



Deep Learning

Activation function

① Linear function

↳ two classes separated by a line

↳ fails over complex data

② Non-linear function

(i) Sigmoid

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

(ii) ReLU

$$\max(0, x)$$

(iii) tanh

$$\tanh(x)$$

(iv) Leaky ReLU

$$\max(0.1x, x)$$

when data gets more complicated we start adding more layers.

but with increasing number of layers the layers stop learning, and in classification, the boundary no longer changes.

Vanishing Gradient problem

during back propagation the gradients gets smaller and smaller and vanishes.

and as

$$\omega' = \omega - \lambda \frac{\partial}{\partial \omega}$$

turns out
to be 0

Remedy ↗



Hence, No Learning

Use activation function that
does not squeeze information

ReLU



Again after a while the boundary does not change again.

Dying ReLU Problem

blocks information < 0 .

Leaky ReLU / ELU

Activation function of Outputs

for classification \rightarrow Softmax function

$\# \text{neurons} = \# \text{classes}$

$$\text{Softmax}(x) = \frac{e^x}{\sum e^x}$$

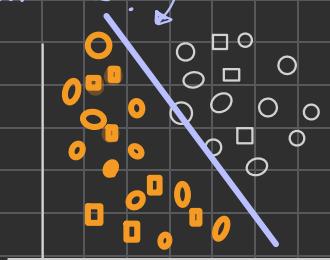
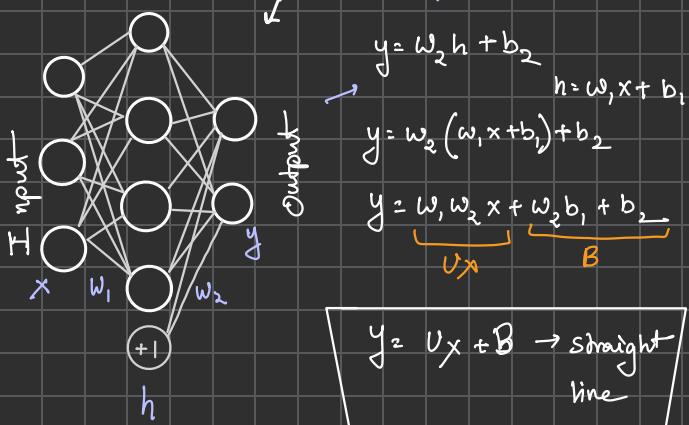
Output probabilities of classes

for regression

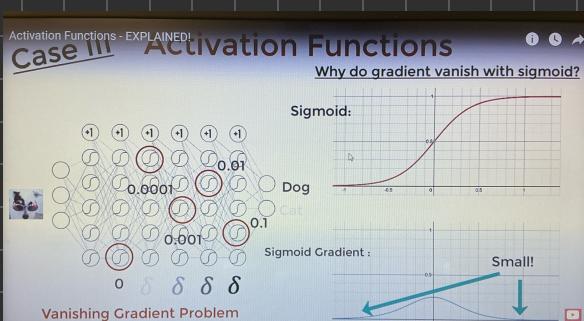
\hookrightarrow Don't need an Activation

function. \rightarrow based on the type of output we want, we can choose the activation function

In depth \rightarrow Why linear function produces a straight line?



② Why Sigmoid causes vanishing gradient?



with larger +ve/-ve value, Sigmoid gradient tends to 0, causing it to minimize.

and with backpropagation, with multiplication it goes further lower \rightarrow making it 0.

due to squeezing nature of Sigmoid function. [takes a number and squeezes it to a range]

One more problem that comes in with increasing # layers. \rightarrow

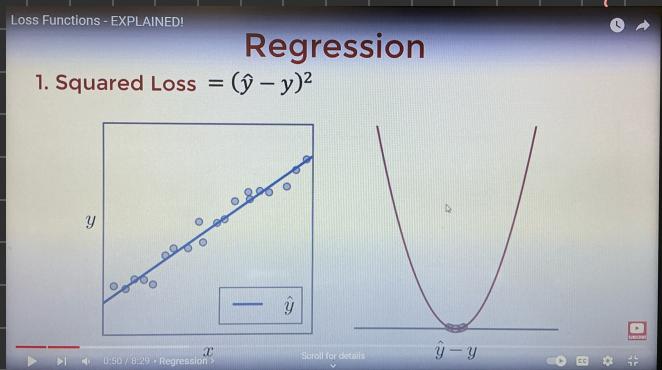
Exploding gradient

When $\frac{\partial L}{\partial x} \rightarrow$ goes very large, we can not converge and

hence we don't reach to minimal loss value \rightarrow underfitting

Loss functions

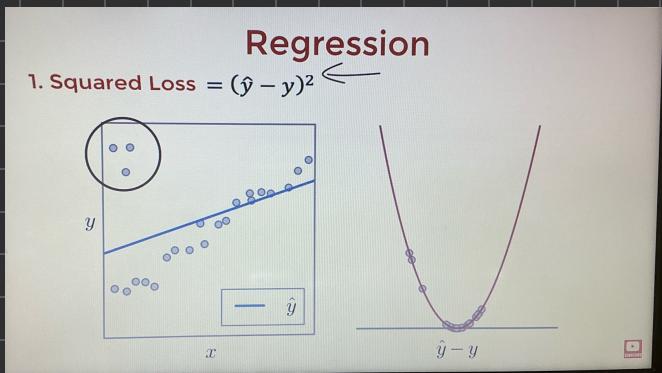
① Regression



the difference is squared or absolute
so, that +ve / -ve differences do not
cancel out.

but with outlier

this performs bad, as the model
wants to fit in the outliers



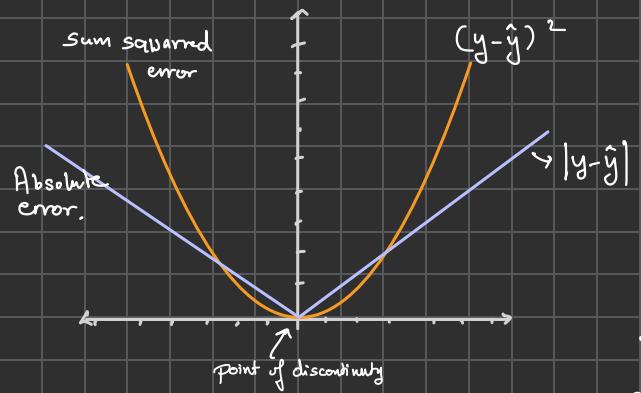
thus

Absolute error

$|\hat{y} - y|$

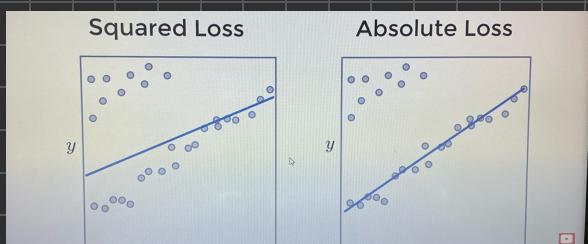
$|\hat{y} - y| \rightarrow$ this might lead to
poor predictions

Support vector Regression



mean absolute error
is not optimized through Gradient descent:
Sub gradients

- One of them considers the outliers
- Other one completely ignores them.
↓
both resulting in poor predictions



Classification

	clauses	probability distribution
A	0.08	
B	0.74	→ max → Hence <u>result</u>
C	0.07	
D	0.11	

we compare it with the ground truth,
and how we compare → depends on the
Loss we use.

① Cross-entropy loss.

↓ we must first understand KL divergence

↓ distance between two probability distributions

$$D_{KL}(P || Q) = \sum P(i) \log \frac{P(i)}{Q(i)}$$

or

$$D_{KL}(P || Q) = \int P(x) \log \frac{P(x)}{Q(x)} dx$$

In depth

Coin 1

$$\begin{cases} 0.5 & \text{Heads} \\ 0.5 & \text{Tails} \end{cases}$$

Coin 2

$$\begin{cases} 0.95 & \text{Heads} \\ 0.05 & \text{Tails} \end{cases}$$

These two coins are very different, but How different are they?
quantitatively

$$\frac{P(\text{observation} | \text{real coin})}{P(\text{observation} | \text{coin 2})}$$

True Coin

$$\begin{cases} P_1 & \text{Heads} \\ P_2 & \text{Tails} \end{cases}$$

Coin 2

$$\begin{cases} q_1 & \text{Heads} \\ q_2 & \text{Tails} \end{cases}$$

trials → H H T H H T H H H T H T

$$P_1 \cdot P_1 \cdot P_2 \cdot P_1 \cdot P_1 \cdot P_2 \cdot P_1 \cdot P_1 \cdot P_2 \cdot P_1 \cdot P_1 \rightarrow P_1^{\#H} P_2^{\#T}$$

$$q_1 \cdot q_2 \cdot q_2 \cdot q_1 \cdot q_1 \cdot q_2 \cdot q_1 \cdot q_1 \cdot q_1 \cdot q_2 \cdot q_1 \cdot q_2 \rightarrow q_1^{\#H} q_2^{\#T}$$

$$\frac{P(\text{observation} \mid \text{real coin})}{P(\text{observation} \mid \text{coin 2})} = \frac{P_1^{\#H} P_2^{\#T}}{q_1^{\#H} q_2^{\#T}}$$



Normalizing it based on sample size

$$= \left[\frac{P_1^{\#H} P_2^{\#T}}{q_1^{\#H} q_2^{\#T}} \right]^{1/N}$$

$$\log \left[\frac{P_1^{\#H} P_2^{\#T}}{q_1^{\#H} q_2^{\#T}} \right]^{1/N} \rightarrow \frac{1}{N} \log \left[\frac{P_1^{\#H} P_2^{\#T}}{q_1^{\#H} q_2^{\#T}} \right]$$

$$= \frac{1}{N} \log P_1^{\#H} + \frac{1}{N} \log P_2^{\#T} - \frac{1}{N} \log q_1^{\#H} - \frac{1}{N} \log q_2^{\#T}$$

$$= \frac{n}{N} \log P_1 + \frac{T}{N} \log P_2 \rightarrow \frac{n}{N} \log q_1 - \frac{T}{N} \log q_2$$



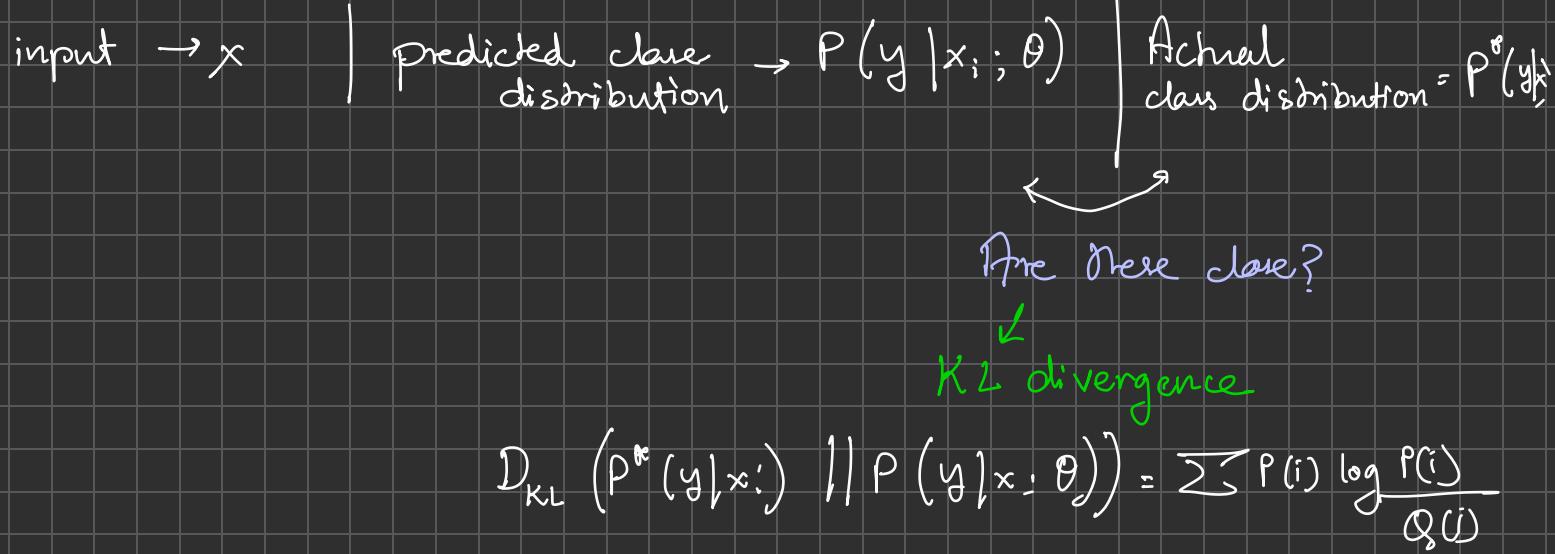
considering sample size tending $\rightarrow \infty$

we assume $\rightarrow P_1$

$$= P_1 \log P_1 + P_2 \log P_2 - P_1 \log q_1 - P_2 \log q_2$$

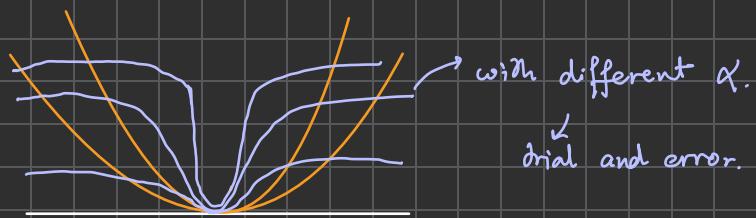
$$= P_1 \log \frac{P_1}{q_1} + P_2 \log \frac{P_2}{q_2}$$

$$= \sum p_i \log \frac{p_i}{q_i}$$



② Hinge Loss \leftarrow used in Support Vector machine

③ for Regression \rightarrow Adaptive loss



Weight initialization

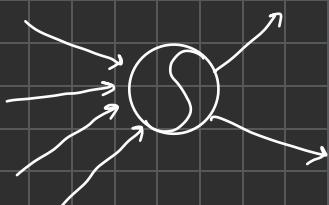
\hookrightarrow to keep in mind

① weight should be small

② should have good variance, each neuron must carry an information within itself.

different neurons will learn different aspects of input.

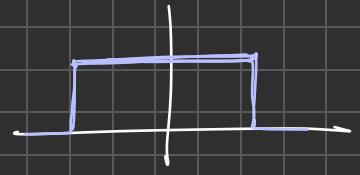
randomly initialization



fan-in $\rightarrow 4$
fan out $\rightarrow 2$

① Uniform initialization

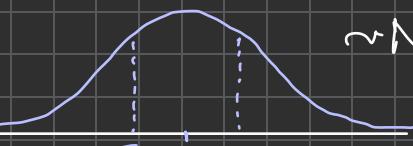
$$\rightarrow \left[-\frac{1}{\sqrt{\text{fanIn}}}, \frac{1}{\sqrt{\text{fanIn}}} \right]$$



② Xavier / glorot initialization

Normal ↴

$$\sim N(0, \frac{2}{\text{fanIn} + \text{fanOut}})$$



$$\rightarrow \text{Uniform} \rightarrow \sim U \left[-\frac{\sqrt{6}}{\sqrt{\text{fanIn} + \text{fanOut}}}, \frac{\sqrt{6}}{\sqrt{\text{fanIn} + \text{fanOut}}} \right]$$



③ He-init {often used with ReLU}

Normal ↴

$$\sim N(0, \sqrt{\frac{2}{\text{fanIn}}})$$

Uniform ↴

$$\sim U \left[-\sqrt{\frac{6}{\text{fanIn}}}, \sqrt{\frac{6}{\text{fanIn}}} \right]$$

Optimizers

→ gradient descent

for every epoch:

$$\omega' = \omega - \alpha \frac{\partial L}{\partial \omega}$$

weight is only updated once after seeing the entire dataset.

and the weight just hovers at its optimal value, rather than actually reaching there.

update the parameters more frequently.

e.g: SGD

↓
but SGD is very noisy, as it gets influenced by every single sample

↙
Mini batch GD

Another way to handle this is by introducing Momentum instead of taking just one current gradient, it also takes in account the previous update to smooth out the oscillations

↙
faster convergence + Less Oscillations

- but
- Too much momentum can overshoot the minima
 - As the optimizer gets closer to minima, the gradient become smaller, → leading to slower convergence., and it can still overshoot.

$$\left\{ \begin{array}{l} w_{t+1} = w_t - v_t \\ v_t = \beta v_{t-1} + \eta \nabla w_t \end{array} \right.$$

↓
SGD + Momentum + Acceleration

Nesterov Accelerated gradient → Instead of updating, solely based on current gradient,

NAG computes the gradient at the predicted future, (where the momentum might take it.) → controlling momentum

$$w_{\text{new}} = w_{\text{old}} - v_{\text{new}}$$

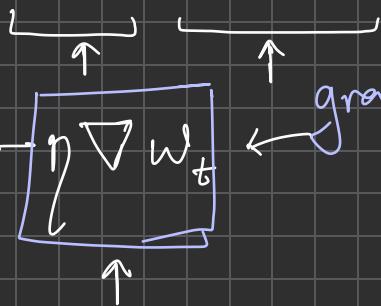
$$v_{\text{new}} = \beta v_{\text{old}} + \eta L (w_{\text{old}} - \beta \cdot v_{\text{old}})$$

Momentum

$$\omega_{t+1} = \omega_t - (\beta v_{t-1} + \gamma \nabla w_t)$$

earlier,

$$\omega_{t+1} = \omega_t$$



NAG →

Next step depends upon

History of Velocity + gradient at that position

Seq,

①

②

for momentum you go sequential.

The only difference with NAG is → we apply momentum and then calculate gradient for the next step point.

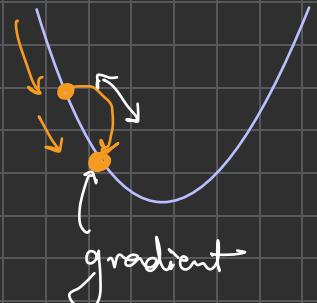
Hence, NAG performs better.

Look ahead → la

$$\omega_{la} = \omega_t - \beta v_{t-1}$$

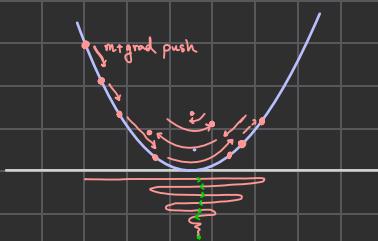
$$v_t = \beta v_{t-1} + \gamma \nabla w_{la}$$

$$\omega_{t+1} = \omega_t - v_t$$

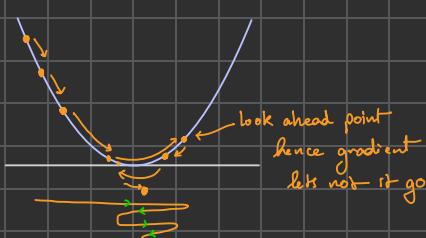


geometric intuition

Momentum

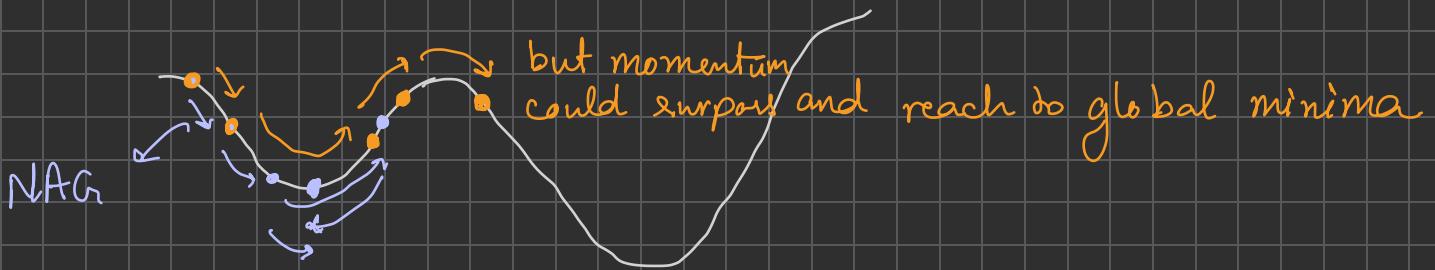


momentum + Acceleration → NAG



Hence, decreasing # oscillations

but the problem with NAG, because it dampens the oscillations
might, saddle point at local minima



Adagrad \rightarrow very similar to SGD, but with adaptive gradient.

$$\text{SGD} \rightarrow \omega_t = \omega_{t-1} - \eta \Delta g_t$$

$$\text{Adagrad} \rightarrow \omega_t = \omega_{t-1} - \eta' \Delta g_t$$

\downarrow different learning rate at each weight

$$\eta' = \frac{\eta}{\sqrt{G_t + \epsilon}}$$

$\sum_{t=1}^{t-1} g_i^2$ \downarrow sum of squared of all prev. gradients

$$\omega_t = \omega_{t-1} - \frac{\eta}{\sqrt{\sum_i g_i^2 + \epsilon}} \cdot \Delta L$$

but if $\sum g_i^2 \rightarrow$ goes very high

\downarrow Slow convergence

RMS prop

{ Root Mean Square Propagation }

RMS Prop → It only considers recent gradients through exponentially decaying average, allowing it to make more reasonable learning rate throughout training.

$$\eta_t = \frac{\eta_{t-1}}{W_{\text{avg}} + \epsilon} \quad \left| \quad W_{\text{avg}} = \alpha W_{\text{avg}}(t-1) + (1-\alpha) \sum \left(\frac{\partial L}{\partial w} \right)^2 \right.$$

Adam Optimizer → It combines two ideas

[Direction & Speed]

[Step Size]

① keep track of average of gradients
↓
to capture the direction in which params are moving.

② keep track of squared average
↓
to adapt the learning rate for each parameters.

① First moving average [Momentum] : Okay, I have been moving down hill for awhile, let's keep going in this direction

$$w_{\text{new}} = 0.9 w_{\text{old}} + 0.1 \text{ current gradient}$$

② Second moving average [learning rate adjustment] : If I am on

$$v_{\text{new}} = 0.999 v_{\text{old}} + 0.001 \cdot \begin{cases} \text{current gradient} \\ \end{cases}^2$$

steep hill → smaller steps
flat → larger steps

③ Adjust the bias → initially because the averages are just starting they might be incorrect (too small). Thus

needs to be corrected

$$\hookrightarrow ④ w_{\text{new}} = w_{\text{old}} - \frac{\text{momentum}}{\sqrt{\text{adjusted step size}} + \epsilon}$$

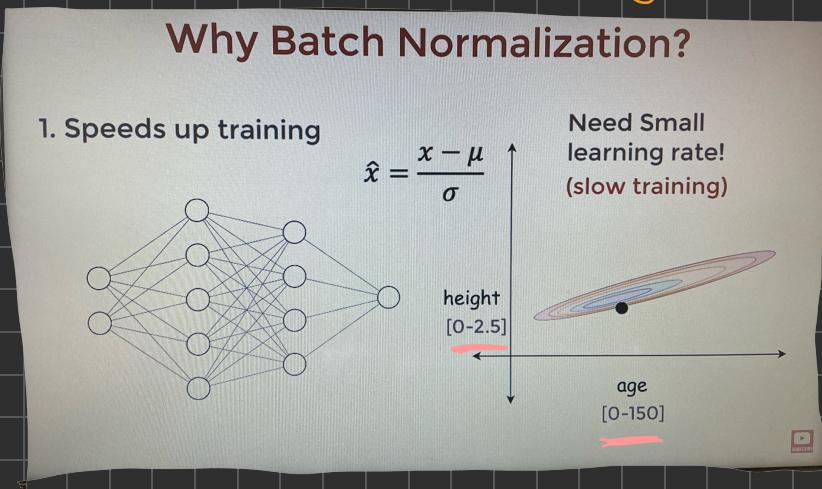
Batch Normalization

Normalizing input speeds up the training process.

Then why not normalize activation layer.

Why batch normalization?

① Speed up training



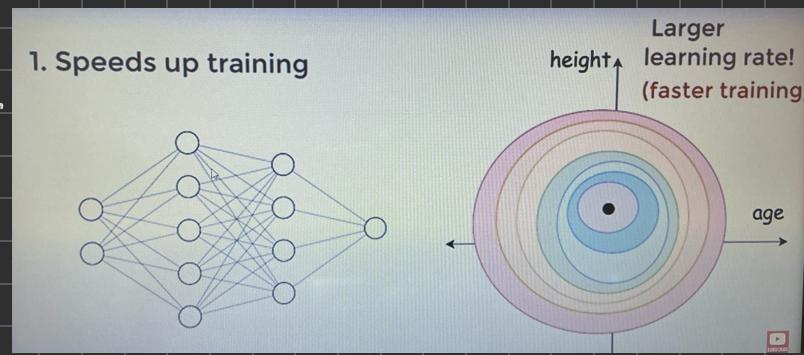
Small variation in Height can greatly change the loss.



and Hence, we would need small learning rate, so that we don't shoot the minimum

↓ with normalizing

With normalizing we can take larger learning rate
thus reaching minimum faster. That said, we can also use adam.



② Allow suboptimal starts → freedom with initial weights, as with batch normalization initial weights are less important.

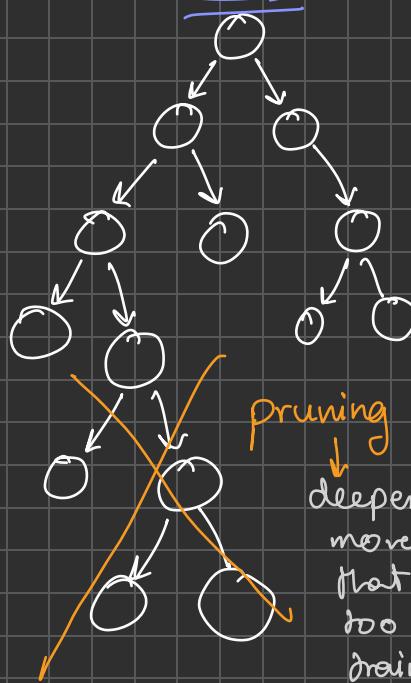
③ little bit → acts as regularizer

Regularization → technique to reduce overfitting

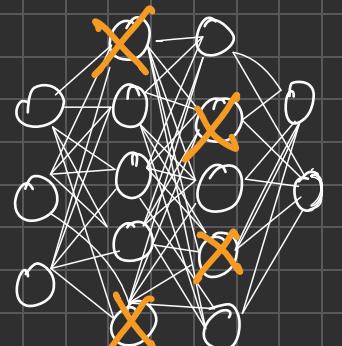
Linear Regression / GLM

$$\sum_i (y - \hat{y})^2 + \lambda \sum_j \beta_j^2$$

Decision Tree



Neural Network



dropouts

2 ways to regularize

① Constraining the model → constraining the flexibility of model or fewer degrees of freedom. → to handle overfitting

Dropout | L₁ / L₂ Reg | Early stopping

Convolutional Neural Network

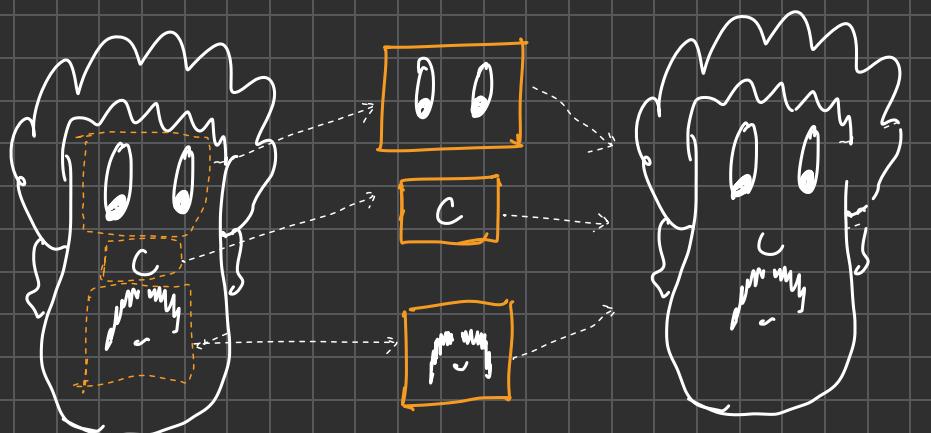
ANN is sufficiently good for image classification of handwritten note. However, if there is a shift or complexity like a current picture ANN by itself would not be sufficiently handle it.

Image Size $\rightarrow 1920 \times 1080 \times 3 \sim \rightarrow 6 \text{ million} \rightarrow 4 \text{ million}$

1st layer 2nd layer
many weights.

- Too much computation.

CNN



convolution /

filter operation

stride = 1

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

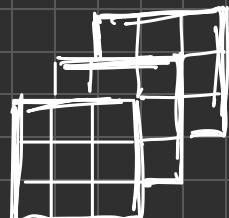
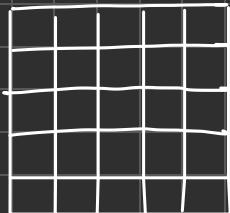
$$\begin{matrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$$\downarrow$$

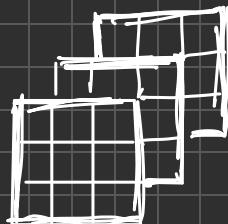
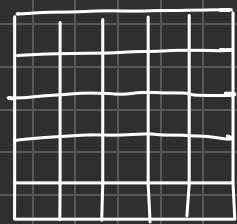
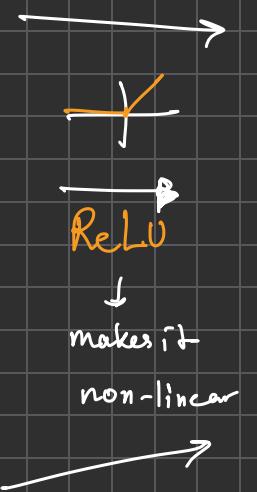
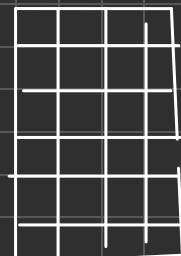
Convolution

-0.11	1	

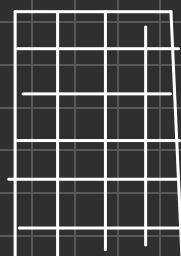
feature map



→



→



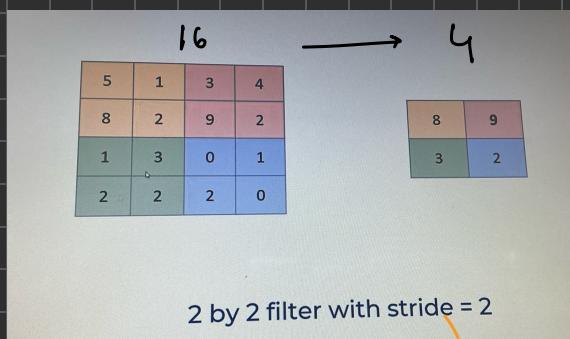
ReLU
makes it non-linear



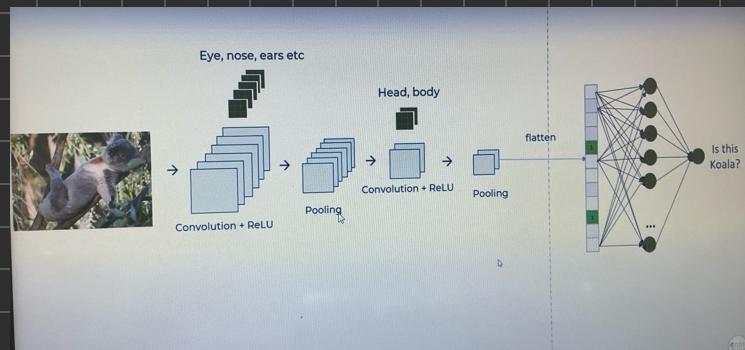
But the problem of computation is still there
of

Pooling → to reduce the dimension

Max pooling



Complete CNN



Convolution Layer

- Connection Sparsity → reduces overfitting
- Conv + Pooling → location invariant feature detection

ReLU

- reduces non-linearity
- faster to compute

Pooling

- reduces dim and hence computation
- reduces overfitting

rotation and thickness can not be handled

↙

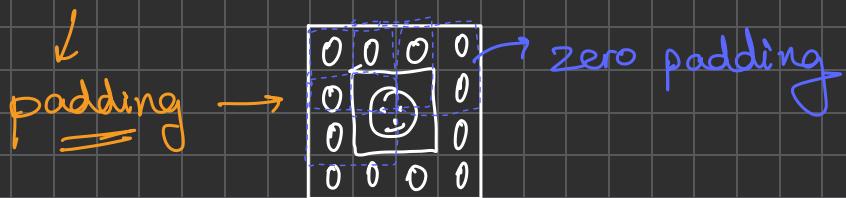
augment training samples → manually

The idea is the weights of the kernel, will be randomly initialized and with forward and backward propagation, iter after iteration, these weights will keep on getting updated. → thus training the model.

Convolution → element wise multiplication + Addition

↓
to capture important element/info from image,
↓
thus, multiple kernels, each carrying essence of individual information

because, we don't want to loose information at the image border.

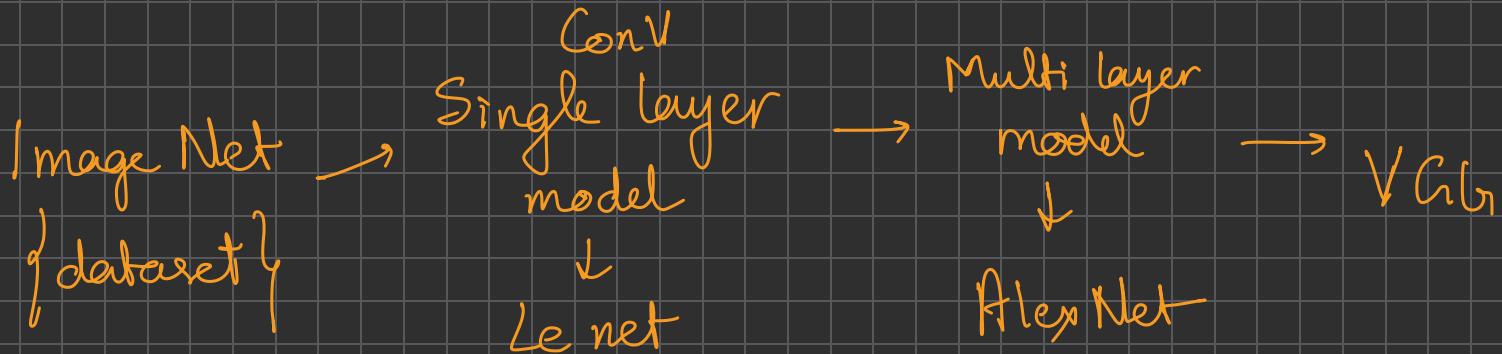


- ① Maintain Output size
- ② Prevent information loss.

strides → with small steps we capture minute details,

while with larger steps, we might miss out on minute details. However, it totally depends on the Scenario.

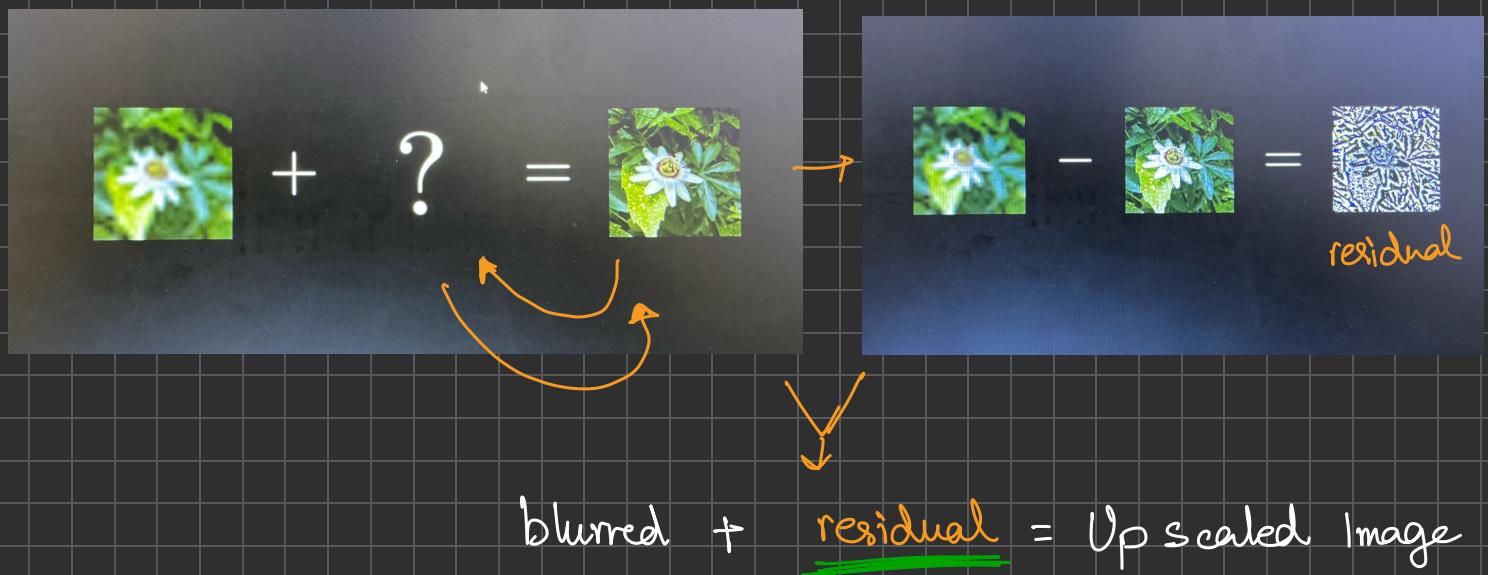
Pooling → by summarizing the context/pixelate/downscaling the image, it solves for → location invariance
scale invariance
rotation invariance



Conv 2D → solves for images

Conv 1D → time-series data

Conv 3D → image + time series data



There was a problem on vanishing gradient,



ResNet → the concept is ; it will skip

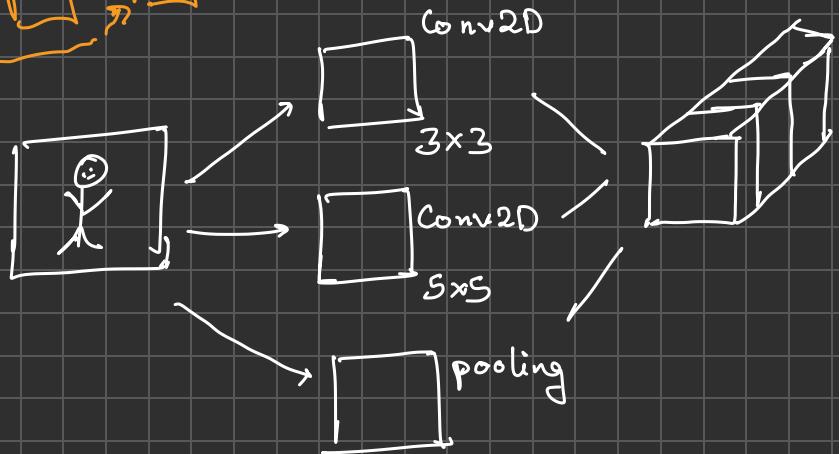
connections thus handling bias mess.

→ ResNet50 → pretrained model

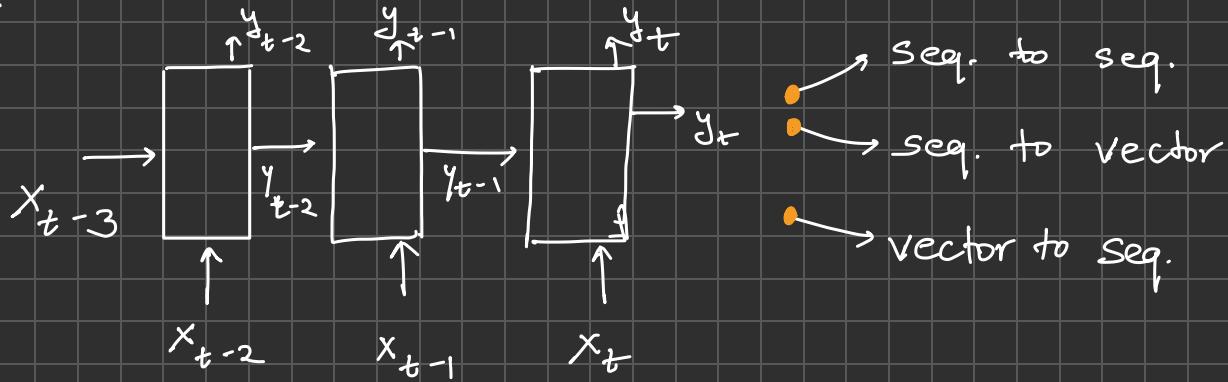
50 residual blocks
↳ each 3/4 CNN



Inception Network →



RNN → Recurrent Neural Network



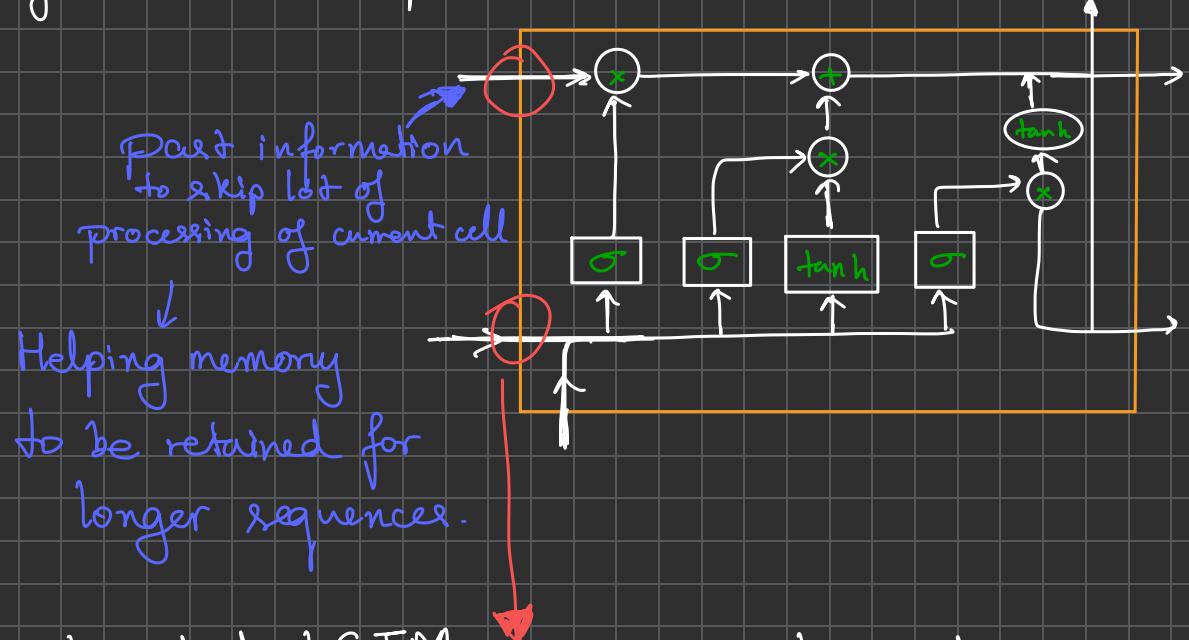
Deals with sequential data

Disadvantages → slow to train

→ fail for long seq. → due to vanishing gradient

LSTM ↴ {long Short term memory}

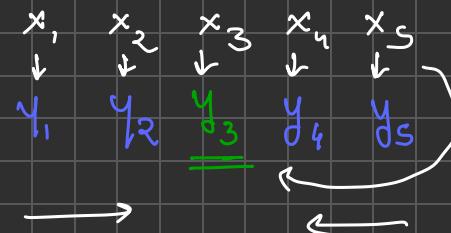
Instead of neurons → replace with LSTM cell



RNNs are slow but LSTMs are even slower, because of complexity. → and the input data needs to be passed sequentially. ↴ Not taking advantage of GPUs at all

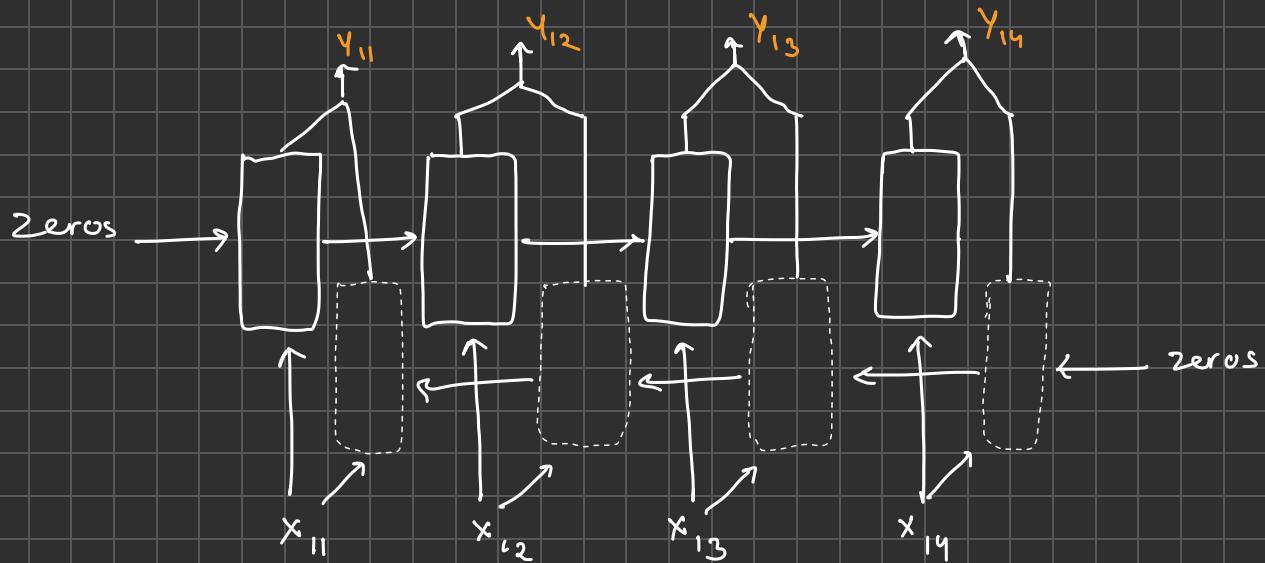
As an alternative we also formed bidirectional RNN.

as in a sentence,



y_3 could be dependent on x_5 .

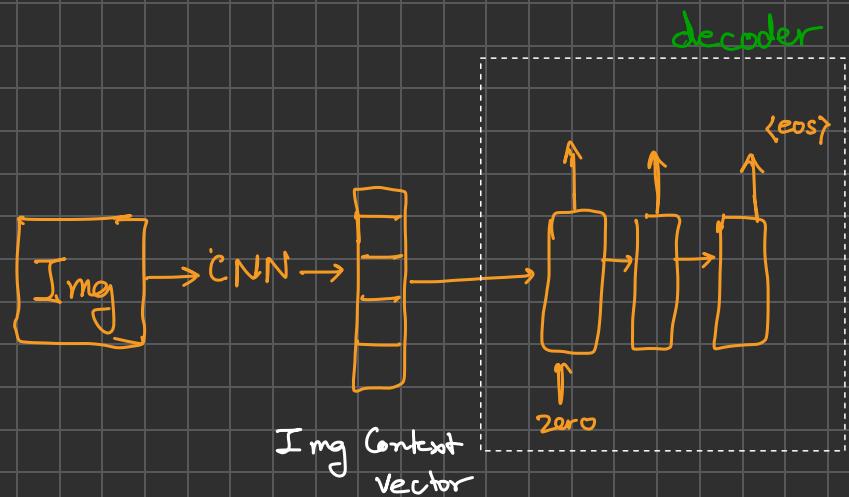
thus, bidirectional RNN



However, for a longer seq. of text, this wasn't a good idea either.

Encoder - Decoder models

① Image Captioning

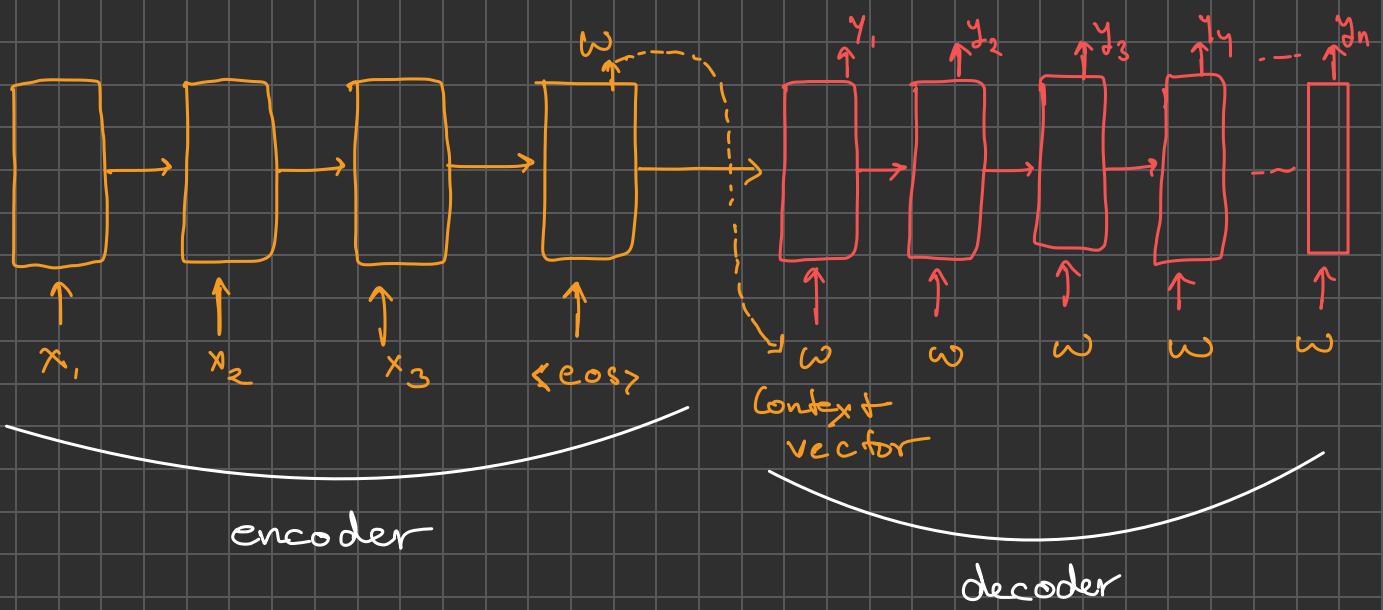


② Seq. to Seq. model → translation

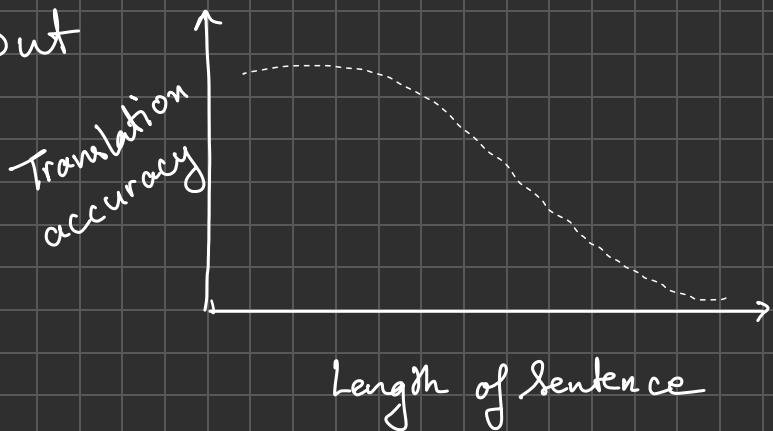
email - Auto reply

code - errors

discussion

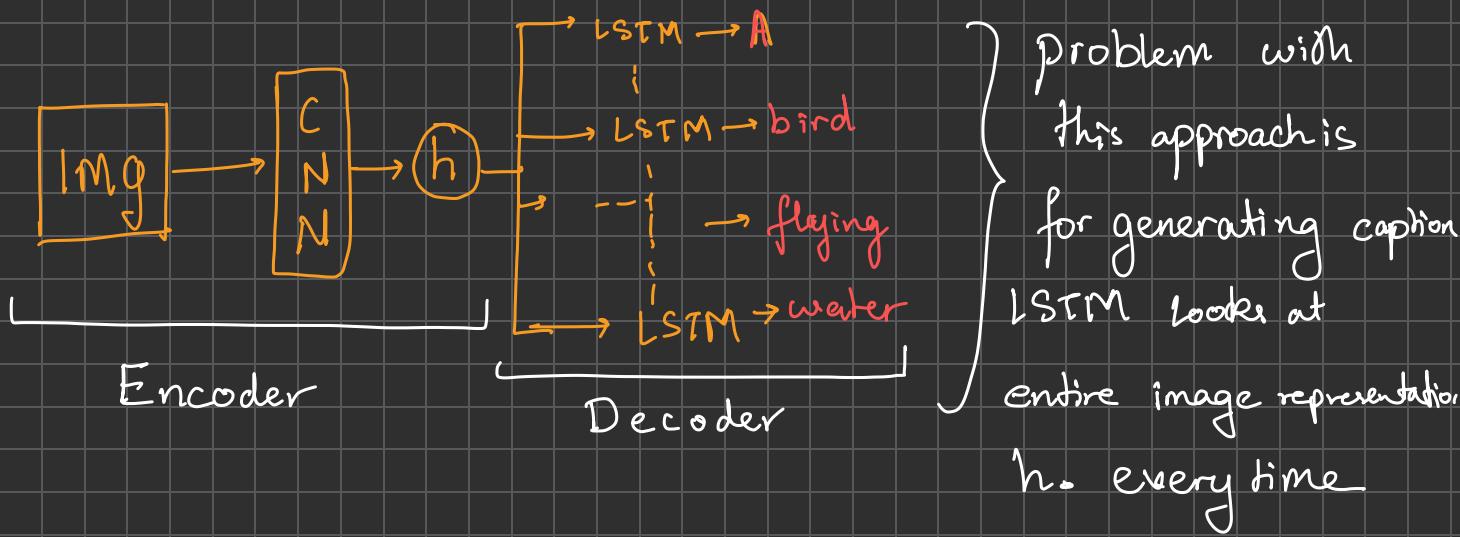


but



for a simple Encoder-Decoder

Attention Mechanism → Giving attention to specific eg: word in a sentence. Just like how we humans do.



to solve this

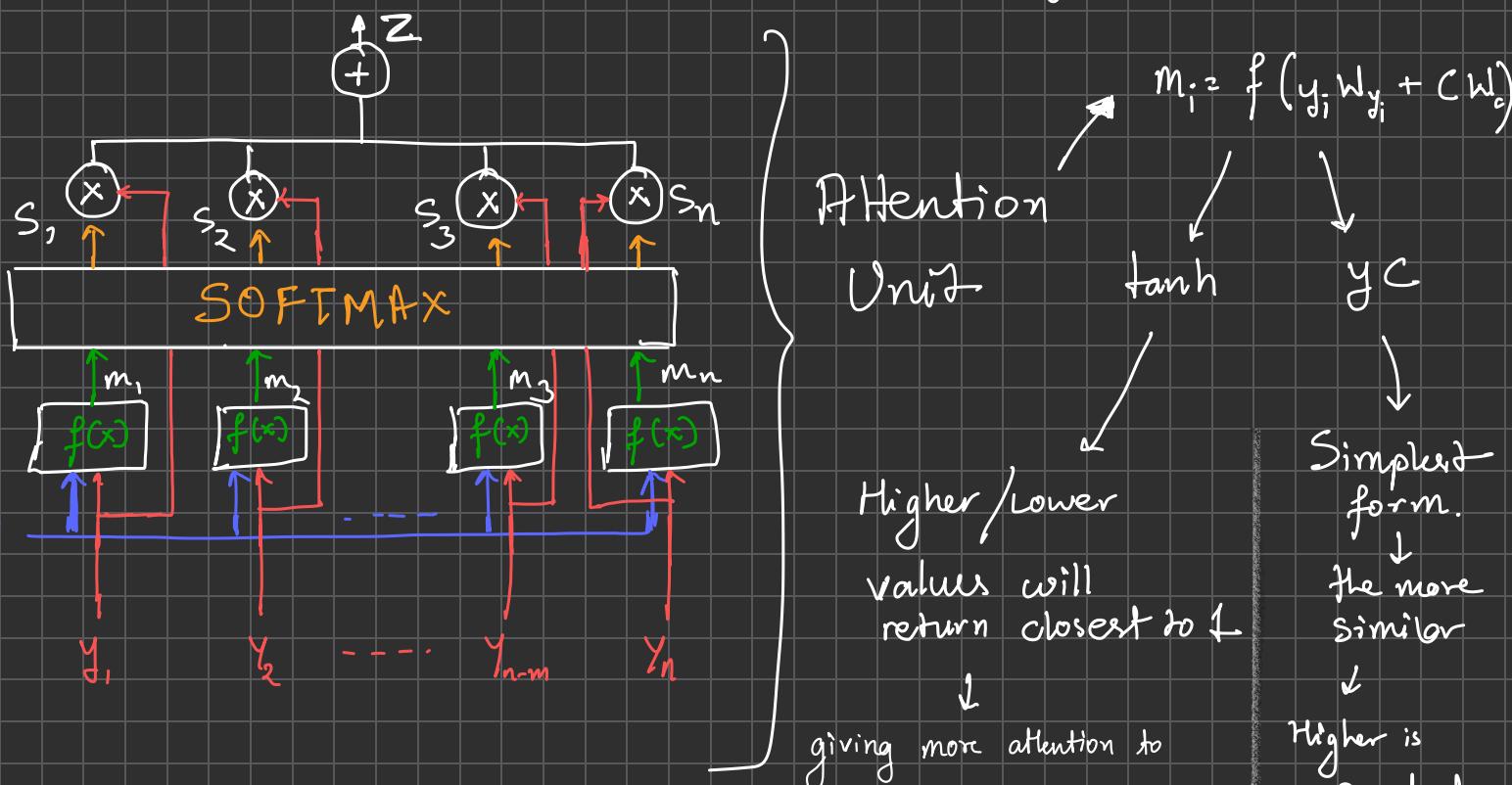


How does it actually decide the region to focus?

$$E(x) = \sum_n P(x = X_n) \cdot X_n \rightarrow \text{inner product of Actual values and their probabilities [Arithmetic Mean]}$$

?

How are these determined? → by context C.



$$s_i = \text{softmax}(m_i)$$

$$= \frac{e^{m_i}}{\sum_n e^{m_n}}$$

$$\rightarrow z = \sum s_i y_i \left\{ \begin{array}{l} \text{inner product} \\ \text{product} \end{array} \right\}$$

Input Representation \rightarrow Embedding Space

↓
each word representing some meaning.

Such that in vector space, words with similar meanings are closer to each other.

We can either pre-train here, or use existing trained models \rightarrow Glove.

$$I \rightarrow [0.3, 0.2, \dots, 0.9]$$

$$am \rightarrow [0.42, 0.16, \dots, 0.24]$$

$$happy \rightarrow [0.29, 0.34, \dots, 0.18]$$

but each word can have different interpretation based on the sentence and position.

defining order of words

Positional encoding.

$$I \rightarrow [0.3, 0.2, \dots, 0.9]$$

$$[0.3, 0.2, \dots, 0.19]$$

$$am \rightarrow [0.42, 0.16, \dots, 0.24]$$

$$+ [0.24, 0.18, \dots, 0.27]$$

$$happy \rightarrow [0.29, 0.34, \dots, 0.18]$$

$$[0.15, 0.16, \dots, 0.94]$$

↓ Contextual embedding vectors

$$I \rightarrow [0.16, \dots, 0.24]$$

$$am \rightarrow [0.14, 0.13, \dots, 0.25]$$

$$happy \rightarrow [-]$$

Encoder Self attention Layer

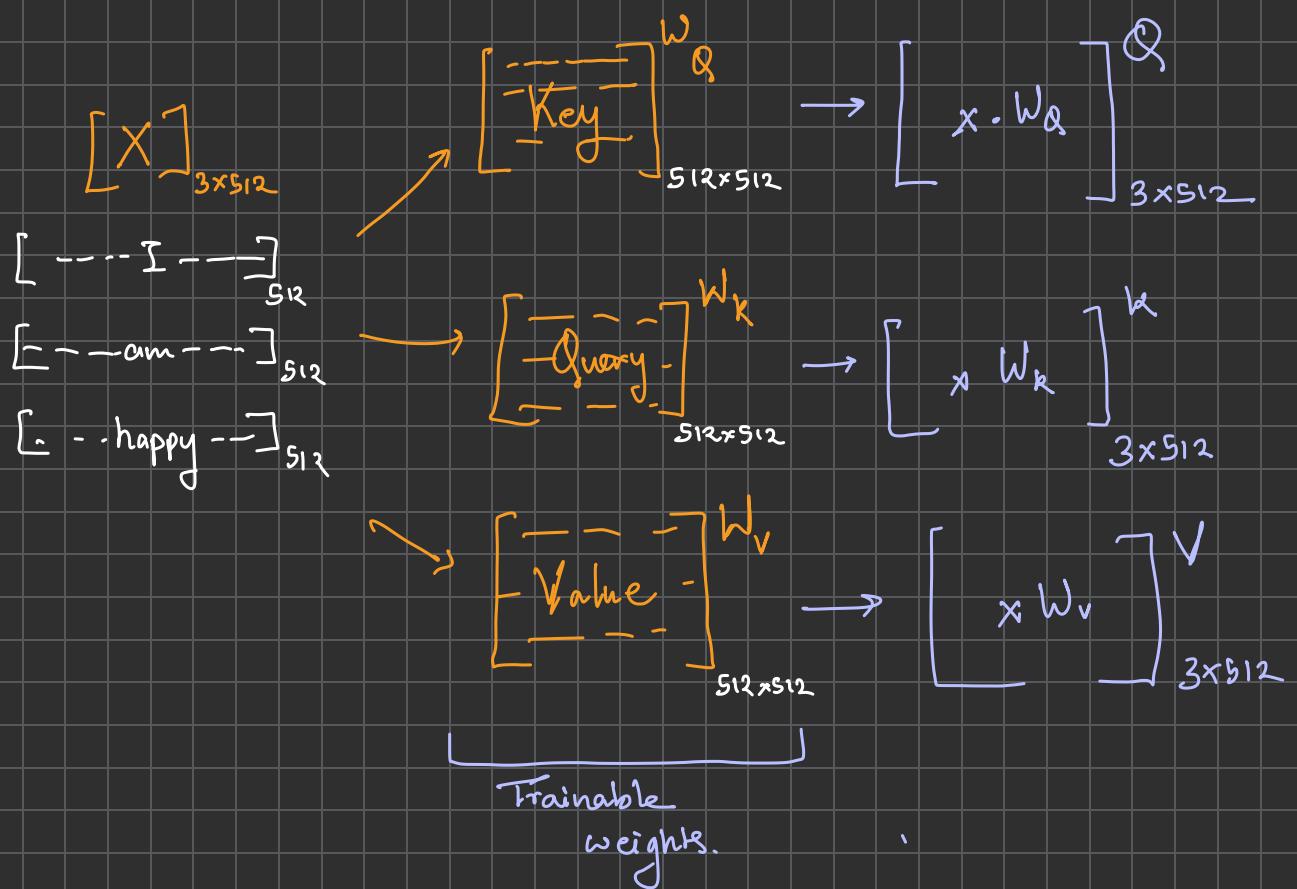
$$I \rightarrow [0.16, \dots, 0.24]$$

$$am \rightarrow [0.14, 0.13, \dots, 0.25]$$

$$happy \rightarrow [-]$$

$$X = \begin{bmatrix} 0.16, \dots, 0.24 \\ 0.14, \dots, 0.25 \\ - \end{bmatrix}$$

3×512



attention score \rightarrow attention score is given by $\rightarrow \text{Attention}(Q, k) = k^T Q$

$$\begin{bmatrix} k_1 \\ k_2 \\ k_3 \end{bmatrix}_{3 \times 512} \times \begin{bmatrix} | & Q \end{bmatrix} = \begin{bmatrix} \cdot \\ A \\ \cdot \end{bmatrix}_{3 \times 3}$$

Attention matrix

multiply because

more similar \rightarrow Higher product \rightarrow more attention

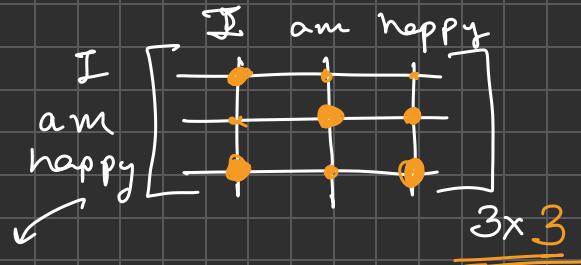
To stabilize the gradients we scale by square root of dim.

$$\begin{bmatrix} A_{11}/\sqrt{3} & \dots \\ \vdots & \vdots \end{bmatrix} \rightarrow \begin{bmatrix} A/\sqrt{dk} \end{bmatrix}$$

Now we send this to a softmax function to get normalized scores. ?

$$\text{Softmax}(A) = \frac{e^{A_{ij}}}{\sum_{j=1}^3 e^{A_{ij}}}$$

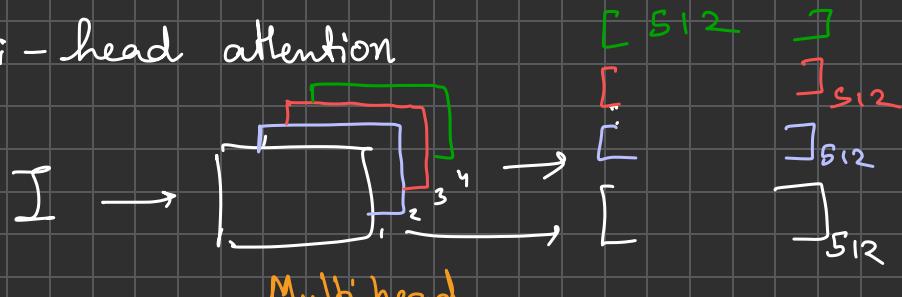
Normalized Attention Matrix



Now we weight it by value

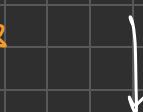
$$A \times V \rightarrow \begin{bmatrix} \dots & I & \dots \\ \dots & am & \dots \\ \dots & happy & \dots \end{bmatrix}_{3 \times 512}$$

For a multi-head attention

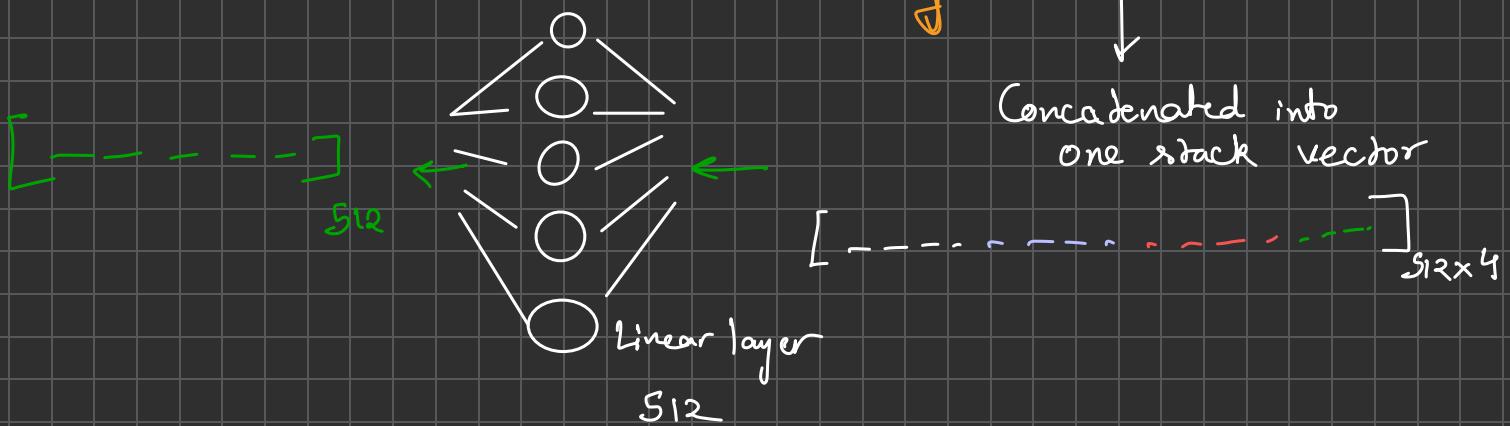


Multi head

attention layers



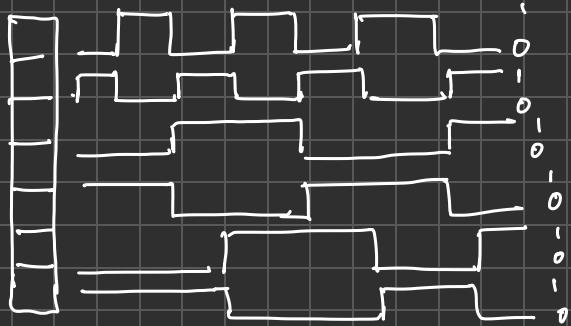
Concatenated into
one stack vector



Positional encoding?

word embedding will carry semantic vector representation, but every word based on the position, can have different interpretation and therefore position embedding.

The idea is we add position encoding as a vector of shape same as embedding layer, so that output remains of some shape



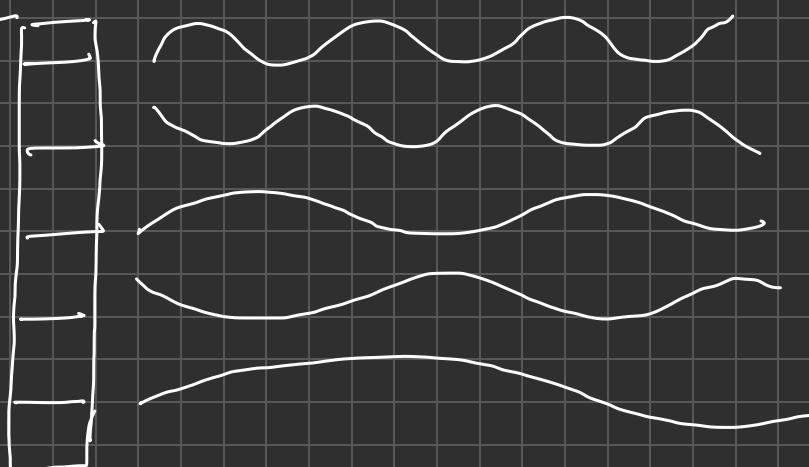
we could have gone ahead with binary values with different frequency for positional encoding.

Vector.

but because for a high-dim vector

or for a larger seq.

we go ahead with smoothed and rotatory trigonometric calculations., ie



thus,

- ① periodic
- ② limited
[-1, +1]

$$\text{for even values} \rightarrow PE(pos, 2i) = \sin\left(\frac{Pos}{1000^{2i/d_model}}\right)$$

$$\text{for odd values} \rightarrow PE(pos, 2i+1) = \cos\left(\frac{Pos}{1000^{2i/d_model}}\right)$$

$i \rightarrow \text{index dim}$ $\left| \begin{array}{l} pos \rightarrow \text{position of word} \\ d_model \rightarrow \text{embedding len} \end{array} \right|$

Self Attention → capture inter-relation

Seq $\rightarrow 6$, $d_{\text{model}} = 512$

$$\text{Attention } (K, Q, V) = \text{softmax} \left(\frac{Q K^T}{\sqrt{d_m}} \right) V$$

$$\text{Softmax} \left(\frac{\begin{matrix} Q \\ 6 \times 512 \end{matrix} \times \begin{matrix} K^T \\ 512 \times 6 \end{matrix}}{\sqrt{512}} \right) = \begin{matrix} \Sigma = 1 \\ 6 \times 6 \end{matrix}$$

$$\begin{matrix} X \\ \times \\ \begin{matrix} V \\ 6 \times 512 \end{matrix} \end{matrix}$$

\downarrow

Attention 6×512 } soft attention matrix

Multi-head Attention

make 3 copies

Input
 6×512
Seq $\times d_{\text{model}}$

$$\begin{aligned} & K \rightarrow \begin{matrix} K \\ 6 \times 512 \end{matrix} \times \begin{matrix} W_K \\ 512 \times 512 \end{matrix} = \begin{matrix} K' \\ 6 \times 512 \end{matrix} \xrightarrow{\substack{d_k \\ \downarrow}} \begin{matrix} K_1 \\ K_2 \\ K_3 \\ \vdots \\ K_n \end{matrix} \text{ Seq} \\ & Q \rightarrow \begin{matrix} Q \\ 6 \times 512 \end{matrix} \times \begin{matrix} W_Q \\ 512 \times 512 \end{matrix} = \begin{matrix} Q' \\ 6 \times 512 \end{matrix} \xrightarrow{\substack{d_k \\ \downarrow}} \begin{matrix} Q_1 \\ Q_2 \\ Q_3 \\ \vdots \\ Q_n \end{matrix} \\ & V \rightarrow \begin{matrix} V \\ 6 \times 512 \end{matrix} \times \begin{matrix} W_V \\ 512 \times 512 \end{matrix} = \begin{matrix} V' \\ 6 \times 512 \end{matrix} \xrightarrow{\substack{d_v \\ \downarrow}} \begin{matrix} V_1 \\ V_2 \\ V_3 \\ \vdots \\ V_n \end{matrix} \end{aligned}$$

to capture different aspect of each word

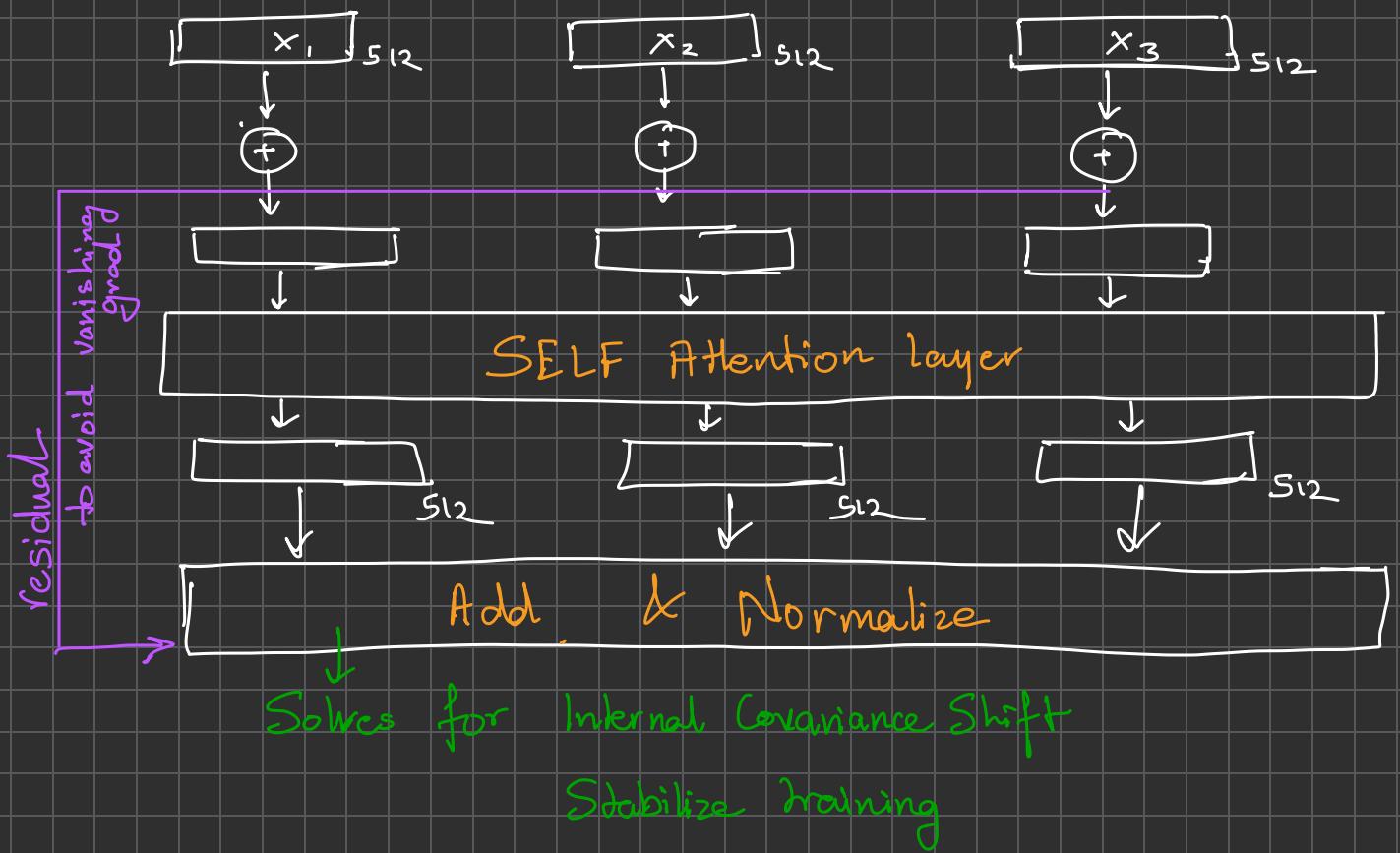
$$\underbrace{\text{split them by } d_{\text{model}}}_{d_{\text{model}}} \rightarrow d_K = d_m / h \quad \# \text{heads}$$

$$\text{Attention } (K, Q, V) = \text{softmax} \left(\frac{K_i^T Q_i}{\sqrt{d_{K_i}}} \right) V_i$$

$$\text{Seq. } \left\{ \begin{matrix} H_1 \\ E \\ A \\ D \\ \vdots \\ H_n \\ E \\ A \\ D \\ \vdots \\ H_n \end{matrix} \right\} \xrightarrow{\substack{d_v \\ \downarrow \\ d_{\text{model}}}} \begin{matrix} H \\ \times \\ W^0 \end{matrix} \quad \begin{matrix} \text{Seq, } h \times d_v \\ \text{or} \\ \text{Seq, } d_m \end{matrix}$$

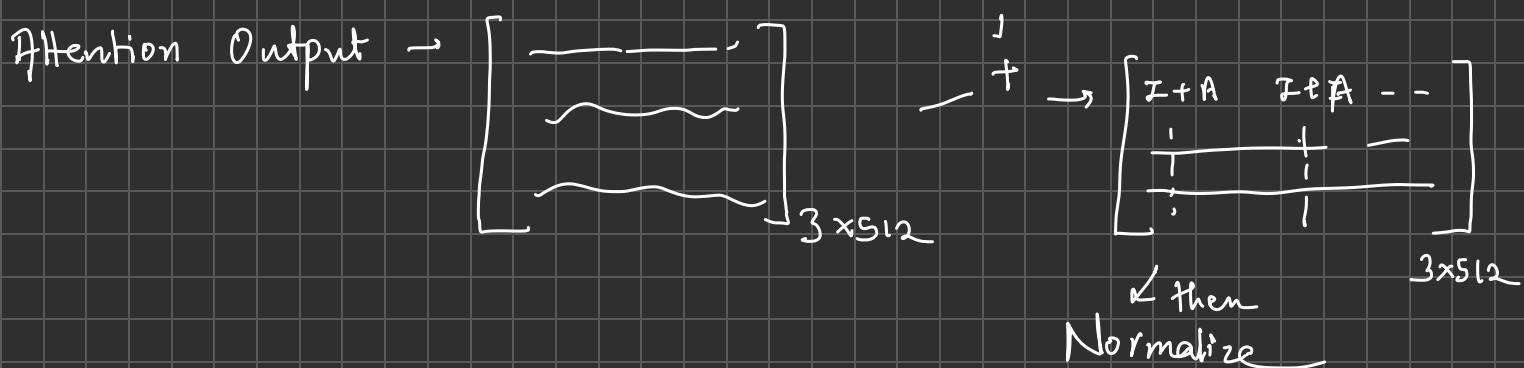
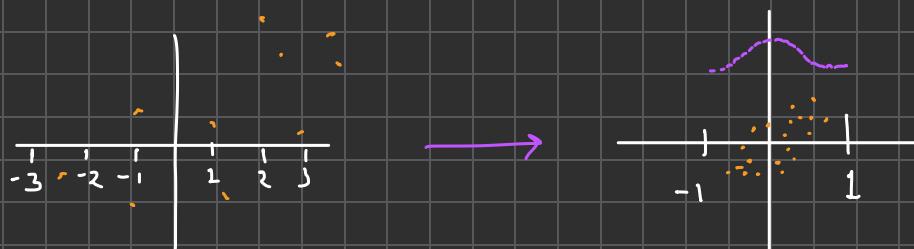
Concat (H_1, H_2, H_3, \dots)

Multihead Attention



$$X' = \gamma [W_i^T x + b_i] \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Normalization}$$

$$Y = \gamma \left[\frac{x' - \mu}{\sigma} \right] + \beta \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{across feature dims.}$$



Decoder

generates → <Start> → <sos>

