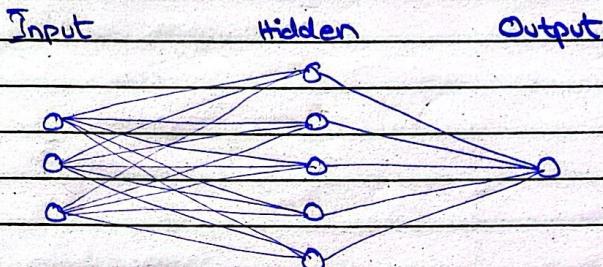


Q.

Artificial Neural Network

An ANN performs many operations in parallel and arranged in layers.



Perceptron :



Activation function

\therefore Each connection from perceptron to other has weight w_j

Computation function

\therefore If $x_0 = 1$

- A neuron with no activation are called Linear Unit.

Layers types :

Convolution	Deconvolution	Fully Connected (FC)	LSTM
Recurrent	Max pooling	Concat	Batch Norm
Dense layers	also called FC.		

Activation Functions :

Input / Hidden
output

ReLU, Leaky ReLU, tanh, sigmoid
softmax, Sigmoid
Multi-class Binary

DATE: ___ / ___ / ___

Advantage of ANN

Disadvantage of ANN

- The ability to learn Non-linear complex relationship.
 - They are good at organizing, preprocessing and categorizing data.
 - Ability to generalize data.
- Lack of rules find optimality based on trial and error.
 - It is also blackbox in nature Dense layers How they work is mystery?
 - every data feed should be numerical.

Example 1 :-

OR

$$\text{Activation Sigmoid} = \frac{1}{1 + e^{-y}}$$

0	0	10
0	1	1

$$\text{OR ReLU } \begin{cases} 1 & \text{if } y > 0 \\ 0 & \text{else} \end{cases}$$

1	0	1
1	1	1

$$y = \gamma_0 + \gamma_1 Q_1 + \gamma_2 Q_2$$

$$y = 1 + 1(0) + 1(0)$$

$$\checkmark \quad y = -1 + 1(0) + 1(0) = -1$$

$$y = -1 + 1(0) + 1(1)$$

$$\checkmark \quad y = -1 + 1(0) + 2(1) = 1$$

$$y = -1 + 1(1) + 2(0)$$

$$\checkmark \quad y = -1 + 2(1) + 2(0) = 1$$

$$\checkmark \quad y = -1 + 2(1) + 2(1) = -1$$

$$x_0 = -1, x_1 = 2, x_2 = 2$$

DATE: ___ / ___ / ___

X

E2.

Perception Rule for XOR Gate

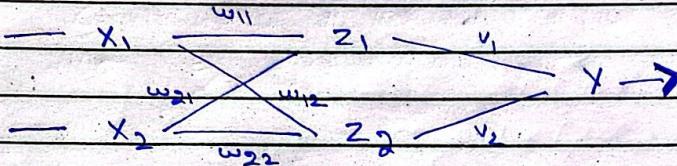
$$y = x_1 \bar{x}_2 + \bar{x}_1 x_2$$

x_1	x_2	y
0	0	0
0	1	1

$$z_1 = x_1 \cdot \bar{x}_2 \quad \therefore \text{ AND } \quad 1 \quad 0 \quad 1$$

$$z_2 = \bar{x}_1 \cdot x_2 \quad 1 \quad 1 \quad 0$$

$$y = z_1 \text{ OR } z_2$$



$$z_1 = x_1 \bar{x}_2 \quad \alpha = 1.5$$

$$x_1 \quad x_2 \quad \bar{x}_2 \quad z_1$$

$$0 \quad 0 \quad 1 \quad 0$$

$$AF = \begin{cases} 1 & \text{if } y \geq 0 \\ 0 & y < 0 \end{cases}$$

$$0 \quad 1 \quad 0 \quad 0$$

$$1 \quad 0 \quad 1 \quad 1$$

$$y = w_1 x_1 + w_2 x_2$$

$$1 \quad 1 \quad 0 \quad 0$$

$$\checkmark \quad y = 1 \cdot 0 + 1 \cdot 0 = 0$$

$$\times \quad y = 1 \cdot 0 + 1 \cdot 1 = 1$$

$$y = 1 \cdot x_1 + (-0.5) x_2$$

$$\left\{ \begin{array}{l} w_{12} = \alpha \theta + \alpha (P\bar{T} - 0) x_2 \\ = 1 + 1.5(0 - 1) 1 \\ = -0.5 \end{array} \right.$$

$$y = 1 \cdot 0 + (-0.5) 1 = -0.5$$

$$\checkmark \quad y = 1 \cdot 1 + (-0.5) 1 = 0.5$$

$$\checkmark \quad y = 1 \cdot 1 + (-0.5) (1) = 0.5$$

DATE: ___ / ___ / ___

$$z_2 = \bar{x}_1 x_2$$

x_1	x_2	\bar{x}_1	z_2
0	0	1	0
0	1	1	1

✓ $y = 1*0 + 1*0 = 0$

$$1 \quad 0 \quad 0 \quad 0$$

✓ $y = 1*0 + 1*1 = 1$

$$1 \quad 1 \quad 0 \quad 0$$

✗ $y = 1*\cancel{0} + 1*0 = 01$

$$w_{12} = \frac{1}{2} + 1.5(0-1)*1 = 0.5$$

$$w_{22} = \frac{1}{2} + 1.5(0-1)*0 = 1$$

✓ $y = \bar{0.5} * 0 + 1 * 0 = -0.5$

✓ $y = (-0.5) * 1 + 1 * 1 = 0.5$

x	y	z_1	z_2	y
0	0	0	0	0
0	1	0	1	1
1	0	1	0	1
1	1	0	0	0

✓ $y_{in} = 1*0 + 1*0 = 0$

✓ $y_{in} = 1*0 + 1*1 = 1$

✓ $y_{in} = 1*1 + 1*0 = 01$

$y_{in} = 1*0 + 1*0 = 0$

$$w_{11} = 1 \quad w_{21} = -0.5$$

$$w_{12} = -0.5 \quad w_{22} = 1$$

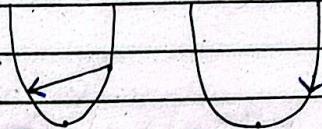
$$v_1 = 1 \quad v_2 = 1$$

1.

(i) Hyperparameters

Learning Rate (α)

This shows how much
longer jump is to be made
to reach min loss.



large learning rate small learning rate.

* If not set properly take huge time to optimize.

Minibatch - convert data to small than pass for each epoch.

Learning depends on data in one epoch. If

Batch size is large more memory is require.

\Rightarrow If batch size is too small with noisy data the convergence will be slow.

(ii) Optimisers

Adam

combine Root mean

SGD

Updated parameter (Q_t)

square and Adaptive

opposite to gradient

gradient algorithm.

$$m_{t+1} = \beta_1 m_t + (1 - \beta_1) \Delta J(Q_t) \quad Q_{t+1} = Q_t - \alpha \Delta J(Q_t)$$

$$v_{t+2} = \beta_2 v_t + (1 - \beta_2) \Delta J(Q_t)$$

$$\hat{m}_{t+1} = \frac{m_{t+1}}{1 - \beta_1^{t+1}} \quad v_{t+1} = \frac{v_{t+1}}{1 - \beta_2^{t+1}} \quad \begin{array}{l} \text{Gradient direction} \\ \text{is opposite to} \\ \text{min loss of} \end{array}$$

$$Q_{t+1} = Q_t - \alpha \sqrt{v_{t+1} + \epsilon} \quad m_{t+1}$$

DATE: ___ / ___ / ___

- * Beta₁, Beta₂ Parameters for Adam.
- and Epsilon.
- * Adam maintain separate learning rate for each parameter and update dynamically based on past gradient.
- * Adam works for wide range of Deep learning tasks.

(iii)

Regularization

Dropout

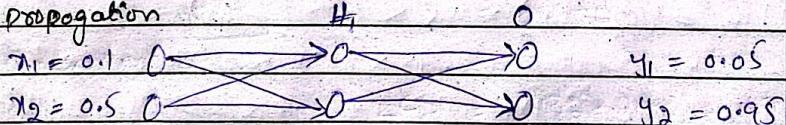
Early stopping

Batch Normalization

- Turn off some connection to learn less \Rightarrow Adjust data to have mean of 0 and std. of 1 \Rightarrow The model keep track of error and which its not improving it stop training.
- learn less to avoid overfitting \Rightarrow In learn help avoid noise which can result in more errors is test.

Example solved for Forward propagation and Backward propagation

Forward



Backward

$$\begin{aligned}
 & \therefore w_2 = 0.2 \quad w_1 = 0.1 \quad w_3 = 0.3 \quad w_4 = 0.4 \quad b_1 = 0.25 \quad b_2 = 0.25 \\
 & H_1 = 0.1 \times 0.1 + 0.5 \times 0.3 + 0.25 = 0.41 \Rightarrow \frac{1}{1+e^{-(0.41)}} = 0.60 \\
 & H_2 = 0.1 \times 0.2 + 0.5 \times 0.4 + 0.25 = 0.47 \Rightarrow \frac{1}{1+e^{-(0.47)}} = 0.61 \\
 & o_1 = 0.60 \times 0.5 + 0.61 \times 0.6 \neq 1.01 \Rightarrow \frac{1}{1+e^{-(1.01)}} = 0.73 \\
 & o_2 = 0.60 \times 0.7 + 0.61 \times 0.8 + 0.25 = 1.26 \Rightarrow \frac{1}{1+e^{-(1.26)}} = 0.77
 \end{aligned}$$

DATE: ___ / ___ / ___

(2)

Error

 y_{org} y_{Pred}

$$E_1 = \frac{1}{2} (0.05 - 0.73492)^2 = 0.23456$$

$$E_2 = \frac{1}{2} (0.95 - 0.77955)^2 = 0.01452$$

$$E_f = 0.24908$$

n=2

(3)

Backpropagation

$$\alpha = 0.6$$

$$\begin{aligned} \delta w_5^{\text{P}} &= (o_1 y_{\text{Pred}} - y_{\text{true}}) * (o_1(1-o_1)) * h_1 \\ &= (0.73 - 0.05)(0.73(1-0.73)) * 0.60 \\ &= 0.08020 \end{aligned}$$

$$\delta w_6^{\text{P}} = (-0.68492)(0.19880)(0.61) = 0.08211$$

$$\delta w_7^{\text{P}} = (0.77955 - 0.95)(0.9(1-0.9)) * 0.60 \\ = -0.01760$$

$$\delta w_8^{\text{P}} = (-0.17044) * (0.17184) * 0.60108 = -0.01802$$

$$w_5' = w_5 - \alpha (\delta w_5)$$

$$= 0.5 - 0.6(0.08020) = 0.45187 \quad \downarrow$$

$$w_6' = 0.6 - 0.6(0.08211) = 0.555073 \quad \downarrow$$

$$w_7' = 0.7 - 0.6(-0.01760) = 0.71056 \quad \uparrow$$

$$w_8' = 0.8 - 0.6(-0.01802) = 0.81081 \quad \uparrow$$

w₅

$$\delta w_1^{\text{P}} = (o_1 - t_1) \times w_5 \times h_1 \times (1-h_1) \times x_1$$

$$= 0.68492 \times 0.198 \times 0.15 \times 0.2397 \times 0.1 = 0.00159$$

w₁

$$\delta w_1 = (-0.17044) * (0.17184) * 0.7 \times 0.23978 \times 0.1 = -0.0049$$

w₇

$$\delta w_1 = 0.00110 \quad (\text{some } \epsilon \text{ and } \epsilon_2)$$

$$= \frac{0.5}{0.61}$$

$$w_1' = w_1 - \alpha (\delta w_1)$$

$$= 0.1 - 0.6(0.00110) = 0.09933 \quad \downarrow$$

$$= 0.510583$$

$$w_1' = 0.19919 \quad \downarrow$$

$$w_3' = 0.29607 \quad \downarrow$$

$$(0.008126 + 0.002457)$$

$$w_4' = 0.39597 \quad \downarrow$$

$$= 0.2956$$

This is for one sample movement forward and updating weights backwards.

DATE: 1/1

Convolution Operations

Normal Convolution

$$\begin{array}{l}
 6 \times 6 \times 3 \quad 3 \times 3 \times 3 \quad (\text{no of filter}) = 4 \times 4 \times 5 \\
 \text{Image} \quad \text{filter} \quad 5 \\
 S=1 \quad p=0 \quad \left[\frac{n+2p-f+1}{p} \right] = 4 \\
 h=1 \quad f=3
 \end{array}$$

Computation Each $3 \times 3 \times 3$ is Multi to each $6 \times 6 \times 3$

$$\begin{array}{l}
 \text{Compute cost} = \# \text{filter para} \times \# \text{filter position} \times \text{No of filters} \\
 2160 = 3 \times 3 \times 3 \times 4 \times 4 \times 5
 \end{array}$$

Depth-Wise Separable Convolution

Depth wise

Point wise

$$\begin{array}{l}
 \bullet \text{ Separate Each channel} \quad 4 \times 4 \times 3 * \underbrace{1 \times 1 \times 3}_{5} = 4 \times 4 \times 5 \\
 \text{apply filter.}
 \end{array}$$

$$6 \times 6 * 3 \times 3 = 4 \times 4$$

$$6 \times 6 * 3 \times 3 = 4 \times 4$$

$$6 \times 6 * 3 \times 3 = 4 \times 4$$

$$CS = 1 \times 1 \times 3 \times 4 \times 4 \times 5$$

stack them

$$\sum CS = 432 + 240 = 672$$

$$CS = 3 \times 3 \times 4 \times 4 \times 3$$

$$= 432$$

But Not unlike Normal they 4×4 are separate has no interlink so we use Pointwise conv after Depthwise.

Convolutionary Neural Networks

Formulas :

1. Convolution layer

$$\text{No of param} = [(f_w \times f_h \times \text{depth}) + 1] \times \text{No of filters}$$

2. Fully Connected (Dense)

$$\text{No of param} = (\text{Input} + 1) \times \text{output size}$$

Note :: for Input layer there is no learning parameters
but the activation is length \times width \times height.

3. Activation Size

$$AS = \text{height} \times \text{width} \times \text{channels}$$

Solution :

for conv1. filter size (5,5) stride (1)

$$\text{so, No of param} = [(5 \times 5 \times 3) + 1] \times 8 = 608$$

3 for No of channels of Input Image

8 for No of filters to be applied to Image.

DATE: / /

For Conv2 filter (5,5) stride (1)

$$so, \text{ No of param} = [(5 \times 5 \times 8) + 1] \times 16 = 3216$$

8 for previous layer depth

16 for filter applied onto them.

For Conv2

~~Activation~~ = $\frac{10 \times 10 \times 16}{size} = 1600$

For Conv1

~~Activation~~ = $\frac{28 \times 28 \times 8}{size} = 6272$

For Pooling 1

~~Activation~~ = $\frac{14 \times 14 \times 8}{size} = 1568$

For FC

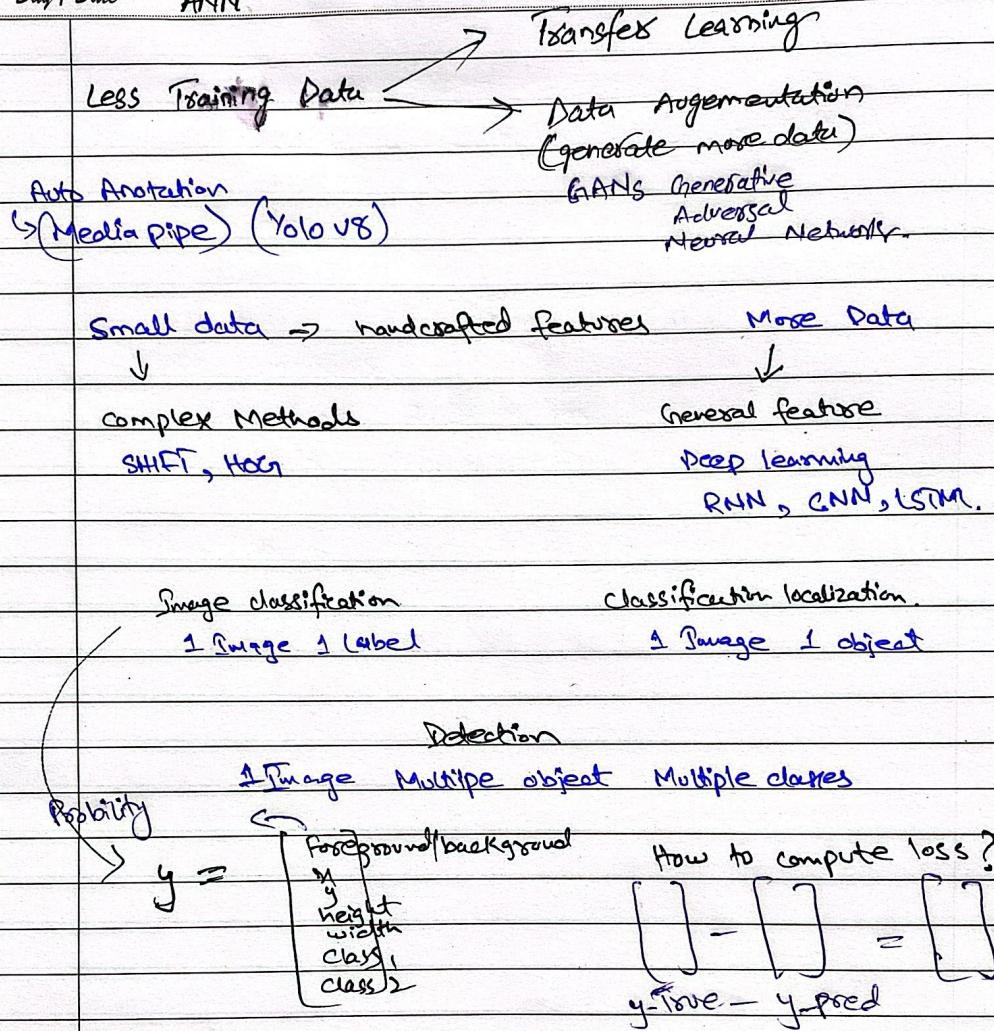
$$\frac{\text{Activation}}{\text{Size}} = 120 \times 1 = 120$$

For ~~FC~~ FC

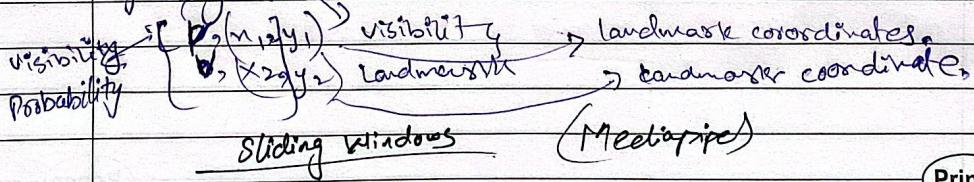
$$\begin{aligned} \text{No of param} &= (\text{Input} + 1) \times \text{output} \\ &= (400 + 1) \times 120 = 48120 \end{aligned}$$

For output (softmax)

$$\begin{aligned} \text{No of param} &= (84 + 1) \times 10 \\ &= 850 \end{aligned}$$



Find all or any one landmark then uses formulas to generate any missing to create 3D mesh.



Slicing Windows

fixed



Random.

Background removal

Any change in background/
Masking will also not work.
Intensity.

To find initial
spot of space
where change
of object

Then introduced
random size to
get exact size
of ytrue and
ypred.

} object are
not of same
size can become small
big by camera view.

Mediapipe

FC to Convolution? So we adopted pointwise convolution.

These convolution output won't go into linear equation
No linearity equation.

Conv \rightarrow $y = m + b$ Conv \rightarrow Conv \rightarrow Predict (pixels)
Conv

\Rightarrow Convolution is better in feature sharing more good.

Yolo Algorithm

\Rightarrow convert each image into grid.

\Rightarrow Make y for each grid. (as before)

\Rightarrow Pass whole image you only look once

\Rightarrow output $(n \times n \times 8) \rightarrow y$ length
grid numbers

Loss computation

Prediction

bound box

$$y_{true} - y_{pred} = loss$$

Intersection over Union

Prince

1.

In You Only Look Once there are several anchor boxes also known defined as prior boxes or default boxes.

V1 2 anchor boxes , No grid formation (single scale)

V2 5 anchor boxes , 13×13 grid (at each scale)

V3 3 anchor boxes (at 3 different scales)

13×13 , 26×26 , 52×52 grids (at each scale)

V4 Same as V3 Mosaic Training

combine multiple grids

V5 same as V4 and improved anchor box initialization.

V6 Introduced " anchor free " approach

uses key point Detection and Non maximum separation

V7 Introduced " Auto Assign Labels "
" dual-teaching "

In dual-teachers model v7 learn

from 2 teachers one is rigid
and other is ease in such manner

knowledge distilled and learning
is consistence .

V8

same
as
YOLO
v3
and
v4

Anchors free

16x16, 32x32, 64x64 grid (at each scale)

Mosaic Training (combine multiple grid)

Improved auto assign label (IoU)

Spatial attention

- (1) Apply Average Pooling and Max Pooling along channel axis and concat them.
- (2) On concat feature descriptor you apply Conv layers. which encodes where to emphasize or suppress.
- (3) Then aggregate channel info into 2D map by using two Pooling operations.

$$MS(F) = \sigma(f_7 \times f_7 ([\text{AvgPool}(F), \text{MaxPool}(F)]))$$

$$MS(F) = \sigma(f_7 \times f_7 ([f_{avg}, f_{max}]))$$

$\sigma \Rightarrow$ Sigmoid function $f_7 \times f_7 \Rightarrow$ CNN operation
with filter 7×7

New use cases :- Instance Segmentation.

Keypoint Detection.

2.

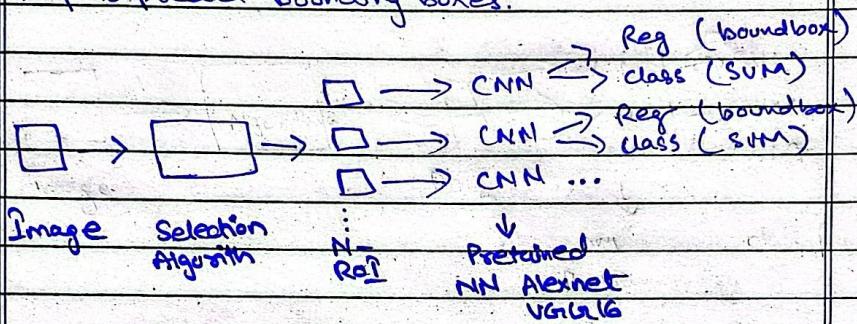
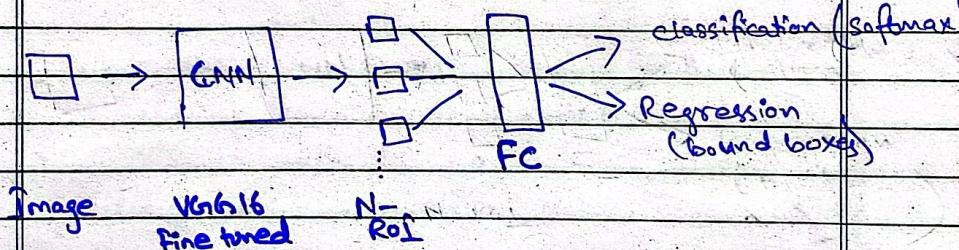
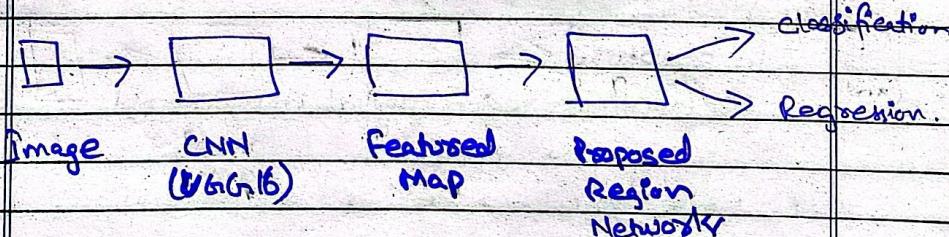
Region Based Conv Neural Network Family

RCNN uses selective search algorithm to propose Region of Interest (RoI) and then uses CNN to classify each RoI.

Fast RCNN uses shared convolution layer to whole image and all ROI rather than dealing with each separately (Reduces computational time)

Faster RCNN uses Region Proposed Networker (RPN) to generate ROI which is faster than selective search algorithm. These networks also uses feature map to predict bounding boxes.

RCNN

RCNN
fasterRCNN
x
selective

Sliding Windows in CNN:

9 Anchors of different scales $1/8, 1/16, 1/32$ are applied to feature map. Then fed to proposed Region Networker.

Non-Maximum Suppression :

on output of RPN it is applied
to remove duplication of detection.

Note:: Bounded Box with max confidence is
chosen.

3.

Key Differences :

1. RCNN uses selective search Algorithm while Yolo uses single neural network to predict bounding boxes.
2. RCNN has two stage approach (Region proposal and classification). Yolo is single stage.
3. RCNN prioritizes accuracy over speed while Yolo prioritized speed over accuracy.
4. Yolo is used on live video stream while RCNN benefits for high accurate decision making on regions.

In structure

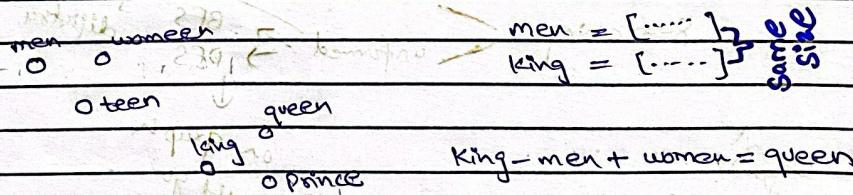
5. RCNN has Region Proposal Network and RoI Pooling extra layer.
While Yolo uses one grid cell generation.

Handle any size of token length as input?
we used chunking (simple trim)

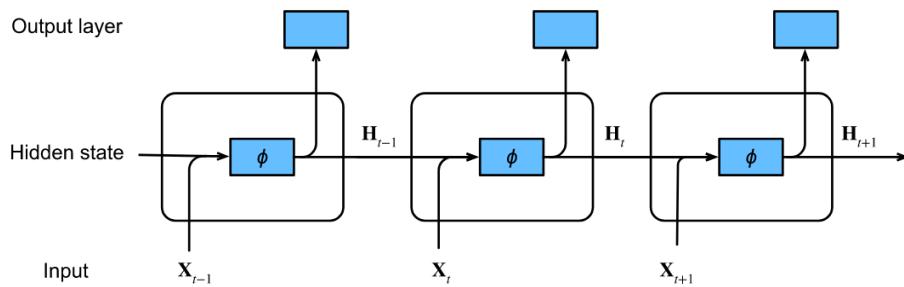
- Attention Mechanisms.
- Retrieval Augmented generation.
- Recurrent Neural Networks.

Embedding layer (`input_dim, output_dim, initializer='uniform', regularization, constraints`)

- ⇒ Categorical Data can not be passed to neural network directly.
- ⇒ Encoding scheme fails when number of words so vector size of each word become very long and it does not capture semantic in sentences.
- ⇒ Embedding layer is a mapping b/w categorical data and its vector representation.
- ⇒ The goal of embedding layer is to arrange embedding in a vector space such that similar words are together.
- ⇒ Not just that if apply operation on these vectors you will get another vector in dimension.



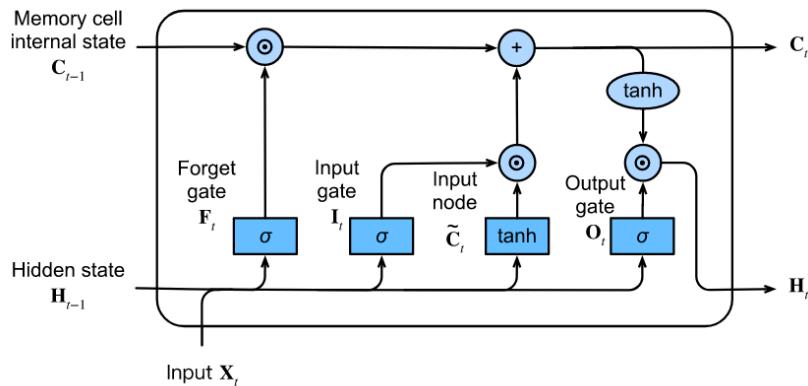
- This is same as game we play in child-hood King is to men and women is to ? queen Relationship
- All LLMs now a days use this Word embedding as base.



ϕ FC layer with activation function \nearrow Copy $\overbrace{->}$ Concatenate

RNN (Recurrent Neural Network)

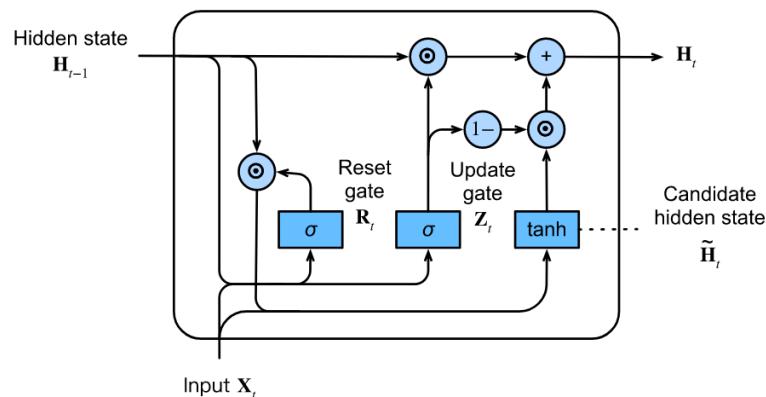
- A type of neural network that processes sequential data
- Uses recurrent connections to capture temporal relationships
- Suffers from vanishing gradient problem, making training difficult



σ FC layer with activation function ○ Elementwise operator \nearrow Copy $\overbrace{->}$ Concatenate

LSTM (Long Short-Term Memory)

- A type of RNN that addresses vanishing gradient problem
- Uses memory cells to store information for long periods
 - Gates (input, output, forget) control information flow
 - More effective for long-term dependencies



σ FC layer with activation function ○ Elementwise operator \nearrow Copy $\overbrace{->}$ Concatenate

GRU (Gated Recurrent Unit)

- A simpler alternative to LSTM
- Uses gates (reset, update) to control information flow
 - Faster training and inference compared to LSTM
 - Less effective for very long sequences

In summary:

- RNN: basic recurrent network, struggles with long-term dependencies.
 - LSTM: more effective for long-term dependencies, but slower.
 - GRU: balances performance and speed, suitable for most tasks.

Backpropogating an LSTM: A Numerical Example

Let us begin by defining our internal weights:

$$W_a = \begin{bmatrix} 0.45 \\ 0.25 \end{bmatrix}, U_a = [0.15], b_a = [0.2]$$

$$W_i = \begin{bmatrix} 0.95 \\ 0.8 \end{bmatrix}, U_i = [0.8], b_i = [0.65]$$

$$W_f = \begin{bmatrix} 0.7 \\ 0.45 \end{bmatrix}, U_f = [0.1], b_f = [0.15]$$

$$W_o = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}, U_o = [0.25], b_o = [0.1]$$

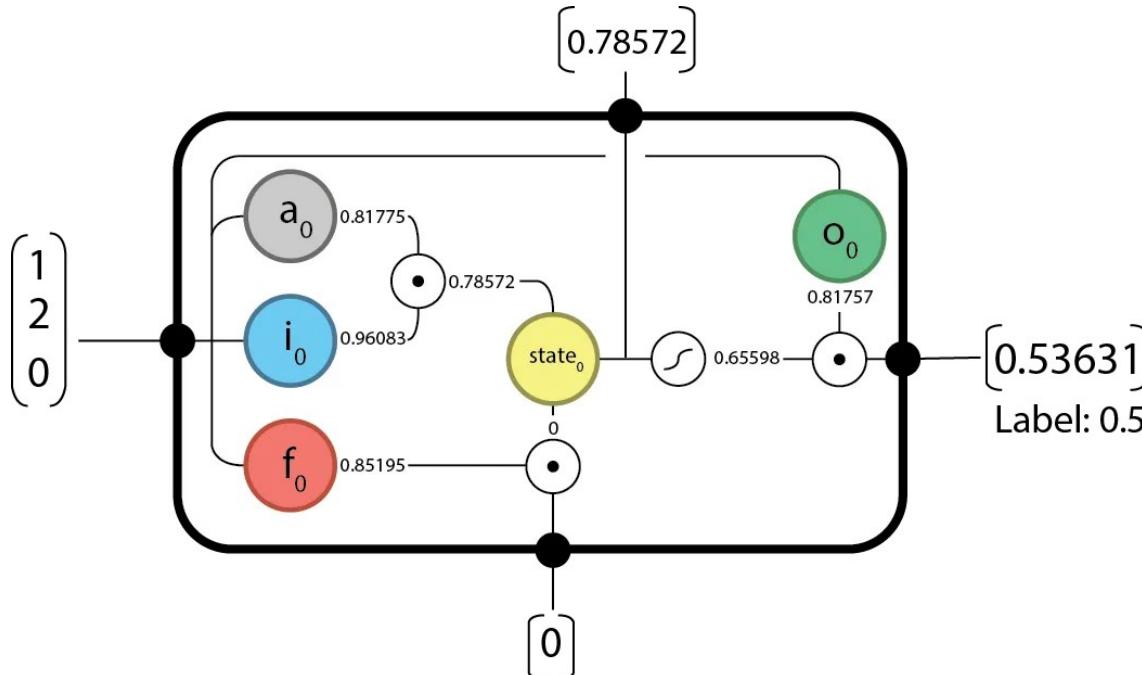
And now input data:

$$x_0 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \text{ with label: } 0.5$$

$$x_1 = \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} \text{ with label: } 1.25$$

I'm using a sequence length of two here to demonstrate the unrolling over time of RNNs.

Forward @ t=0



$$a_0 = \tanh(W_a \cdot x_0 + U_a \cdot out_{-1} + b_a) = \tanh([0.45 \ 0.25] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.15] [0] + [0.2]) = 0.81775$$

$$i_0 = \sigma(W_i \cdot x_0 + U_i \cdot out_{-1} + b_i) = \sigma([0.95 \ 0.8] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.8] [0] + [0.65]) = 0.96083$$

$$f_0 = \sigma(W_f \cdot x_0 + U_f \cdot out_{-1} + b_f) = \sigma([0.7 \ 0.45] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.1] [0] + [0.15]) = 0.85195$$

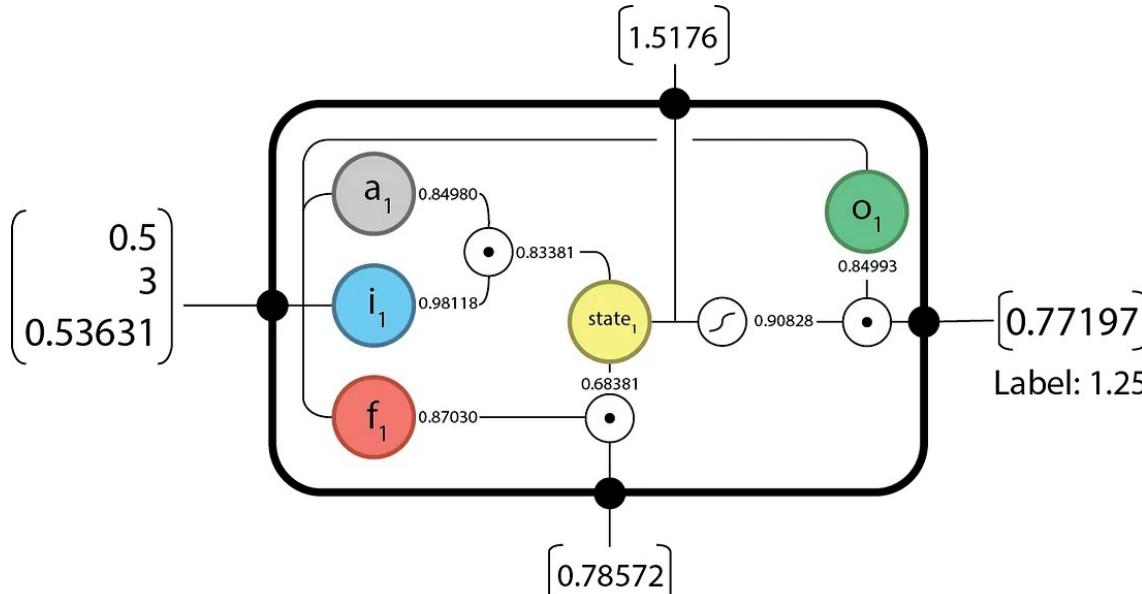
$$o_0 = \sigma(W_o \cdot x_0 + U_o \cdot out_{-1} + b_o) = \sigma([0.6 \ 0.4] \begin{bmatrix} 1 \\ 2 \end{bmatrix} + [0.25] [0] + [0.1]) = 0.81757$$

$$state_0 = a_0 \odot i_0 + f_0 \odot state_{-1} = 0.81775 \times 0.96083 + 0.85195 \times 0 = 0.78572$$

$$out_0 = \tanh(state_0) \odot o_0 = \tanh(0.78572) \times 0.81757 = 0.53631$$

From here, we can pass forward our state and output and begin the next time-step.

Forward @ t=1



$$a_1 = \tanh(W_a \cdot x_1 + U_a \cdot out_0 + b_a) = \tanh([0.45 \ 0.25] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.15] [0.53631] + [0.2]) = 0.84980$$

$$i_1 = \sigma(W_i \cdot x_1 + U_i \cdot out_0 + b_i) = \sigma([0.95 \ 0.8] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.8] [0.53631] + [0.65]) = 0.98118$$

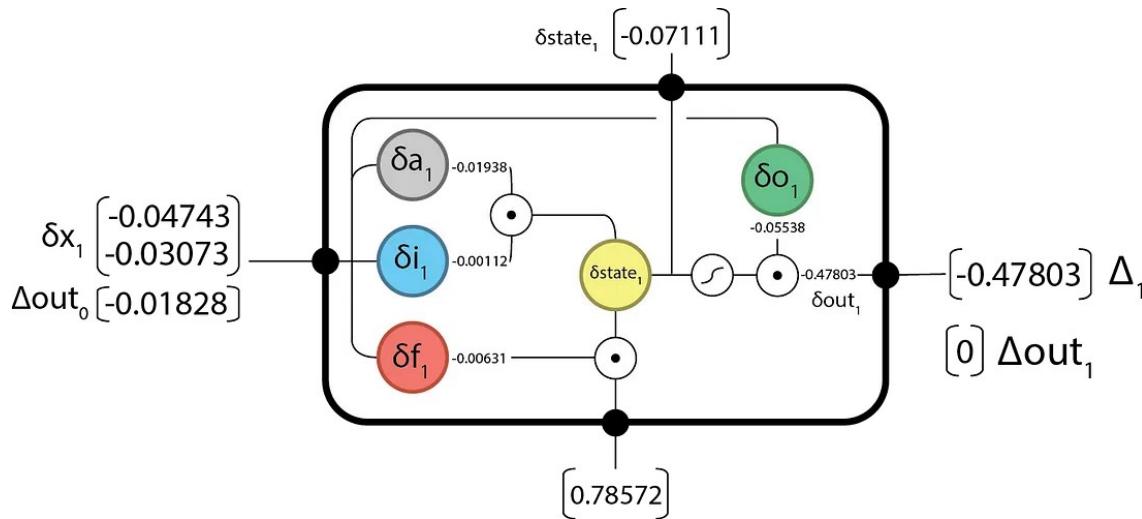
$$f_1 = \sigma(W_f \cdot x_1 + U_f \cdot out_0 + b_f) = \sigma([0.7 \ 0.45] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.1] [0.53631] + [0.15]) = 0.87030$$

$$o_1 = \sigma(W_o \cdot x_1 + U_o \cdot out_0 + b_o) = \sigma([0.6 \ 0.4] \begin{bmatrix} 0.5 \\ 3 \end{bmatrix} + [0.25] [0.53631] + [0.1]) = 0.84993$$

$$\begin{aligned} state_1 &= a_1 \odot i_1 + f_1 \odot state_0 = 0.84980 \times 0.98118 + 0.87030 \times 0.78572 = 1.5176 \\ out_1 &= \tanh(state_1) \odot o_1 = \tanh(1.5176) \times 0.84993 = 0.77197 \end{aligned}$$

And since we're done our sequence we have everything we need to begin backpropogating.

Backward @ t=1



First we'll need to compute the difference in output from the expected (label).

Note for this we'll be using L2 Loss:

$$E(x, \hat{x}) = \frac{(x - \hat{x})^2}{2}$$

The derivate w.r.t. x is:

$$\partial_x E(x, \hat{x}) = x - \hat{x}$$

So,

$$\Delta_1 = \partial_x E = 0.77197 - 1.25 = -0.47803$$

$\Delta out_1 = 0$ because there are no future time-steps.

$$\delta out_1 = \Delta_1 + \Delta out_1 = -0.47803 + 0 = -0.47803$$

$$\delta state_1 = \delta out_1 \odot o_1 \odot (1 - \tanh^2(state_1)) + \delta state_2 \odot f_2 = -0.47803 \times 0.84993 \times (1 - \tanh^2(1.5176)) + 0 \times 0 = -0.07111$$

$$\delta a_1 = \delta state_1 \odot i_1 \odot (1 - a_1^2) = -0.07111 \times 0.98118 \times (1 - 0.84980^2) = -0.01938$$

$$\delta i_1 = \delta state_1 \odot a_1 \odot i_1 \odot (1 - i_1) = -0.07111 \times 0.84980 \times 0.98118 \times (1 - 0.98118) = -0.00112$$

$$\delta f_1 = \delta state_1 \odot state_0 \odot f_1 \odot (1 - f_1) = -0.07111 \times 0.78572 \times 0.87030 \times (1 - 0.87030) = -0.00631$$

$$\delta o_1 = \delta out_1 \odot \tanh(state_1) \odot o_1 \odot (1 - o_1) = -0.47803 \times \tanh(1.5176) \times 0.84993 \times (1 - 0.84993) = -0.05538$$

$$\delta x_1 = W^T \cdot \delta gates_1$$

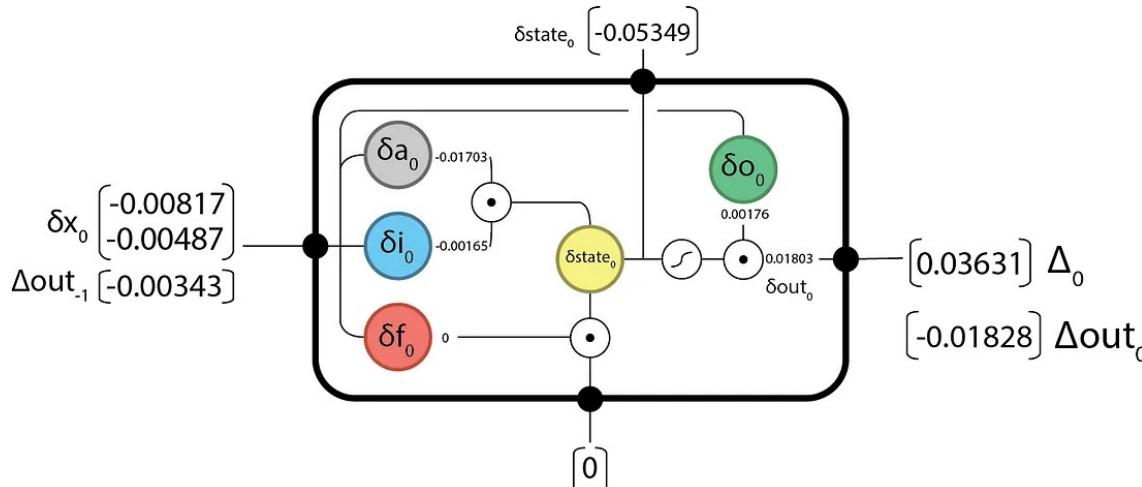
$$= \begin{bmatrix} 0.45 & 0.95 & 0.70 & 0.60 \\ 0.25 & 0.80 & 0.45 & 0.40 \end{bmatrix} \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} = \begin{bmatrix} -0.04743 \\ -0.03073 \end{bmatrix}$$

$$\Delta out_0 = U^T \cdot \delta gates_1$$

$$= \begin{bmatrix} 0.15 & 0.80 & 0.10 & 0.25 \end{bmatrix} \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} = -0.01828$$

Now we can pass back our Δout and continue on computing...

Backward @ t=0



$$\Delta_0 = \partial_x E = 0.53631 - 0.5 = 0.03631$$

$$\Delta out_0 = -0.01828, \text{ passed back from T=1}$$

$$\delta out_0 = \Delta_0 + \Delta out_0 = 0.03631 + -0.01828 = 0.01803$$

$$\delta state_0 = \delta out_0 \odot o_0 \odot (1 - \tanh^2(state_0)) + \delta state_1 \odot f_1 = 0.01803 \times 0.81757 \times (1 - \tanh^2(0.78572)) + -0.07111 \times 0.87030 = -0.05349$$

$$\delta a_0 = \delta state_0 \odot i_0 \odot (1 - a_0^2) = -0.05349 \times 0.96083 \times (1 - 0.81775^2) = -0.01703$$

$$\delta i_0 = \delta state_0 \odot a_0 \odot i_0 \odot (1 - i_0) = -0.05349 \times 0.81775 \times 0.96083 \times (1 - 0.96083) = -0.00165$$

$$\delta f_0 = \delta state_0 \odot state_{-1} \odot f_0 \odot (1 - f_0) = -0.05349 \times 0 \times 0.85195 \times (1 - 0.85195) = 0$$

$$\delta o_0 = \delta out_0 \odot \tanh(state_0) \odot o_0 \odot (1 - o_0) = 0.01803 \times \tanh(0.78572) \times 0.81757 \times (1 - 0.81757) = 0.00176$$

$$\delta x_0 = W^T \cdot \delta gates_0$$

$$= \begin{bmatrix} 0.45 & 0.95 & 0.70 & 0.60 \\ 0.25 & 0.80 & 0.45 & 0.40 \end{bmatrix} \begin{bmatrix} -0.01703 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} = \begin{bmatrix} -0.00817 \\ -0.00487 \end{bmatrix}$$

$$\Delta out_{-1} = U^T \cdot \delta gates_1$$

$$= \begin{bmatrix} 0.15 & 0.80 & 0.10 & 0.25 \end{bmatrix} \begin{bmatrix} -0.01703 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} = -0.00343$$

And we're done the backward step!

Now we'll need to update our internal parameters according to whatever solving algorithm you've chosen.
I'm going to use a simple Stochastic Gradient Descent (SGD) update with learning rate: $\lambda=0.1$.

We'll need to compute how much our weights are going to change by:

$$\begin{aligned} \delta W &= \sum_{t=0}^T \delta gates_t \otimes x_t \\ &= \begin{bmatrix} -0.01703 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} [1.0 \ 2.0] + \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} [0.5 \ 3.0] = \begin{bmatrix} -0.02672 & -0.0922 \\ -0.00221 & -0.00666 \\ -0.00316 & -0.01893 \\ -0.02593 & -0.16262 \end{bmatrix} \\ \delta U &= \sum_{t=0}^{T-1} \delta gates_{t+1} \otimes out_t \\ &= \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} [0.53631] = \begin{bmatrix} -0.01039 \\ -0.00060 \\ -0.00338 \\ -0.02970 \end{bmatrix} \\ \delta b &= \sum_{t=0}^T \delta gates_{t+1} \\ &= \begin{bmatrix} -0.01703 \\ -0.00165 \\ 0 \\ 0.00176 \end{bmatrix} + \begin{bmatrix} -0.01938 \\ -0.00112 \\ -0.00631 \\ -0.05538 \end{bmatrix} = \begin{bmatrix} -0.03641 \\ -0.00277 \\ -0.00631 \\ -0.05362 \end{bmatrix} \end{aligned}$$

And updating out parameters based on the SGD update function:

$$W^{new} = W^{old} - \lambda * \delta W^{old}$$

$$W_a = \begin{bmatrix} 0.45267 \\ 0.25922 \end{bmatrix}, U_a = [0.15104], b_a = [0.20364]$$

$$W_i = \begin{bmatrix} 0.95022 \\ 0.80067 \end{bmatrix}, U_i = [0.80006], b_i = [0.65028]$$

$$W_f = \begin{bmatrix} 0.70031 \\ 0.45189 \end{bmatrix}, U_f = [0.10034], b_f = [0.15063]$$

$$W_o = \begin{bmatrix} 0.60259 \\ 0.41626 \end{bmatrix}, U_o = [0.25297], b_o = [0.10536]$$

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$

$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$

Every row in the \mathbf{X} matrix corresponds to a word in the input sentence. We again see the difference in size of the embedding vector (512, or 4 boxes in the figure), and the q/k/v vectors (64, or 3 boxes in the figure)

$$\text{softmax} \left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V} = \mathbf{Z}$$

The Beast With Many Heads

The paper further refined the self-attention layer by adding a mechanism called “multi-headed” attention. This improves the performance of the attention layer in two ways:

1. It expands the model’s ability to focus on different positions. Yes, in the example above, \mathbf{z}_1 contains a little bit of every other encoding, but it could be dominated by the actual word itself. If we’re translating a sentence like “The animal didn’t cross the street because it was too tired”, it would be useful to know which word “it” refers to.
2. It gives the attention layer multiple “representation subspaces”. As we’ll see next, with multi-headed attention we have not only one, but multiple sets of Query/Key/Value weight matrices (the Transformer uses eight attention heads, so we end up with eight sets for each encoder/decoder). Each of these sets is randomly initialized. Then, after training, each set is used to project the input embeddings (or vectors from lower encoders/decoders) into a different representation subspace.

