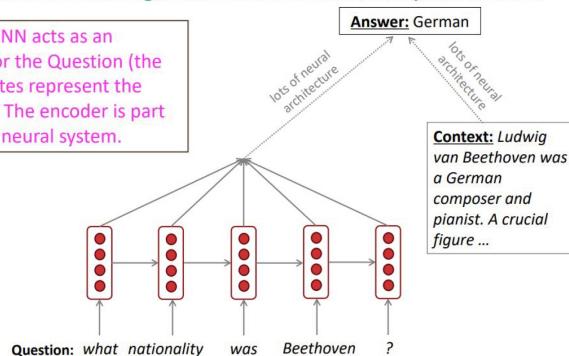
e.g., sentiment classification



RNNs can be used as an encoder module

e.g., question answering, machine translation, many other tasks!

Here the RNN acts as an **encoder** for the Question (the hidden states represent the Question). The encoder is part of a larger neural system.



Bidirectional RNNs

Standard RNNs only consider past information (left-to-right).

Bidirectional RNNs (BiRNNs) use:

A forward RNN (left-to-right)

A backward RNN (right-to-left)

Then concatenate their outputs.

Why?

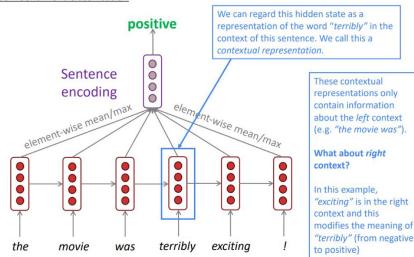
Some tasks require future context as well as past.

Example: In the word "read," context matters for pronunciation

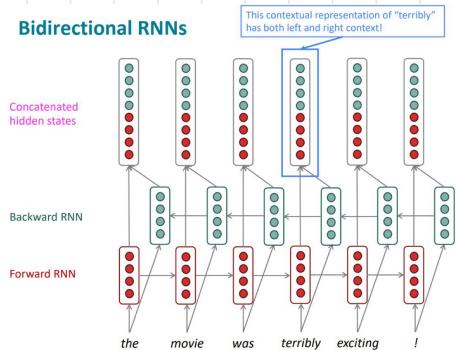
— "I will read" vs. "I have read".

Bidirectional and Multi-layer RNNs: motivation

Task: Sentiment Classification



Bidirectional RNNs



Bidirectional RNNs

On timestep t :

This is a general notation to mean "compute one forward step of the RNN" – it could be a simple RNN or LSTM computation.

$$\text{Forward RNN } \vec{h}^{(t)} = \text{RNN}_{\text{FW}}(\vec{h}^{(t-1)}, \mathbf{x}^{(t)})$$

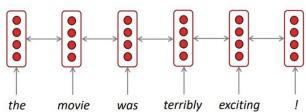
$$\text{Backward RNN } \overleftarrow{h}^{(t)} = \text{RNN}_{\text{BW}}(\overleftarrow{h}^{(t+1)}, \mathbf{x}^{(t)})$$

$$\text{Concatenated hidden states } \mathbf{h}^{(t)} = [\vec{h}^{(t)}, \overleftarrow{h}^{(t)}]$$

Generally, these two RNNs have separate weights

We regard this as "the hidden state" of a bidirectional RNN. This is what we pass on to the next parts of the network.

Bidirectional RNNs: simplified diagram



Bidirectional RNNs

- Note: bidirectional RNNs are only applicable if you have access to the entire input sequence
 - They are **not** applicable to Language Modeling, because in LM you *only* have left context available.
- If you do have entire input sequence (e.g., any kind of encoding), **bidirectionality is powerful** (you should use it by default).
- For example, **BERT** (Bidirectional Encoder Representations from Transformers) is a powerful pretrained contextual representation system **built on bidirectionality**.
 - You will learn more about **transformers**, including BERT, in a couple of weeks!

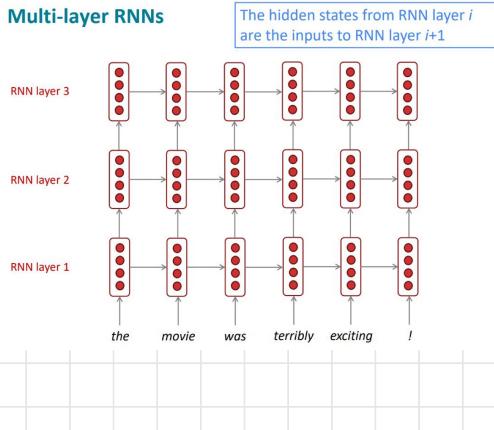
Multi-layer RNNs (a.k.a. Deep RNNs)

A Multi-layer RNN is simply an RNN where the output from one RNN layer becomes input to another RNN layer.

This enables the network to learn more complex representations.

Pros

- Better feature extraction from complex sequential data.
- Used in advanced language models and speech recognition systems.



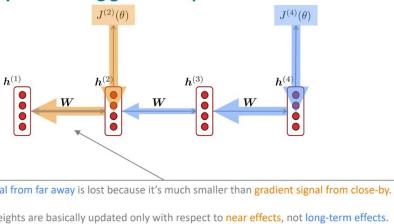
Multi-layer RNNs in practice

- Multi-layer or stacked RNNs allow a network to compute **more complex representations**
 - they work better than just have one layer of high-dimensional encodings!
 - The **lower RNNs** should compute **lower-level features** and the **higher RNNs** should compute **higher-level features**.
- High-performing RNNs are usually multi-layer (but aren't as deep as convolutional or feed-forward networks)
- For example: In a 2017 paper, Britz et al. find that for Neural Machine Translation, **2 to 4 layers** is best for the encoder RNN, and **4 layers** is best for the decoder RNN
 - Often 2 layers is a lot better than 1, and 3 might be a little better than 2
 - Usually, **skip-connections/dense-connections** are needed to train deeper RNNs (e.g., **8 layers**)
- Transformer-based networks (e.g., BERT) are usually deeper, like **12 or 24 layers**.
 - You will learn about Transformers later; they have a lot of skipping-like connections

LSTM

1

Why is vanishing gradient a problem?



How to fix the vanishing gradient problem?

- The main problem is that *it's too difficult for the RNN to learn to preserve information over many timesteps.*
- In a vanilla RNN, the hidden state is constantly being rewritten

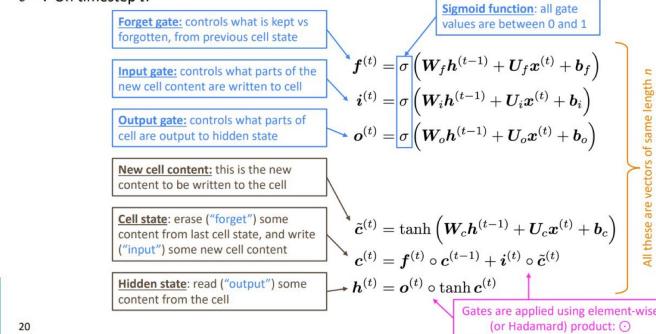
$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_x x^{(t)} + b)$$
- Could we design an RNN with separate **memory** which is added to?

Long Short-Term Memory RNNs (LSTMs)

- On step t , there is a **hidden state** $h^{(t)}$ and a **cell state** $c^{(t)}$
 - Both are vectors length n
 - The cell stores **long-term information**
 - The LSTM can **read**, **erase**, and **write** information from the cell
 - The cell becomes conceptually rather like RAM in a computer
- The selection of which information is erased/written/read is controlled by three corresponding **gates**
 - The gates are also vectors of length n
 - On each timestep, each element of the gates can be **open** (1), **closed** (0), or somewhere in-between
 - The gates are **dynamic**: their value is computed based on the current context

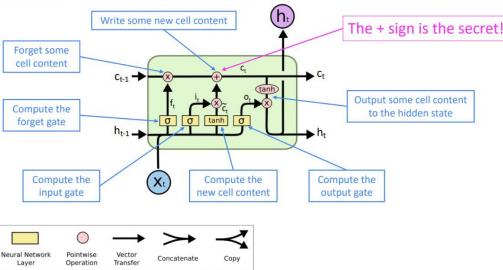
Long Short-Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t :



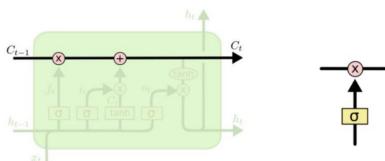
Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:

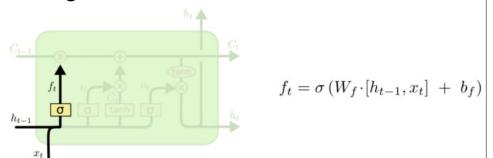


Conveyor Belt

The cell state is kind of like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions. It's very easy for information to just flow along it unchanged

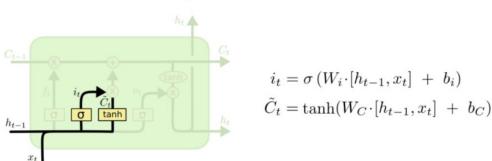


Forget Gate



The sigmoid layer outputs numbers between zero and one, describing how much of each component should be let through. A value of zero means "let nothing through," while a value of one means "let everything through!"

Input Gate Layer

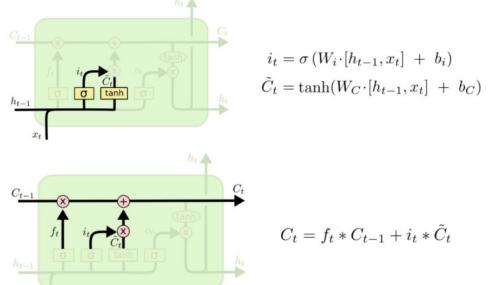


The next step is to decide what new information we're going to store in the cell state

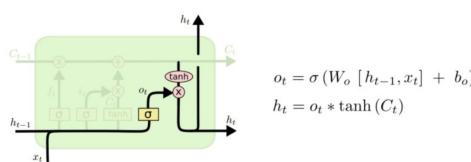
This has two parts. First, a sigmoid layer called the "input gate layer" decides which values we'll update

Next, a tanh layer creates a vector of new candidate values, , that could be added to the state. In the next step, we'll combine these two to create an update to the state

Update Gate Layer + Memory Cell



Output Gate



$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

How does LSTM solve vanishing gradients?

- The LSTM architecture makes it **much easier** for an RNN to **preserve information over many timesteps**
 - e.g., if the forget gate is set to 1 for a cell dimension and the input gate set to 0, then the information of that cell is preserved indefinitely.
 - In contrast, it's harder for a vanilla RNN to learn a recurrent weight matrix W_h that preserves info in the hidden state
 - In practice, you get about 100 timesteps rather than about 7

Is vanishing/exploding gradient just an RNN problem?

- No! It can be a problem for all neural architectures (including **feed-forward** and **convolutional**), especially **very deep ones**.
 - Due to chain rule / choice of nonlinearity function, gradient can become vanishingly small as it backpropagates
 - Thus, lower layers are learned very slowly (i.e., are hard to train)
- Another solution: lots of new deep feedforward/convolutional architectures **add more direct connections** (thus allowing the gradient to flow)

For example:

- Residual connections aka "ResNet"
- Also known as **skip-connections**
- The **identity connection** **preserves information** by default
- This makes **deep** networks much **easier to train**

ReLU

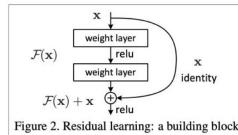


Figure 2. Residual learning: a building block.

LSTMs: real-world success

- In **2013–2015**, LSTMs started achieving state-of-the-art results
 - Successful tasks include handwriting recognition, speech recognition, machine translation, parsing, and image captioning, as well as language models
 - **LSTMs** became the **dominant approach** for most NLP tasks
- Now (**2019–2023**), **Transformers** have become dominant for all tasks
 - For example, in **WMT** (a Machine Translation conference + competition):
 - In WMT 2014, there were 0 neural machine translation systems (!)
 - In **WMT 2016**, the summary report contains "RNN" 44 times (and these systems won)
 - In WMT 2019: "RNN" 7 times, "Transformer" 105 times

MACHINE TRANSLATION

10

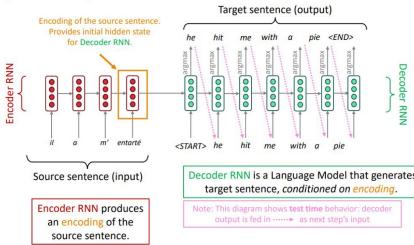
What is Neural Machine Translation?

- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single end-to-end neural network*
- The neural network architecture is called a **sequence-to-sequence** model (aka seq2seq) and it involves **two RNNs**

Machine Translation (MT) is the process of automatically converting text from one language to another using neural networks.

Neural Machine Translation (NMT)

The sequence-to-sequence model



Sequence-to-sequence is versatile!

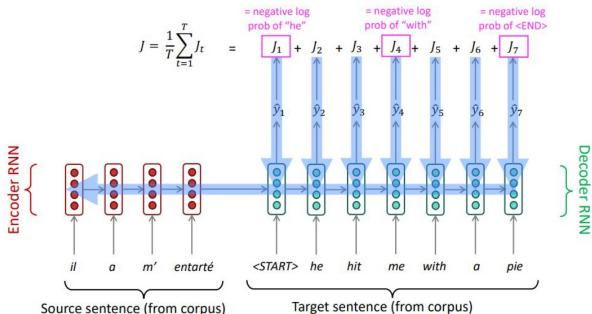
- The general notion here is an **encoder-decoder** model
 - One neural network takes input and produces a neural representation
 - Another network produces output based on that neural representation
 - If the input and output are sequences, we call it a seq2seq model
- Sequence-to-sequence is useful for *more than just MT*
 - Many NLP tasks can be phrased as sequence-to-sequence:
 - **Summarization** (long text → short text)
 - **Dialogue** (previous utterances → next utterance)
 - **Parsing** (input text → output parse as sequence)
 - **Code generation** (natural language → Python code)

Neural Machine Translation (NMT)

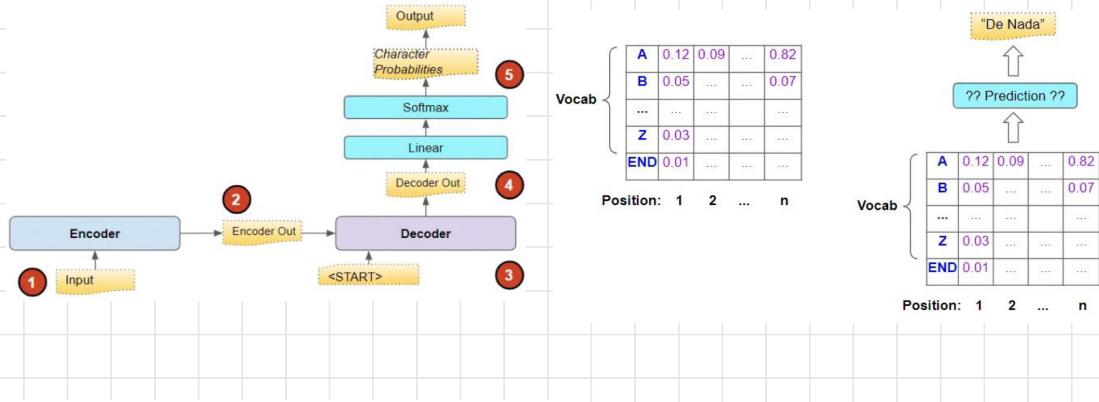
- The **sequence-to-sequence** model is an example of a **Conditional Language Model**
 - **Language Model** because the decoder is predicting the next word of the target sentence y
 - **Conditional** because its predictions are *also* conditioned on the source sentence x
- NMT directly calculates $P(y|x)$:
$$P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots P(y_T|y_1, \dots, y_{T-1}, x)$$

Probability of next target word, given target words so far and source sentence x
- **Question:** How to train an NMT system?
- **(Easy) Answer:** Get a big parallel corpus...
 - But there is now exciting work on “unsupervised NMT”, data augmentation, etc.

Training a Neural Machine Translation system

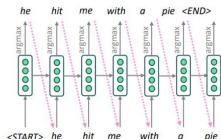


Seq2seq is optimized as a single system. Backpropagation operates "end-to-end".



Greedy decoding

- We saw how to generate (or "decode") the target sentence by taking argmax on each step of the decoder

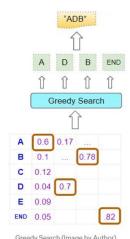


- This is **greedy decoding** (take most probable word on each step)
- Problems with this method?**

Problems with greedy decoding

- Greedy decoding has no way to undo decisions!
 - Input: il a m' entarté (the hit me with a pie)
 - he _____
 - he hit _____
 - he hit a _____

(whoops! no going back now...)



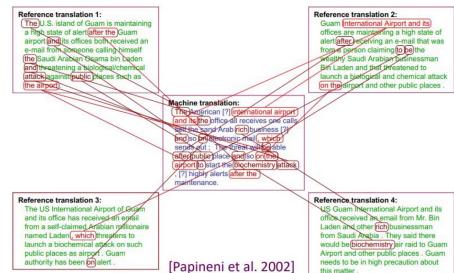
Greedy Search (Image by Author)

How do we evaluate Machine Translation?

BLEU (Bilingual Evaluation Understudy)

- BLEU compares the machine-written translation to one or several human-written translation(s), and computes a similarity score based on:
 - n-gram precision (usually for 1, 2, 3 and 4-grams)
 - Plus a penalty for too-short system translations
- BLEU is **useful** but **imperfect**
 - There are many valid ways to translate a sentence
 - So a **good** translation can get a **poor** BLEU score because it has low n-gram overlap with the human translation ☺

BLEU score against 4 reference translations



BLEU Score (Bilingual Evaluation Understudy)

A widely used evaluation metric in machine translation for comparing a predicted sentence to a reference sentence.

How BLEU Score Works:

Compares n-grams (e.g., unigrams, bigrams) between the predicted and reference sentences.
Uses clipped precision to prevent overcounting repeated words.

Reference Text: It was raining heavily today

Predicted Text: It It It is raining heavily

candidate

Unigram Precision:

Total predicted unigrams: 6
Matching unigrams from reference: 3
(only one "It", "raining", and "heavily" count)
Clipped precision = $3 / 6 = 0.5$

→ 6 predicted words

Bigram Precision:

Reference bigrams:
["It was", "was raining", "raining heavily", "heavily today"]
Predicted bigrams:
["It It", "It It", "It is", "is raining", "raining heavily"]
Matching bigrams: only "raining heavily"
Clipped precision = $1 / 5 = 0.2$

Final BLEU Score Calculation:

BLEU=Brevity Penalty×Precision Score

Where:

r = reference length = 5

c = candidate length = 6

Brevity Penalty (BP) = 1 (since candidate is longer)

Combined score (from example) ≈ 0.316



Usually we use N=4 and $w_n = 1/4$, but in this case we will use N=2 and $w_n = 1/2$

$$\text{Brevity Penalty} = \begin{cases} 1, & \text{if } c > r \\ e^{(1-r/c)}, & \text{if } c \leq r \end{cases}$$

r = number of words in the reference text

c = number of words in the predicted text

Brevity penalty cannot be larger than 1, even if the predicted text is larger than the reference text.

r = 5 c = 6

Bleu = 0.316

NMT: perhaps the biggest success story of NLP Deep Learning?

Neural Machine Translation went from a **fringe research attempt** in **2014** to the **leading standard method** in **2016**

- **2014:** First seq2seq paper published
- **2016:** Google Translate switches from SMT to NMT – and by 2018 everyone has



- This is amazing!

- SMT systems, built by **hundreds** of engineers over **many years**, outperformed by NMT systems trained by a **small group** of engineers in a **few months**

ATTENTION MECHANISM

Why?

RNNs process sequences in order and encode the input into a single fixed-length context vector, which is limiting for long sentences.

Solution:

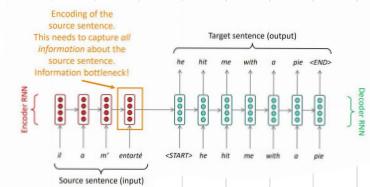
Attention allows the decoder to focus on different parts of the input sequence at each time step.

It helps the model decide "what to look at" while translating each word.

Attention

- Attention provides a solution to the bottleneck problem.
- **Core idea:** on each step of the decoder, **use direct connection to the encoder to focus on a particular part** of the source sequence

ATTENTION



Attention is a *general* Deep Learning technique

- **More general definition of attention:**

- Given a set of vector **values**, and a vector **query**, **attention** is a technique to compute a weighted sum of the values, dependent on the query.

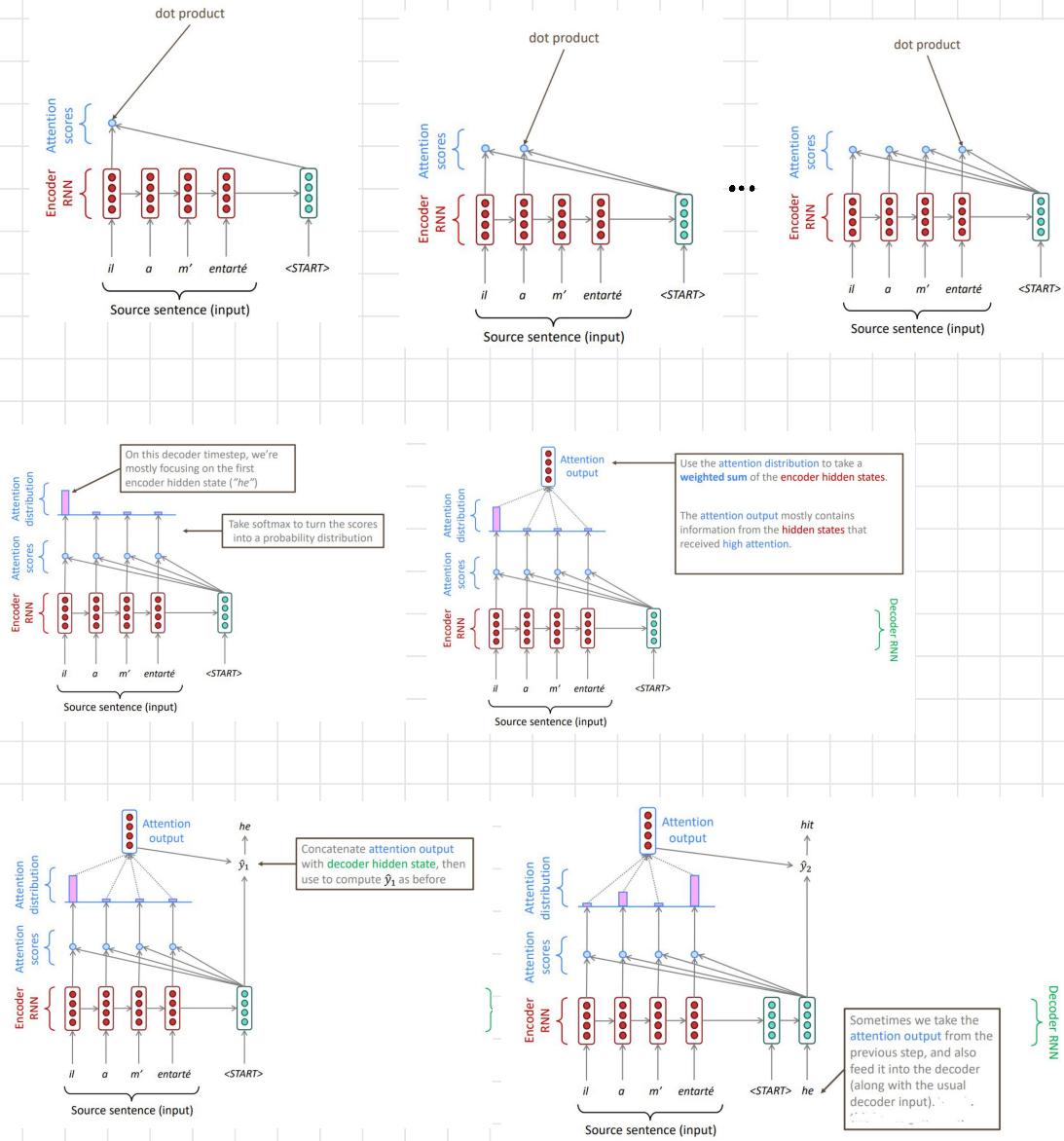
Intuition:

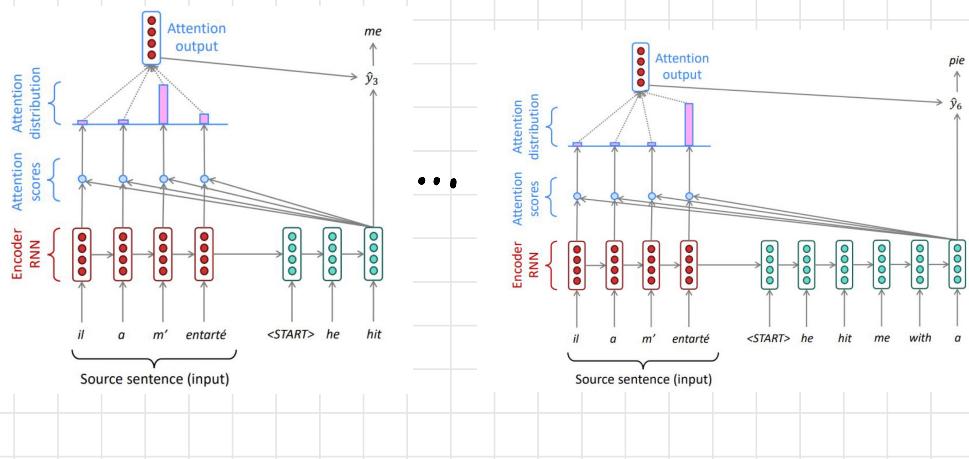
- The weighted sum is a **selective summary** of the information contained in the values, where the query determines which values to focus on.
- Attention is a way to obtain a **fixed-size representation of an arbitrary set of representations** (the values), dependent on some other representation (the query).

Upshot:

- Attention has become the powerful, flexible, general way pointer and memory manipulation in all deep learning models. A new idea from after 2010! From NMT!

Sequence-to-sequence with attention





Decoder RNN

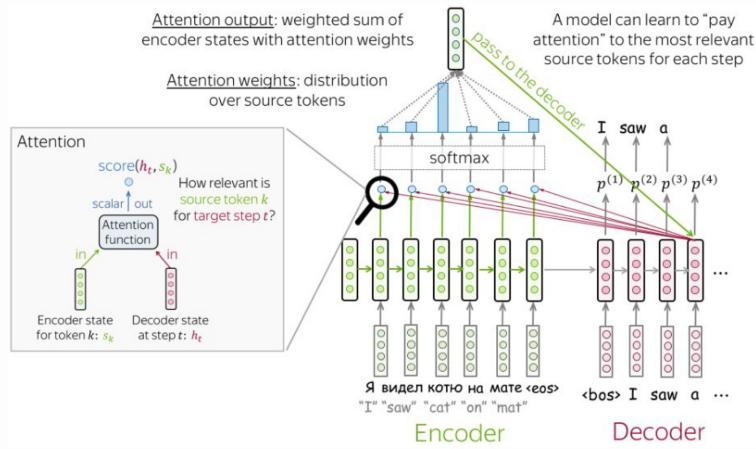
Attention is great!



- Attention significantly **improves NMT performance**
 - It's very useful to allow decoder to focus on certain parts of the source
- Attention provides **more "human-like" model** of the MT process
 - You can look back at the source sentence while translating, rather than needing to remember it all
- Attention **solves the bottleneck problem**
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention **helps with the vanishing gradient problem**
 - Provides shortcut to faraway states
- Attention provides **some interpretability**
 - By inspecting attention distribution, we see what the decoder was focusing on
 - We get (soft) **alignment for free!**
 - This is cool because we never explicitly trained an alignment system
 - The network just learned alignment by itself

il	hi	hi	me	pi	pi
a	hi	me	pi	pi	pi
m'	me	pi	pi	pi	pi
entarté	pi	pi	pi	pi	pi

19



Attention output
↑
(weighted sum)

$$c(t) = a_1^{(t)} s_1 + a_2^{(t)} s_2 + \dots + a_m^{(t)} s_m = \sum_{k=1}^m a_k^{(t)} s_k$$

"source context for decoder step t "

Attention weights
↑
(softmax)

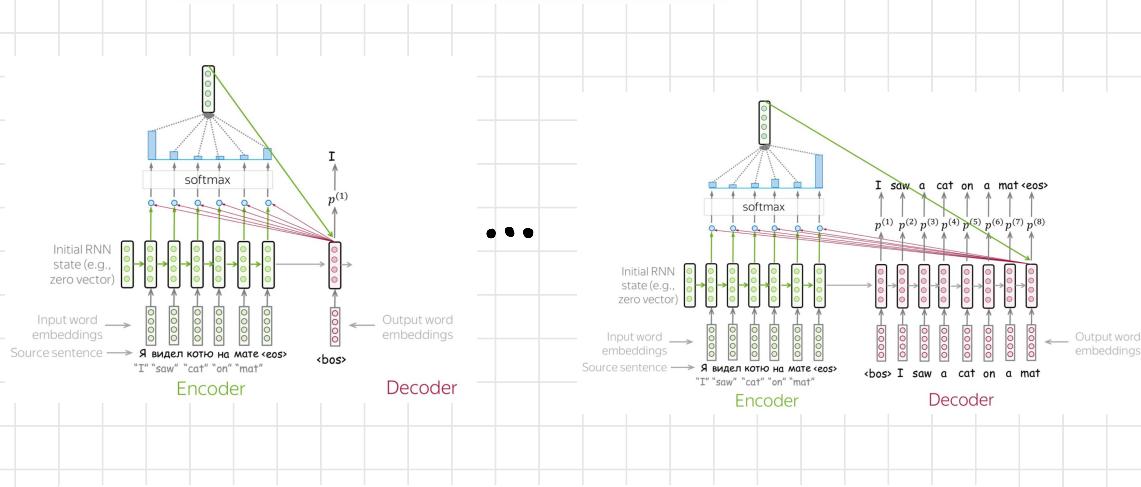
$$a_k^{(t)} = \frac{\exp(score(h_t, s_k))}{\sum_{i=1}^m \exp(score(h_t, s_i))}, k = 1..m$$

"attention weight for source token k at decoder step t "

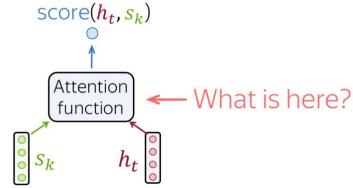
Attention scores
↑
(score(h_t, s_k), $k = 1..m$)

"How relevant is source token k for target step t ?"

Attention input s_1, s_2, \dots, s_m h_t
all encoder states one decoder state



HOW TO COMPUTE ATTENTION SCORE?



1. Dot-product

$$h_t^T \times [s_k]$$

$$\text{score}(h_t, s_k) = h_t^T s_k$$

2. Bilinear

$$h_t^T \times [W] \times [s_k]$$

$$\text{score}(h_t, s_k) = h_t^T W s_k$$

3. Multi-Layer Perceptron

$$w_2^T \times \tanh \left[W_1 \times [h_t] \times [s_k] \right]$$

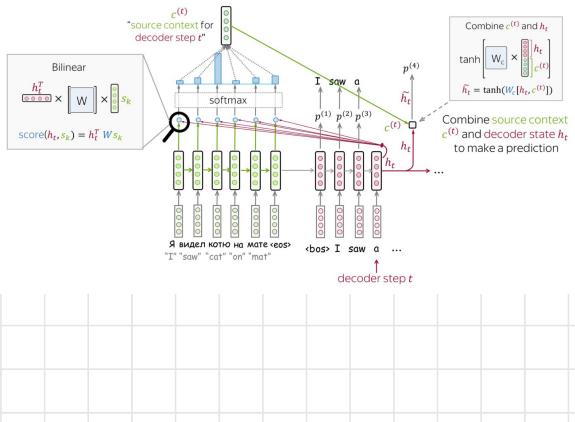
$$\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1[h_t, s_k])$$

Fast and easy to implement.
Used in basic attention models.

bilinear function (aka "Luong attention")

multi-layer perceptron (aka
"Bahdanau attention") - the method
proposed in the original paper

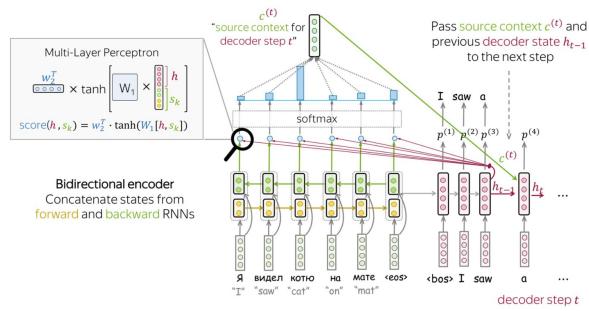
LUONG MODEL



Luong Model:

Also computes attention, but assumes the decoder already has full context and computes attention post hoc.
Uses dot or bilinear scoring.
More efficient, works well with longer sequences.

BAHDANAU MODEL



Bahdanau Model:

Encoder produces hidden states for each input token.
Attention scores are computed for all encoder states.
Context vector is formed and passed to the decoder.
The decoder predicts the next word using this context + its own state.

Key Point: Alignment is learned jointly with translation.

- Attention always involves:

1. Computing the **attention scores** $e \in \mathbb{R}^N$
2. Taking softmax to get **attention distribution** α :

$$\alpha = \text{softmax}(e) \in \mathbb{R}^N$$

There are multiple ways to do this

3. Using attention distribution to take weighted sum of values:

$$a = \sum_{i=1}^N \alpha_i h_i \in \mathbb{R}^{d_1}$$

thus obtaining the **attention output** a (sometimes called the **context vector**)