

Matrix Calculus:-

$y = f(x) \rightarrow x, y$ can be scalars, vectors or matrices

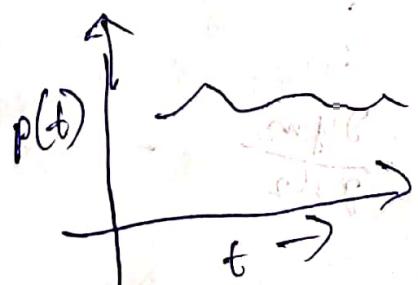
Let y be a scalar function of 3 independent variables.

$$\text{i.e. } y = f(x, y, z) \cdot \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x^T \\ y^T \\ z^T \end{bmatrix} = p$$

$$\nabla f = \frac{\partial f}{\partial x} \hat{x} + \frac{\partial f}{\partial y} \hat{y} + \frac{\partial f}{\partial z} \hat{z} = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} & \frac{\partial f}{\partial z} \end{bmatrix}^T$$

Derivatives with vectors.

$\vec{y} \in \mathbb{R}^m$ w.r.t scalar x , $\vec{y} = [y_1, y_2, \dots, y_m]^T$



$$p(t) = \begin{bmatrix} x(t) \\ y(t) \end{bmatrix}$$

$$\frac{dp(t)}{dt} = \text{Velocity}$$

$$\frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_n}{\partial x} \end{bmatrix}$$

2) Diff. scalar by a vector $y \in \mathbb{R}$, $x \in \mathbb{R}$.

$$\text{So, } x = [x_1, x_2, \dots, x_n]^T$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix}$$

3) Both x and y are vectors $y \in \mathbb{R}^m$, $x \in \mathbb{R}^n$

$$y = [y_1, y_2, \dots, y_m]^T$$

$$x = [x_1, x_2, \dots, x_n]^T$$

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

Jacobian Matrix

Derivatives with Matrix -

case 1:-

Matrix by scalar

$$\frac{\partial y}{\partial x} = \left(\frac{\partial y_{11}}{\partial x}, \frac{\partial y_{12}}{\partial x}, \dots, \frac{\partial y_{1n}}{\partial x} \right)^T + \left(\frac{\partial y_{21}}{\partial x}, \frac{\partial y_{22}}{\partial x}, \dots, \frac{\partial y_{2n}}{\partial x} \right)^T + \dots + \left(\frac{\partial y_{m1}}{\partial x}, \frac{\partial y_{m2}}{\partial x}, \dots, \frac{\partial y_{mn}}{\partial x} \right)^T$$

case 2:-

$x \in \mathbb{R}^{m \times n}$, y is a scalar

$$\frac{\partial y}{\partial x} = \left(\frac{\partial y}{\partial x_{11}}, \frac{\partial y}{\partial x_{12}}, \dots, \frac{\partial y}{\partial x_{1n}} \right)^T + \left(\frac{\partial y}{\partial x_{21}}, \frac{\partial y}{\partial x_{22}}, \dots, \frac{\partial y}{\partial x_{2n}} \right)^T + \dots + \left(\frac{\partial y}{\partial x_{m1}}, \frac{\partial y}{\partial x_{m2}}, \dots, \frac{\partial y}{\partial x_{mn}} \right)^T$$

$$\nabla y = Ax$$

$$\frac{\partial y}{\partial z} = A \frac{\partial x}{\partial z}$$

$$2) \lambda = y^T x = \left(x_1 y_1 + x_2 y_2 + \dots + x_n y_n \right)$$

$$\frac{\partial \lambda}{\partial z} = \left(x_1 \frac{\partial y_1}{\partial z} + y_1 \frac{\partial x_1}{\partial z} \right) + \left(x_2 \frac{\partial y_2}{\partial z} + y_2 \frac{\partial x_2}{\partial z} \right) + \dots$$

$$= \left(x_1 \frac{\partial y_1}{\partial z} + x_2 \frac{\partial y_2}{\partial z} + x_3 \frac{\partial y_3}{\partial z} + \dots \right) +$$

$$\left(y_1 \frac{\partial x_1}{\partial z} + y_2 \frac{\partial x_2}{\partial z} + \dots \right)$$

$$= x^T \frac{\partial y}{\partial z} + y^T \frac{\partial x}{\partial z}$$

Gradient Descent

Loss function $J(f(x^{(l)}, \theta), y)$

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} J$$

$$J = \frac{1}{N} \sum_{i=1}^N J(f(x^{(i)}, \theta), y^{(i)})$$

Stochastic gradient Descent :-

Select one instance randomly.

$$J = J(f(x^{(r)}, \theta), y^{(r)})$$

Mini Batch $\rightarrow \{x^{(1)}, \dots, x^{(m)}\}$

$$J = \frac{1}{m} \sum_{i=1}^m J(f(x^{(i)}, \theta), y^{(i)})$$

$$\rightarrow \epsilon_k = (1-\lambda) \epsilon_0 + \lambda \epsilon_T$$

$$\lambda = \frac{k}{T}$$

Deep Learning

9/3/23

Gradient Descent

sufficient condition on the learning rate for the convergence of the gradient descent or SGD.

$$\rightarrow \sum_{k=1}^{\infty} \epsilon_k = \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \epsilon_k^2 < \infty$$

common practice :- use a decaying value of ϵ_k .

$$\boxed{\epsilon_k = (1-\alpha) \epsilon_0 + \alpha \epsilon_t} \quad \alpha = \frac{k}{T}$$

$\epsilon_0 < c \epsilon_0$ (typically $0.01 \epsilon_0$)

Consider a scalar function,

$$\phi = \phi(x, y, z)$$

$$\nabla = \hat{i} \frac{\partial}{\partial x} + \hat{j} \frac{\partial}{\partial y} + \hat{k} \frac{\partial}{\partial z}$$

$$\nabla \phi = \hat{i} \frac{\partial \phi}{\partial x} + \hat{j} \frac{\partial \phi}{\partial y} + \hat{k} \frac{\partial \phi}{\partial z}$$

Let, $\vec{OP} = \vec{r}$

$$\vec{OQ} = \vec{r} + d\vec{r}$$

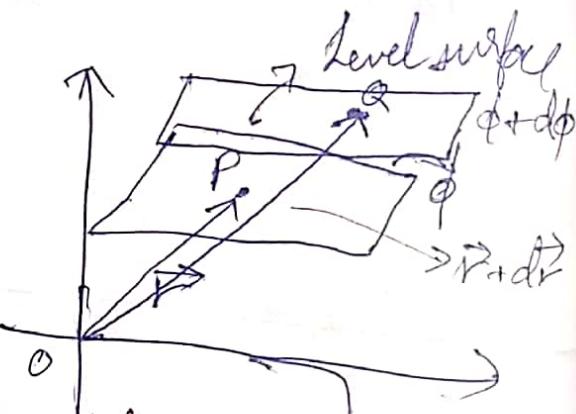
then shifting \vec{r} to $\vec{r} + d\vec{r}$

Shifting the level surface from ϕ to $\phi + d\phi$.

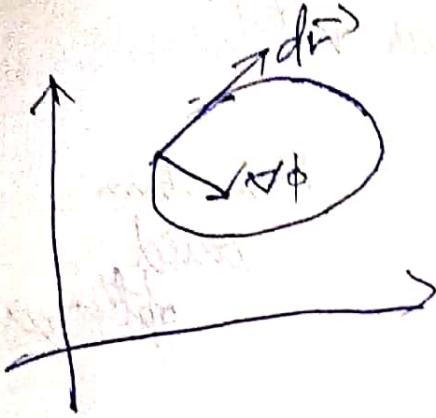
$$d\phi = \frac{\partial \phi}{\partial x} dx + \frac{\partial \phi}{\partial y} dy + \frac{\partial \phi}{\partial z} dz$$

$$= \left(\hat{i} \frac{\partial \phi}{\partial x} + \hat{j} \frac{\partial \phi}{\partial y} + \hat{k} \frac{\partial \phi}{\partial z} \right) \cdot (\hat{i} dx + \hat{j} dy + \hat{k} dz)$$

$$= \nabla \phi \cdot d\vec{r}$$



If we change P to Q on the same surface ϕ . Then $d\phi = 0$. Then $\nabla \phi \cdot d\vec{r} = 0$



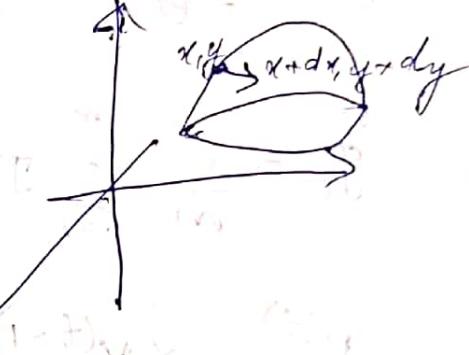
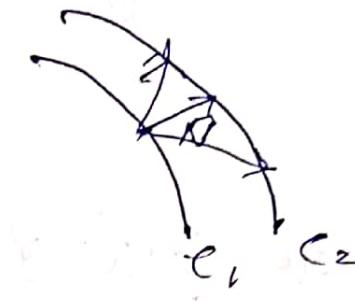
$\nabla\phi$ points in a direction normal to the level surface

$$\vec{\nabla}\phi = |\vec{\nabla}\phi| \cdot \hat{N}$$

$$d\phi = \vec{\nabla}\phi \cdot d\vec{r} = |\vec{\nabla}\phi| \cdot |d\vec{r}| \cdot \cos\theta$$

$$= |\vec{\nabla}\phi| d\vec{N}$$

$$\frac{d\phi}{d\vec{N}} = |\vec{\nabla}\phi|$$

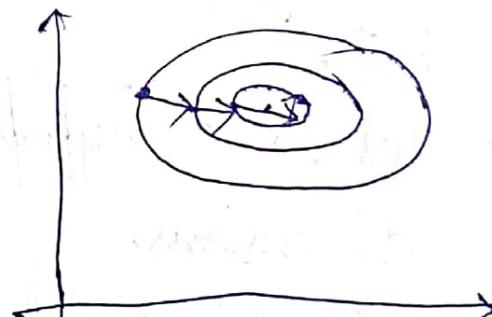


Challenges of Gradient Descent:-

- Choice of the proper learning rate or step size.
- Predefined learning rate for all directions.
- Problem of saddle points

10/03/2023

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta g_t$$



→ We assume a momentum force given as $v^{(t)}$

$$v^{(t)} \leftarrow \alpha v^{(t-1)} - \eta g_t$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + v^{(t)}$$

Momentum
based
optimizer

$\theta^{(0)}$ ← Randomly initialized parameters.

$$v^{(0)} = \theta$$

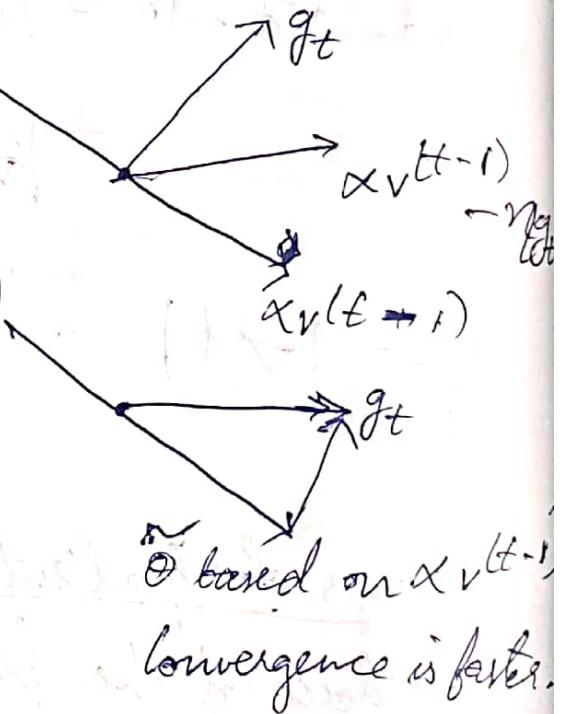
→ Nesterov Acceleration of Gradient :-

$$\tilde{\theta}^{(t)} = \theta^{(t-1)} + \alpha v^{(t-1)}$$

$$\tilde{g}_t = \frac{1}{m} \sum_{i=1}^m \nabla f(x^{(i)}, \tilde{\theta}), y^{(i)})$$

$$v^{(t)} = \alpha v^{(t-1)} - \eta \tilde{g}_t$$

$$\theta^{(t+1)} = \tilde{\theta}^{(t)} + v^{(t)}$$



Problem with momentum based optimizer :-

- The hyperparameters α & η , determine the learning rate.
- Use same hyperparameters values across all dimensions.
- Better if tuned according to gradient.

Sohm :- Adaptive hyperparameter settings.

- Adagrad (Technique)

→ Adagrad :-

$$g_t = \frac{1}{m} \sum \nabla_{\theta} J(f(x, \theta^{(t)}), y)$$

~~$$s_t = \sum_{\tau=1}^t g_{\tau} \cdot g_{\tau}^T$$~~

~~$$\theta^{(t+1)} = \theta^{(t)} - \frac{\eta}{\sqrt{s_t + r_t}} g_t$$~~

$$g_t = \begin{bmatrix} \frac{g_1}{\sqrt{s_1 + r_1}} \\ \frac{g_2}{\sqrt{s_2 + r_2}} \\ \vdots \\ \frac{g_t}{\sqrt{s_t + r_t}} \end{bmatrix}$$

RMS Prop :-

Discard the history from the extreme past.

~~$$r_{t+1} = \beta r_t + (1-\beta) s_t$$~~

$$r_0 = 0.$$

~~$$s_t = \sum_{\tau=1}^t g_{\tau} \cdot g_{\tau}^T$$~~

$$r_1 = (1-\beta) s_0$$

$$r_2 = \beta(1-\beta) s_0 + (1-\beta) s_1$$

$$r_3 = \beta^2(1-\beta) s_0 + \beta(1-\beta) s_1 + (1-\beta) s_2$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\eta}{\sqrt{r_t + s_t}} g_t$$

This creates a bias based on value of β .

$$\hat{s}_t = \frac{s_t}{1-\beta} \quad \hat{r}_t = \frac{r_t}{1-\beta}$$

Adam (Adaptive Momentum Optimizer)

$$s_t = \beta_1 s_{t-1} + (1 - \beta_1) g_t \quad \hat{s}_t = \frac{s_t}{1 - \beta_1}$$

$$r_t = \beta_2 r_{t-1} + (1 - \beta_2) g_t \circ g_t \quad \hat{r}_t = \frac{r_t}{1 - \beta_2}$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\eta}{\sqrt{\epsilon I + \hat{r}_t}} \cdot \hat{s}_t$$

→ Batch Normalization

happens when the distribution of the input features vary across the mini batches.

internal covariate shift

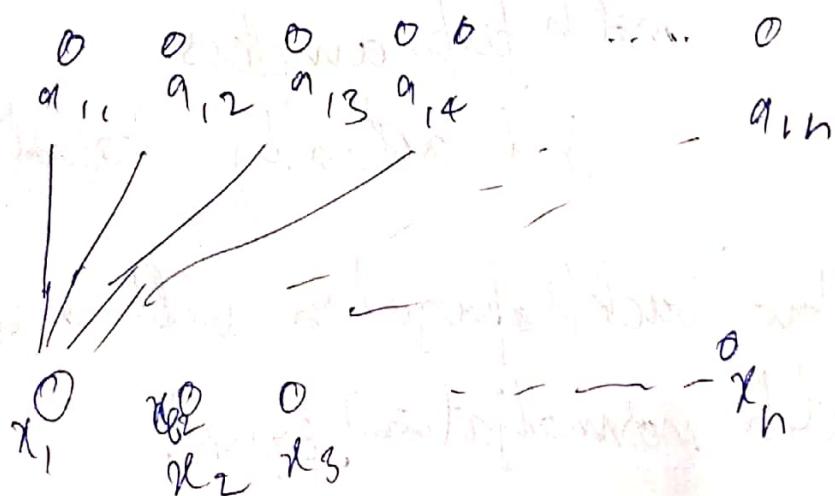
$$f_2((f_1(x_i, \theta_1), \theta_2))$$

$$f_2(u, \theta_2)$$

Soln: It would help if the pre-activations at each layer is normalized.

↳ many ways, they proposed gaussian normalization.

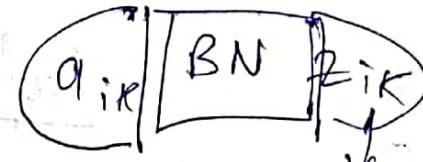
→ It is answered by standardizing the pre-activation at each unit.



$$\hat{a}_{ik} = a_{ik} - \mathbb{E}[a_{ik}]$$

$$\text{scale} \rightarrow \frac{1}{\text{std}} \cdot \hat{a}_{ik}$$

$$\text{SE}[a_{ik}]$$



some scaling and transformation

$$x_1, x_2, \dots, x_m$$

↓
no. of minibatches

→ Once batch normalization is done, re-normalization required (as scaling and transformation had been done).

$$y^{(k)} = \gamma^{(k)} a_{ik} + \beta^{(k)}$$

↓
doing rescaling
at least
larger

not hyperparameters

but actually learnt.

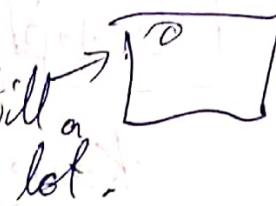
- **HW** → How backpropagation will happen for batch normalization? eqns?

↳ What if we do it at each layer?
As:- no effect of normalization would be visible.

→ CNN

- (i) images can be very large, so feed forward networks can't handle.
- (ii) images have higher dimensional ~~relations~~ relations in pixels, spatial dependence.
- (iii) transformations localization

convolution serves
better in this case.



still a
lot.

20/03/23

CNN properties

- 1) Locality
- 2) Stationarity
- 3) Hierarchy

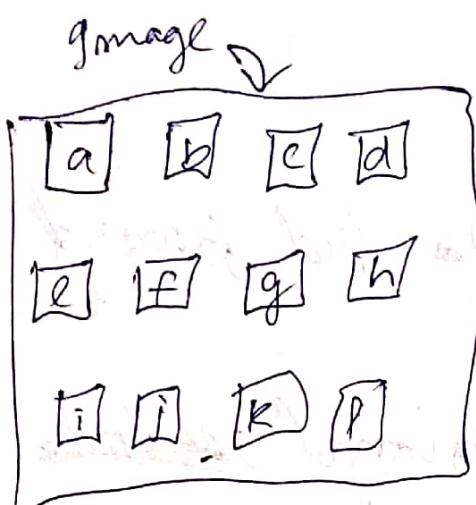
CNN tasks

① Convolution layer (A Predefined size matrix)

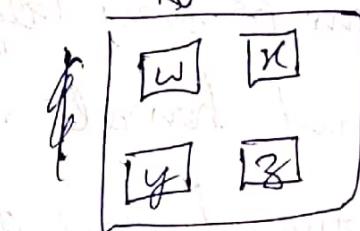
→ Kernel (Filters)

(Element-wise

multiplication)



Kernel



Resultant

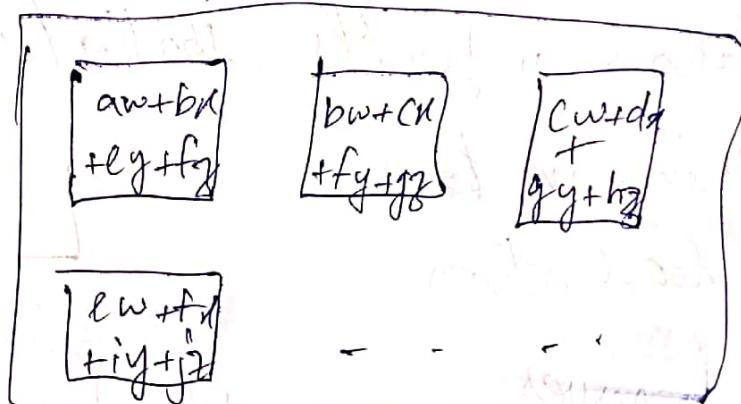
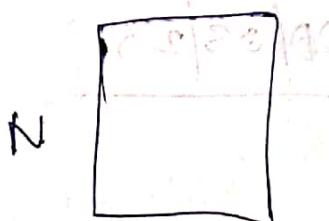
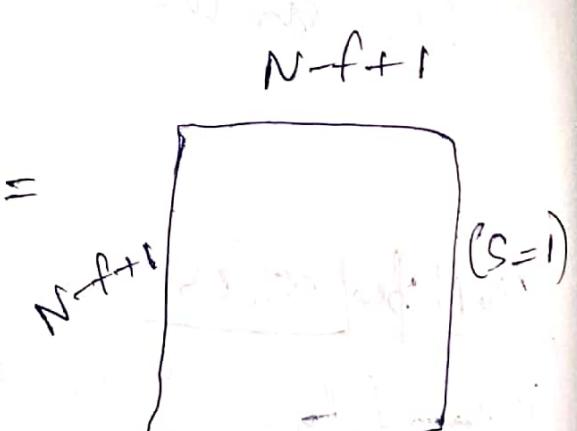


Image
N



Kernel
f

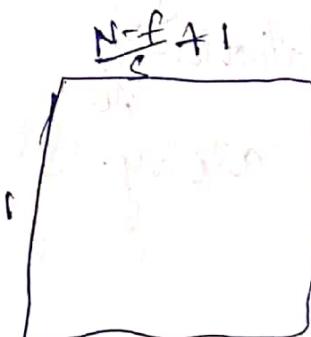
*



(S=1)

Stride = (S)

$\frac{N-f+1}{S}$



Stride is the length by which the kernel is shifted.

$$I = \begin{bmatrix} a \\ b \\ c \\ d \\ e \\ f \end{bmatrix} \quad F = \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix}$$

$$a w + b x + c y + d z$$

$$= I^T \cdot F = I \cdot F$$

→ Advantages:-

- * Sparse interaction
- * Parameter sharing

The no. of inputs in the lower layer, which affects an output unit in the higher layer is called the receptive field of the output unit.

CNN is also known as DNN, as the loss in receptive field is covered up by increasing the no. of layers in the network.

Also, the parameters are shared (which captures locality.)

Invariance :- (Output shouldn't change much).

eg:- Translational invariance, scaling invariance, Rotational invariance.

$$f(g(I)) = f(I) \text{ (Invariance) (The order}$$

$$f(g(I)) = g(f(I)) \text{ (equivariance)}$$

of the layers
but shouldn't
matter).

CNN layers are neither fully invariant nor fully equivariant.

24/03/23

Deep Learning

Alexnet Architecture

Input :- $227 \times 227 \times 3$

Conv 1 :- 96 filters of dimension $11 \times 11 \times 3$ stride = 4

$$\frac{N-f}{s} + 1 \geq SS$$

$$\frac{227-11}{4} + 1$$

$$= SS \times SS \times 96$$

(Output size)

Parameters

$$(11 \times 11 \times 3) \times 96 \approx 35K$$

Max-pool - 3×3 filters of stride 2

$$\frac{(SS-3)}{2} + 1 = 27 \times 27 \times 96$$

Normalization

Conv 2 :- 256 filters of dimension $3 \times 3 \times 96$,
stride = 1. (Padding = 2)

$$(27+4-5)+1 = 27 \times 27 \times 256$$

 ~~$(27+5) \times 27 = 12 \times 12 \times 256$~~

Parameters :-

$$(5 \times 5 \times 96) \times 256$$

Max-Pool-2

$$3 \times 3 \text{ filters of stride } = 2 \quad \frac{(27-3)}{2} + 1 = 13 \times 13 \times 256$$

Normalization-2

Conv-3

384 filters 3×3 , stride 1, pad 1.

$$(13+2-3) + 1 = 13 \times 13 \times 384$$

(output size)

Parameters:-

$$(3 \times 3 \times 256) \times 384$$

~~384~~ Conv-4

384 filters 3×3 , strides = 1, pad = 1.

Output size

$$13 \times 3 \times 384$$

Parameters:-

$$(3 \times 3 \times 384) \times 384$$

Layer-5 :- 256 filters: 3×3 , stride = 1, pad = 1

$13 \times 13 \times 256$ (output size)

Parameters :-

$(3 \times 3 \times 384) \times 256$

Max pool - 3 :- 3×3 filters, stride = 2

$\frac{(13-3)}{2} + 1 = 6 \times 6 \times 256$

FC - 1 :- 4096 neurons

(Output size) (Parameters)

4096 9216×4096

FC - 2 :- 4096 neurons

(Output size) (Parameters)

4096 4096×4096

FC - 3 :- 1000 neurons

(Parameters)

4096×1000

Problems :-

25/03/23

→ It was a bit shallow, one of the first architectures to use deep neural network.

→ Vgg Net

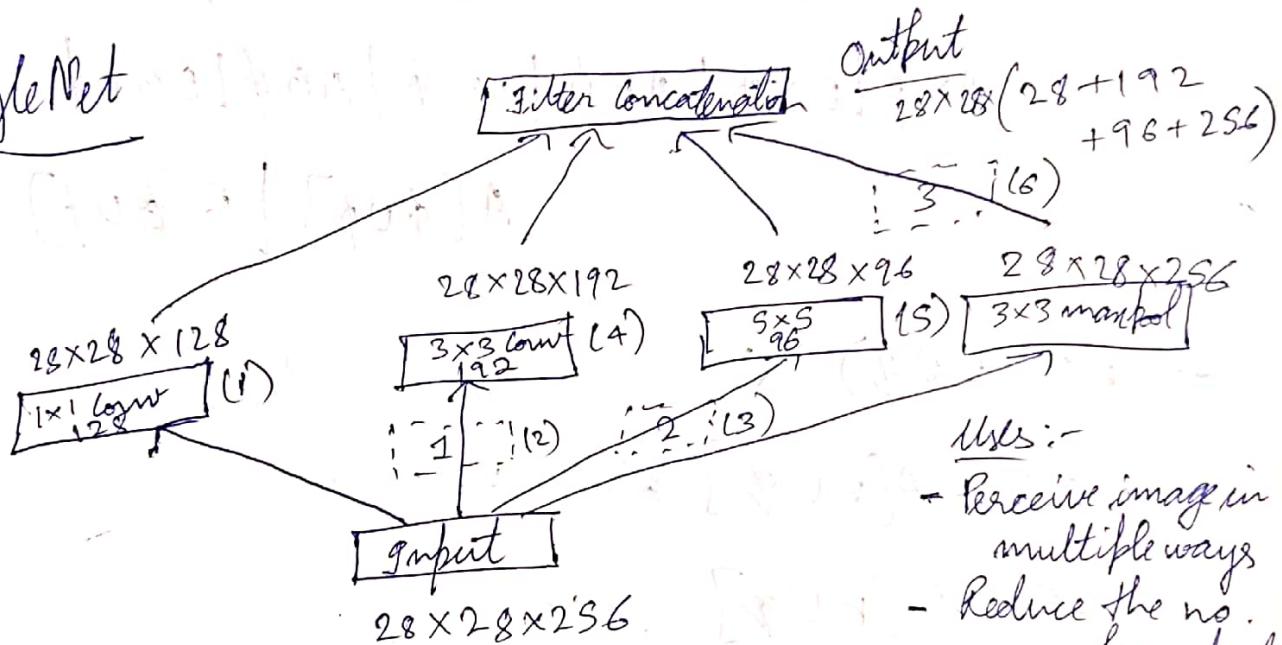
- 138 million parameters

- 19 layers

→ Complexity of the CNN's

Measured by the # Floating Point Operations (FLOPs)

→ GoogleNet



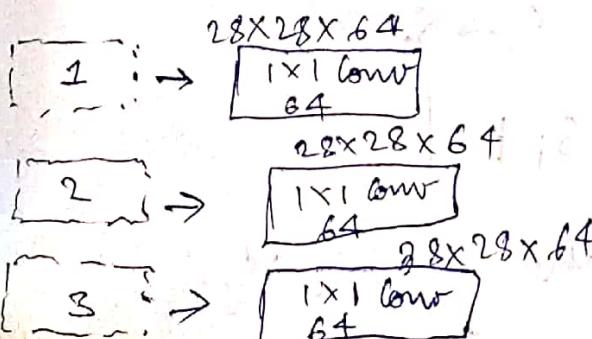
Consider the MULTS

$$1) 28 \times 28 \times 256 \times 1 \times 1 \times 128$$

$$2) 28 \times 28 \times 256 \times 3 \times 3 \times 192$$

$$3) 28 \times 28 \times 256 \times 5 \times 5 \times 96$$

$$\text{Total FLOPs} = 854 \text{ M}$$



- Uses:-
- Perceive image in multiple ways
 - Reduce the no. of computations

$$1) 28 \times 28 \times 256 \times 1 \times 1 \times 128$$

$$2) 28 \times 28 \times 256 \times 1 \times 1 \times 64$$

$$3) 28 \times 28 \times 256 \times 1 \times 1 \times 64$$

$$4) 28 \times 28 \times 64 \times 3 \times 3 \times 192$$

$$5) 28 \times 28 \times 64 \times 5 \times 5 \times 96$$

$$6) 28 \times 28 \times 256 \times 1 \times 1 \times 64$$

$$\text{Mults} = 358 \text{ M}$$

Deep Learning

27/03/23

CNN:- Parameters shared across grids

RNN:- Parameters shared across long sequences.

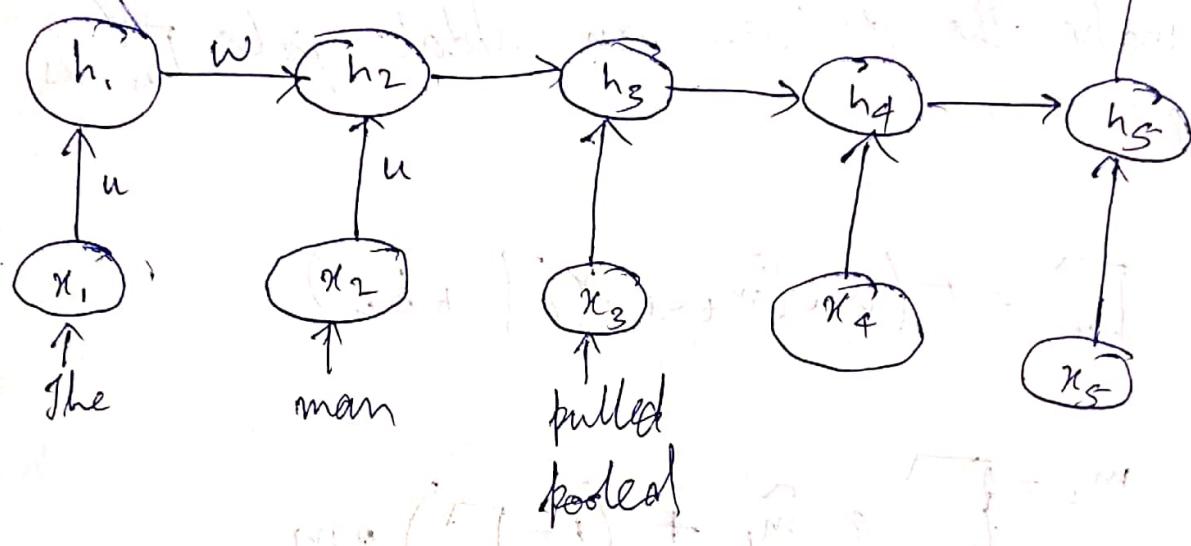
28/03/23

Deep learning

The man ~~pulled~~ the cart.

The man ~~pooled~~ the cart

$P(y_1, y_2, \dots, y_n)$ Prob score
 $\begin{bmatrix} P_C \\ P_I \end{bmatrix}$ \uparrow whether the sentence is correct or not.



$$h_t = \sigma_h (u \vec{x}_t + w h_{t-1} + b_h)$$

$$\hat{y} = \sigma_y (v h_t + b_y)$$

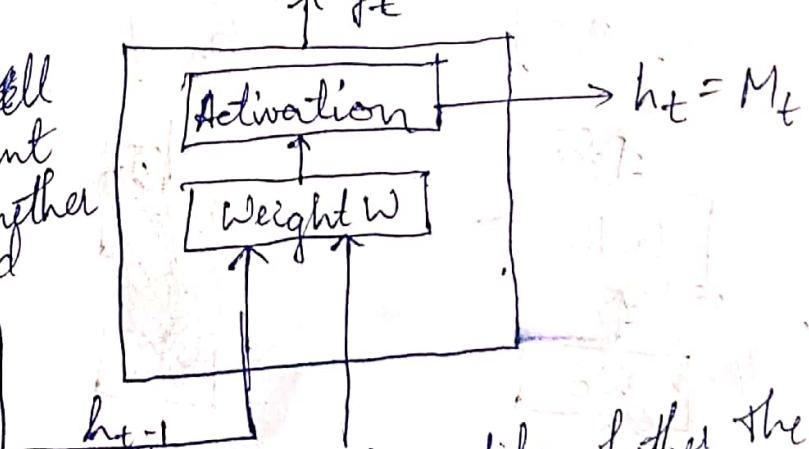
$$L = y \log \hat{y} + (1 - y) \log (1 - \hat{y})$$

GRU

Maintains a memory cell that contains relevant information like whether the noun ~~was~~ was singular/plural.

Maintains a memory cell that contains relevant information ($M_t = h_t$)

Vanilla RNN



x_t like whether the noun mentioned was singular/plural

iii) At every step, the GRU needs to decide whether m_t needs to be replaced by an alternate candidate value \tilde{m}_t

$$\tilde{m}_t = \tanh(W_c[m_{t-1}, x_t] + b_m)$$

iv) To make the decision, an update gate I_u is used

$$I_u = \sigma(W_u[m_{t-1}, x_t] + b_u)$$

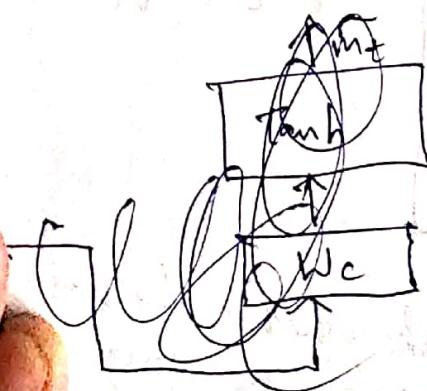
$$m_t = I_u \circ \tilde{m}_t + (1 - I_u) \circ m_{t-1}$$

$\circ \rightarrow$ hadamard product

$$\vec{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

$$\vec{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

$$\vec{a} \circ \vec{b} = \begin{bmatrix} a_1 \cdot b_1 \\ a_2 \cdot b_2 \\ \vdots \\ a_n \cdot b_n \end{bmatrix}$$



Attention-Mechanism

DL

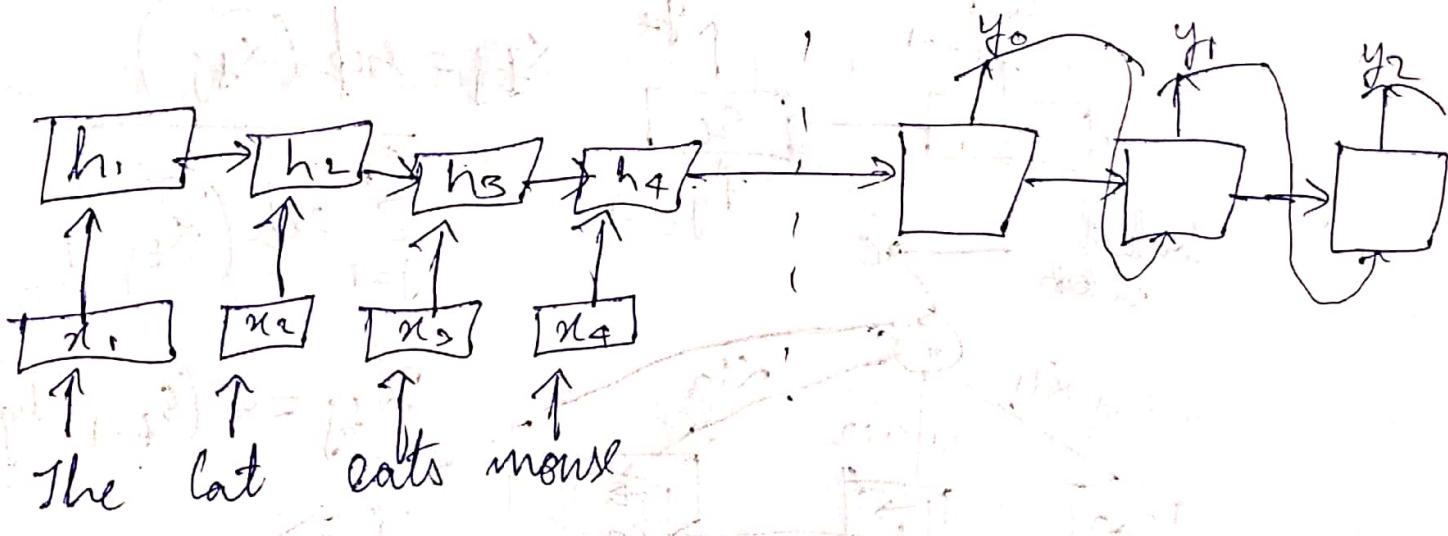
10/04/23

Neural machine translation by jointly learning to align and translate, Bahdanau et al ICLR 2015.

The cat eats mouse

Billi chunha khati hair

Akkatta Aklathinanat



RNN based encoder

$$x = (x_1, x_2, \dots, x_{T_n})$$

$$h_t = f(h_{t-1}, x_t)$$

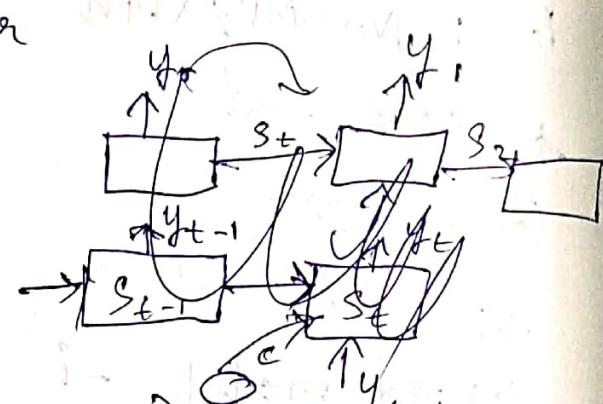
$$c = g(h_1, h_2, \dots, h_{T_n})$$

Encoder

Decoder

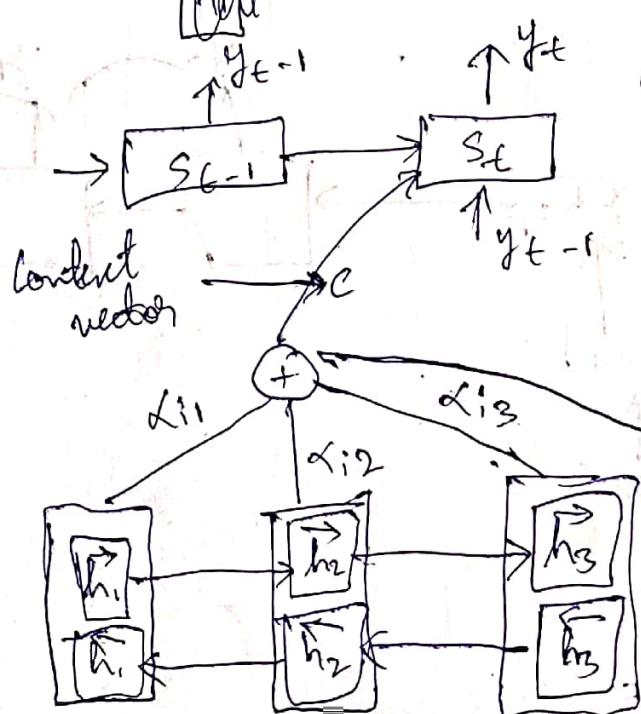
$$p(y) = \prod_{t=1}^T p(y_t)$$

$$p(y_t | d(y_1, y_2, \dots, y_{t-1}, y))$$



$$p(y_t | d(y_1, y_2, \dots, y_{t-1}, y)) = g(y_{t-1}, s_t, c)$$

$$c^{(t)} = \sum_{j=1}^{T_x} \alpha_{tj} h_j$$

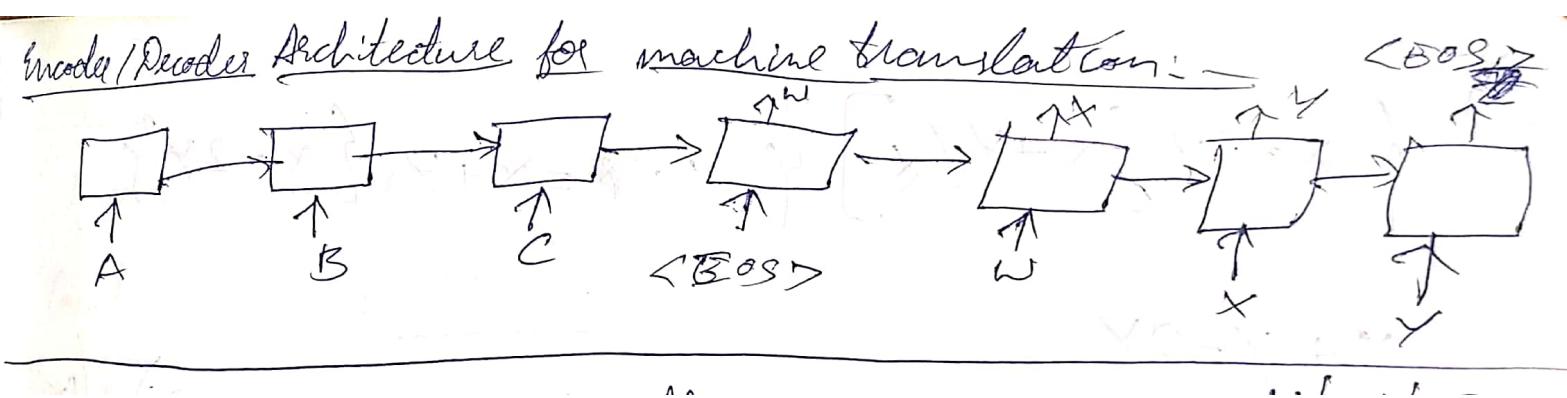


$$\alpha_{tj} = \exp(\tilde{\alpha}_{tj})$$

$$\tilde{\alpha}_{tj} = \sum_{j=1}^{T_x} \exp(\tilde{\alpha}_{tj})$$

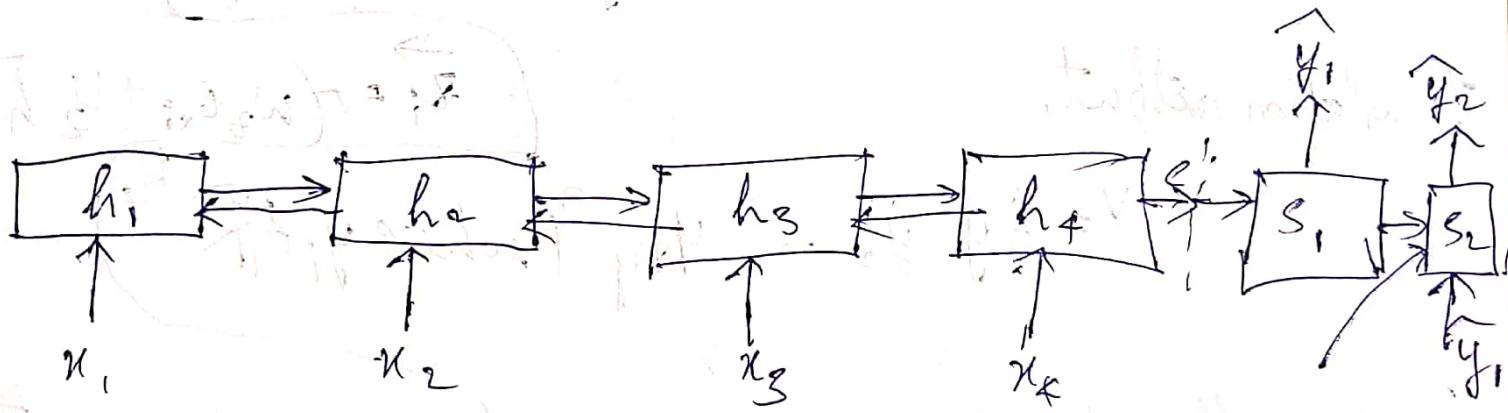
$$\tilde{\alpha}_{tj} = f(s_{t-1}, h_j)$$



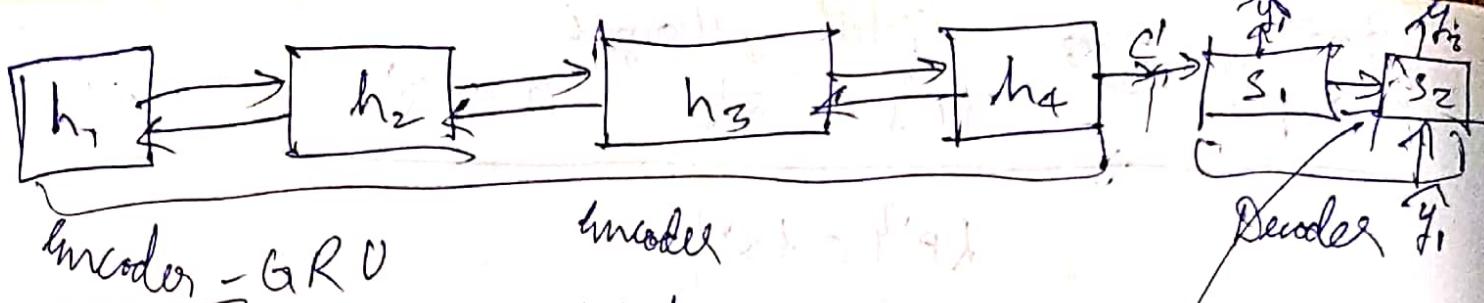


Deep Learning

13/04/23



Problem:- In case of blankary sentences, 'c' became the bottle neck.



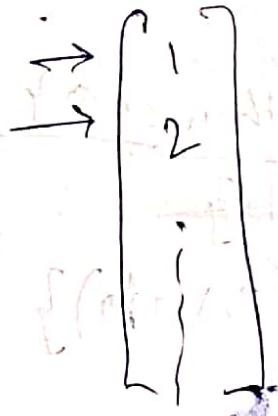
Encoder - GRU

source sentence - Word token

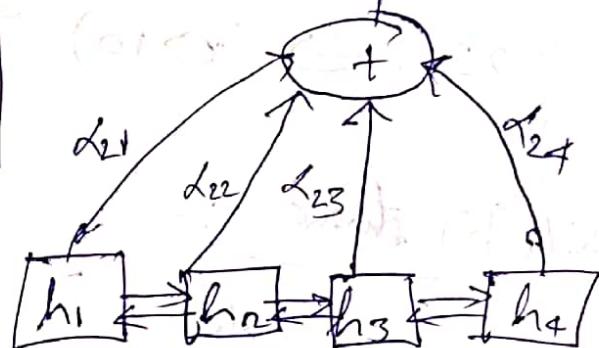
1 of K vectors

Construct a vocabulary of K_x

words



1 of $K \rightarrow$ one hot encoding
over K_x words



$$h_i = (1 - \vec{z}_i) \circ h_{i-1} + \vec{z}_i \vec{h}_i$$

↑
Update gate

$$h_i = \tanh(w_{2i} x_i + u_i \vec{r}_{i-1} \vec{h}_{i-1})$$

↑
Candidate cell
↑
Relevance gate

Input in $x = \{x_1, x_2, \dots, x_{T_H}\}$ where $x_i \in \mathbb{R}^{K_x}$

Translation output,

$$y = \{y_1, y_2, \dots, y_{T_Y}\} \text{ where } y_i \in \mathbb{R}^{K_y}$$

$$\vec{z}_i = \sigma(w_2 x_i + v_2 \vec{h}_{i-1})$$

Encoder :- forward states

$$\vec{r}_i = \sigma(w_r x_i + v_r \vec{h}_{i-1})$$

same equation for backward states

Decoder:-

$$S_i = (1 - z_i) \circ S_{i-1} + z_i \circ S_{i-1}^*$$

$$S_i = \tanh(w^T y_i + U [r_i \circ S_{i-1}] + C c_i)$$

candidate

$$c_i = \sum_{j=1}^{T_y} \alpha_{ij} h_j$$

$$z_i^* = \sigma(w^T y_{i-1} + U_2 S_{i-1} + C_2 c_i)$$

$$r_i^2 = \sigma(w^T y_{i-1} + U_r S_{i-1} + C_r c_i)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j=1}^{T_y} \exp(e_{ij})}$$

$$e_{ij} = v_a^T \tanh(w_a S_{i-1} + u_a h_j)$$

$$e_{ij} = v_a^T \tanh(w_a S_{i-1} + u_a h_j)$$

Prob of output

$$P(y_i | S_i, y_{i-1}, c_i) \propto \exp(y_i^T w_i)$$

$$t_i \sim u_s S_{i-1} + v_o E y_{i-1} + c_o c_i$$

Blue score :- Performance measure in machine translation

Consider 2 sentences:-

- The ball is blue (Predicted)
 - The ball has a blue colour (Target)
 - The ball the ball \rightarrow 100% overlap
 - The ball colour blue has
- Should not have high blue score.

Step 1

N-gram

1-gram "The", "ball", "is", "blue"

2-gram "The ball", "ball is", "is blue".

3-grams "The ball is", "ball is blue".

Step 2

Precision = # of correct predicted words

Total predicted words

Precision (Pred/Target) = $\frac{3}{4}$ (The, ball, blue)

Prob But there can be repetition, (The, The, The, The)

Solution

use clipped precision instead of precision, or
use multiple target sentences.

→ to find the blue score, find the precision for

(P_1) 1-gram, (P_2) 2-gram, (P_3) 3-gram, ..., n -gram.

(Higher the value of n , better the ~~blue score~~ result)

Significance

→ find geometric average precision score

$$= \exp \left(\frac{1}{N} \sum_{n=1}^N w_n \log p_n \right)$$

$$= \exp \left(\frac{1}{N} \sum_{n=1}^N \log p_n^{1/N} \right)$$

$$= \exp \left(\log \left(p_1^{1/N} \cdot p_2^{1/N} \cdots p_N^{1/N} \right) \right)$$

$$= p_1^{1/N} \cdot p_2^{1/N} \cdot p_3^{1/N} \cdots p_N^{1/N}$$

$$\text{for } (N=3) \quad \text{GAP} = p_1^{1/3} \cdot p_2^{1/3} \cdot p_3^{1/3}$$

Problem
High probab.
for smaller

→ Find brevity penalty
(penalize shorter sentences)

$$= \begin{cases} 1, & \text{length} > r \\ e^{(r - \text{length})}, & \text{length} \leq r \end{cases}$$

→ predicted length

Black score = GA P. beauty penalty