

Matrix calculusdependent \leftarrow independent

$$Y = f(X)$$

\hookrightarrow can be scalar/vector/matrix

Let $y = f(x, y, z)$

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

① Derivative of a vector wrt scalar,

let $y \in \mathbb{R}^m$ wrt scalar $x \in \mathbb{R}$

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

$$\frac{\partial y}{\partial x} = \left[\frac{\partial y_1}{\partial x}, \frac{\partial y_2}{\partial x}, \dots, \frac{\partial y_m}{\partial x} \right]^T$$

② Derivative of a scalar wrt vector*

scalar $y \in \mathbb{R}$ wrt vector $x \in \mathbb{R}^m$

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

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

③ Derivative of a vector w.r.t vector *

Let $\mathbf{y} \in \mathbb{R}^m$ & $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{y} = [y_1, y_2, \dots, y_m]^T \quad \mathbf{x} = [x_1, x_2, \dots, x_n]^T$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \text{(Jacobian matrix)} \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

④ Derivative of a matrix by a scalar

Let $\mathbf{Y} \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}$ is scalar

$$\mathbf{Y} = \begin{bmatrix} Y_{11} & Y_{12} & \dots & Y_{1n} \\ Y_{21} & Y_{22} & \dots & Y_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{m1} & Y_{m2} & \dots & Y_{mn} \end{bmatrix} = \{Y_{ij}\}_{m \times n}$$

$$\frac{\partial \mathbf{Y}}{\partial x} = \begin{bmatrix} \frac{\partial Y_{11}}{\partial x} & \frac{\partial Y_{12}}{\partial x} & \dots & \frac{\partial Y_{1n}}{\partial x} \\ \frac{\partial Y_{21}}{\partial x} & \frac{\partial Y_{22}}{\partial x} & \dots & \frac{\partial Y_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial Y_{m1}}{\partial x} & \frac{\partial Y_{m2}}{\partial x} & \dots & \frac{\partial Y_{mn}}{\partial x} \end{bmatrix}$$

* Numerator layout: Follow the dimension of numerator. If denominator has a dimension, the overall derivative is transposed, with dim. of denominator.

(5) Derivative of a scalar w.r.t matrix*

$y \in \mathbb{R}$ is scalar & $X \in \mathbb{R}^{m \times n}$ is a matrix

$$\frac{\partial y}{\partial X} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{21}} & \dots & \frac{\partial y}{\partial x_{m1}} \\ \frac{\partial y}{\partial x_{12}} & \frac{\partial y}{\partial x_{22}} & \dots & \frac{\partial y}{\partial x_{m2}} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y}{\partial x_{1n}} & \frac{\partial y}{\partial x_{2n}} & \dots & \frac{\partial y}{\partial x_{mn}} \end{bmatrix}$$

Identities

$$① \quad y = Ax \Rightarrow \frac{\partial y}{\partial z} = A \frac{\partial x}{\partial z}$$

$$② \quad \alpha = y^T x = [x_1 y_1 + x_2 y_2 + \dots + x_n y_n]$$

$$\begin{aligned} \frac{\partial \alpha}{\partial z} &= \left(x_1 \frac{\partial y_1}{\partial z_1} + y_1 \frac{\partial x_1}{\partial z_1} \right) + \left(x_2 \frac{\partial y_2}{\partial z_2} + y_2 \frac{\partial x_2}{\partial z_2} \right) + \dots \\ &= \left(x_1 \frac{\partial y_1}{\partial z_1} + x_2 \frac{\partial y_2}{\partial z_2} + \dots \right) + \left(y_1 \frac{\partial x_1}{\partial z_1} + y_2 \frac{\partial x_2}{\partial z_2} + \dots \right) \end{aligned}$$

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

Gradient Descent

Loss funⁿ $J(f(x^{(i)}, \theta), y^{(i)})$

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

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

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

Cons - computation time ↑↑

↳ for each update, all points in dataset are to be considered & corresponding losses are added

↓
Solⁿ

Stochastic GD

↳ select 1 instance randomly & use it for updating θ , let $x^{(r)}$ be the random instance

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

Cons

↳ at each iteration, L is highly fluctuating leading to slow convergence

Solⁿ → Mini-Batching

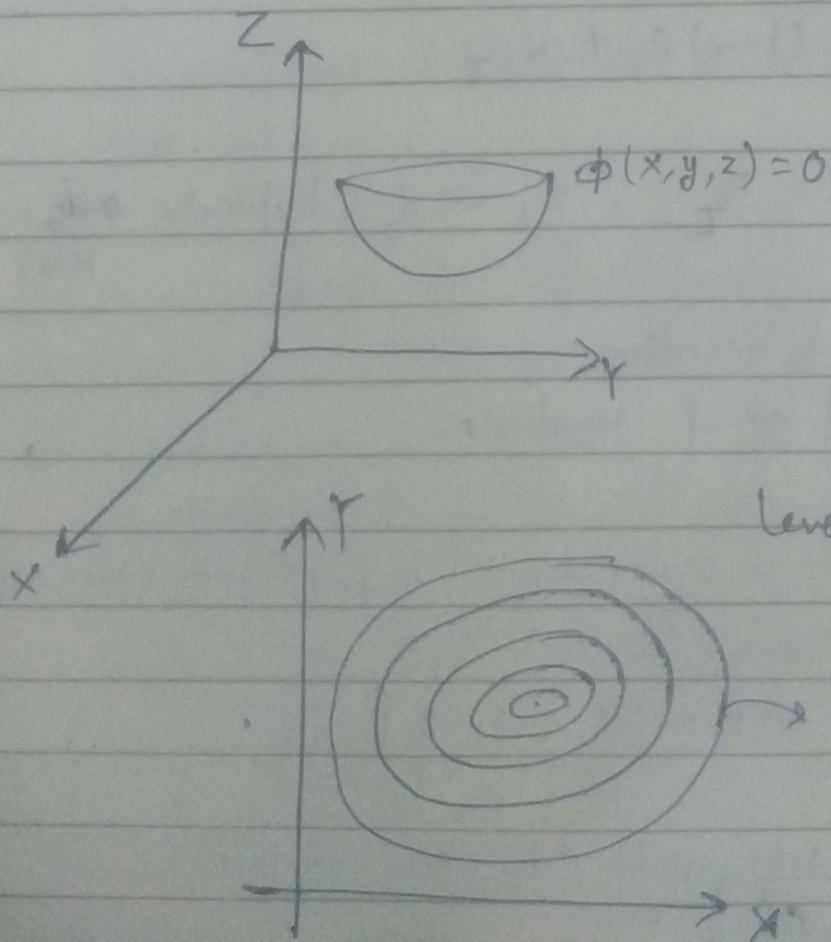
↳ select m instances (sample)
instead of whole N

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

Mini-batch = $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

ε -bottleneck

↳ $\varepsilon \uparrow$ or const. overshooting the minima
 $\varepsilon \downarrow$ slow convergence



$$\varepsilon_k = (1-\alpha)\varepsilon_0 + \alpha\varepsilon_T ; \quad \alpha = \frac{k}{T}$$

DL - 09/03.

Mini-Batch GD

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

$$\theta^{(t+1)} = \theta^{(t)} - \varepsilon \hat{g}_t$$

on ε

Sufficient cond'n \uparrow for the convergence of stochastic GD,

$$\left\{ \begin{array}{l} \sum_{k=1}^{\infty} \varepsilon_k = \infty \quad \& \quad \sum_{k=1}^{\infty} \varepsilon_k^2 < \infty \end{array} \right. \rightarrow \varepsilon_k \text{ should be fractional ; } k \text{ is arbitrary}$$

Common practice, — use an adaptive (decaying) value of ε

$$\varepsilon_k = (1-\alpha) \varepsilon_0 + \alpha \varepsilon_T$$

$$\text{where } \alpha = \frac{k}{T}, \quad \varepsilon_T \ll \varepsilon_0 \text{ (typically 0.001 or less)}$$

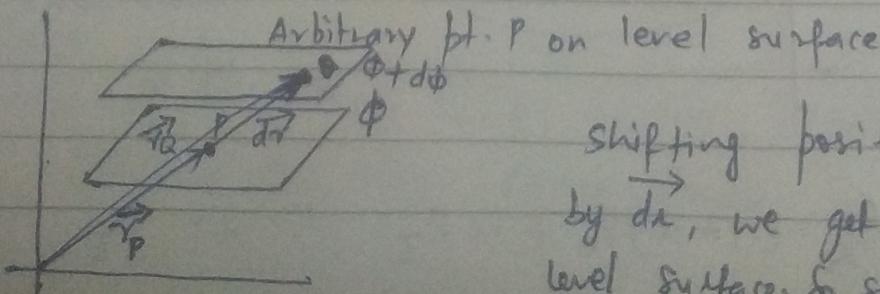
k = no of iteration

T = total no of iterations

Let $\phi(x, y, z)$ be a scalar funcn (a funcn with scalar output)

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

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



Shifting position of P by \vec{dr} , we get Q in another level surface. So surface is $\phi + d\phi$

Let $\vec{OP} = \vec{r}$, $\vec{OQ} = \vec{r} + \vec{dr}$

Then shifting \vec{r} to $\vec{r} + \vec{dr}$ shifts the level surface from ϕ to $d\phi$.

$d\phi$ = change in level surface

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

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

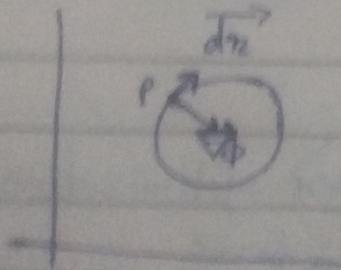
$$\Rightarrow d\phi = \vec{\nabla}\phi \cdot \vec{dr}$$

If we change P to Q s.t Q is on same LS as P,
then $d\phi = 0$.

$$d\phi = 0 \Rightarrow \vec{\nabla}\phi \cdot \vec{dr} = 0$$

$\Rightarrow \vec{\nabla}\phi$ & \vec{dr} are lar to each other

Thus, gradient is lar to \vec{dr} (tangent to LS at P)



$\vec{\nabla}\phi$ points in a direction normal to the LS.

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

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

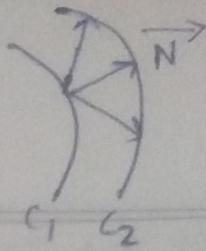
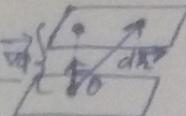
↳ unit vector
normal to

$\theta = \text{angle b/w } \vec{\nabla\phi} \text{ & } \vec{d\pi}$

(normal)

min^m dist. b/w = $d\vec{N} = \underline{d\vec{\pi} \cos \theta}$

the 2 L.S



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

$$\Rightarrow d\phi = |\vec{\nabla\phi}| \cdot |\vec{d\pi}| \cdot \cos \theta^*$$

$$\Rightarrow d\phi = |\vec{\nabla\phi}| \cdot |\vec{dN}| \Rightarrow \boxed{\frac{d\phi}{|\vec{dN}|} = |\vec{\nabla\phi}|}$$

∴

\vec{dN} is the normal vector

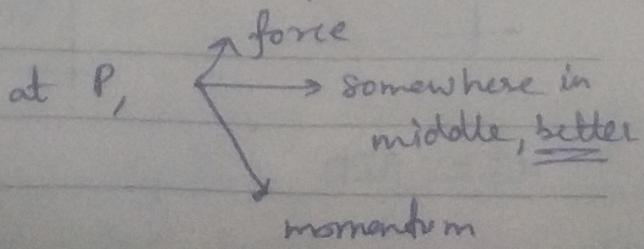
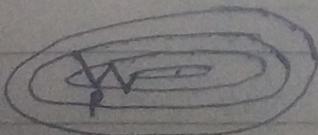
Gradient is the direction of max^m change, i.e., the direction in which ϕ changes most rapidly.

For minimization, one must come in direction opposite to gradient for steepest descent. That's why we use $-\hat{g}_t$ in GD, with ϵ for its own reasons.

Challenges to QD

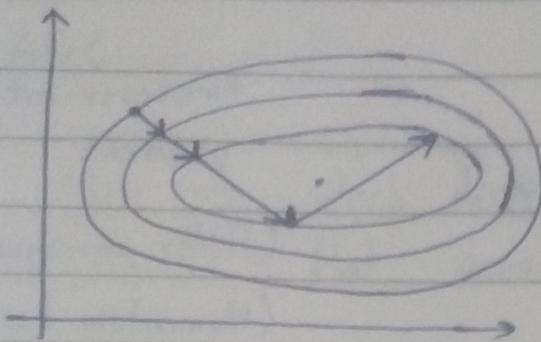
- choice of proper ϵ
 - ↓ ϵ slow convergence, many iterations
 - ↑ ϵ overshooting minima, oscillations
- predefined ϵ for all directions
- saddle points - slope (+ve) in 1 direction & (-ve) in 1 direction

Handling oscillations → momentum-based optimization techniques



* In DL, random initialisation is seldom. Pre-trained set of parameters is used instead. However, in theory (here), it may be considered random.

DL-10/03



Challenge - oscillating convergence

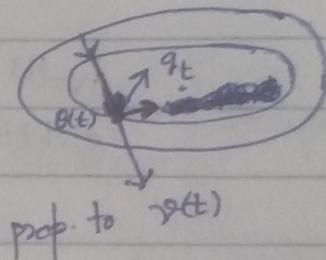
$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \frac{\eta_t}{\epsilon} g_t$$

Solⁿ - momentum-based technique

Assume that there is a momentum given by $v(t)$,

$$v^{(t)} \leftarrow \alpha v^{(t-1)} - \frac{\eta_t}{\epsilon} g_t$$

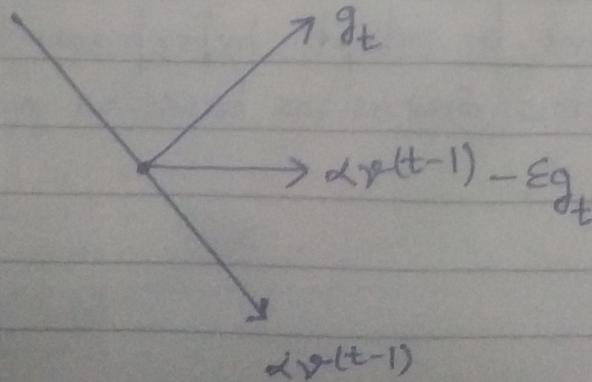
$$\theta^{(t+1)} \leftarrow \theta^{(t)} + \alpha v^{(t)}$$



Momentum-based optimizers handle oscillations better.

$$\theta^{(0)} \leftarrow \text{randomly initialized parameter}^*$$

Nesterov accⁿ gradient



Nesterov accⁿ

gradient

(Faster convergence)

(Momentum based)

(intermediate
set of parameters,
 $\tilde{\theta}$ & \tilde{g}_t)

$\tilde{\theta}$ based on $\alpha v^{(t-1)}$

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

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

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

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

Challenges to momentum-based optimizers

- ↳ hyperparameters α & ϵ determine learning rate
- ↳ uses same values of α & ϵ across all dimensions
 - ↳ better if tuned w.r.t gradient

Solⁿ - adaptive hyperparameter settings

↳ AdaGrad

- ↳ tries to adapt hyperparameter values across dimensions based on gradient

* element-wise product (\circ)
let $g_t \in \mathbb{R}^m$ then $g_t \circ g_t = \begin{bmatrix} g_1^2 \\ g_2^2 \\ \vdots \\ g_m^2 \end{bmatrix}$ & so on.

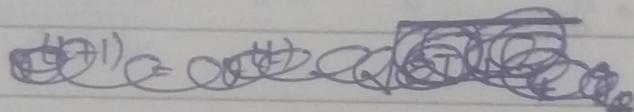
AdaGrad

Idea → steep descent \Rightarrow ~~high~~ ^{low} values for ~~grad~~ ^{hyperplane}
flat descent \Rightarrow high " "

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

$$\gamma_t = \sum_{T=1}^t g_T \circ g_T^*$$

$\theta^{(t+1)} = \theta^{(t)} - \eta g_t$ is now changed to,

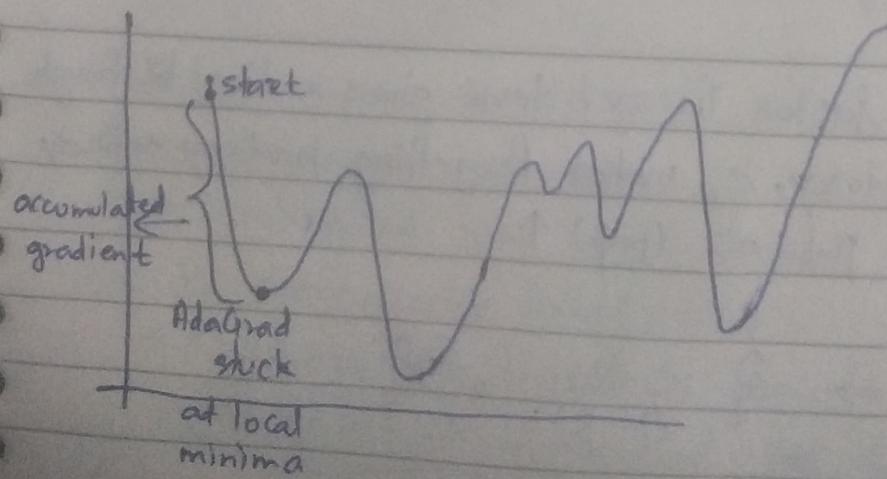


$$g_t = \begin{bmatrix} g_1^{(1)} \\ g_2^{(2)} \\ \vdots \\ g_m^{(m)} \end{bmatrix} \text{ changes to}$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\eta w_t}{\sqrt{\epsilon I + \gamma_t}}$$

$$g_t = \begin{bmatrix} g_1^{(1)} / \sqrt{\epsilon_1 + \gamma_1} \\ g_2^{(2)} / \sqrt{\epsilon_2 + \gamma_2} \\ \vdots \\ g_m^{(m)} / \sqrt{\epsilon_m + \gamma_m} \end{bmatrix}$$

Challenge → gets stuck in local minima
as it accumulates gradient over time



Improvement - discard the history from the extreme past



RMS Prop.

RMS Prop

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

$$s_{t+1} = \beta s_t + (1-\beta) s_t ; s_0 = 0$$

$$s_t = \sum_{\tau=1}^t g_\tau \circ g_\tau \quad \theta^{(t+1)} = \theta^{(t)} - \frac{\eta g_t}{\sqrt{s_t}}$$

$$s_1 = (1-\beta) s_0, \quad s_2 = \beta(1-\beta) s_0 + (1-\beta) s_1$$

$$s_3 = \beta^2(1-\beta) s_0 + \beta(1-\beta) s_1 + (1-\beta) s_2$$

If $\beta \downarrow$ priority is given to the s_t values which are closer to current time, ie, history gets lessened over time. All game is in powers of β getting small & small as β is small.

The $(1-\beta)$ factor in each term gives a bias, though, to the history, ie, makes forgetting history entirely difficult. $\beta \downarrow \Rightarrow (1-\beta) \uparrow$

$$\text{Correction} \rightarrow \hat{s}_t = \frac{s_t}{(1-\beta)}, \quad \hat{s}_t = \frac{s_t}{(1-\beta)}$$

Idea behind Adam (Adaptive Moment) Technique

Adam optimizer

Corrections

$$s_t = \beta_1 s_{t-1} + (1-\beta_1) g_t \quad \left\{ \begin{array}{l} \hat{s}_t = \frac{s_t}{(1-\beta_1)} \\ \hat{s}_{t-1} = \frac{s_{t-1}}{(1-\beta_1)} \end{array} \right.$$

$$\hat{\pi}_t = \beta_2 \hat{\pi}_{t-1} + (1-\beta_2) g_t \circ g_t \quad \left\{ \begin{array}{l} \hat{\pi}_t = \frac{\pi_t}{(1-\beta_2)} \\ \hat{\pi}_{t-1} = \frac{\pi_{t-1}}{(1-\beta_2)} \end{array} \right.$$

~~Optimizer equations~~

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

DL - 16/03

Batch Normalization

↪ not an optimization technique, booster of optimization process indeed (like, catalyst)

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

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

Consider mini-batch $B = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

↪ mini-batch GD

↪ compromise b/w computation complexity (GD)
& convergence time (SGD)

↪ challenge - mini-batch may not be a representative of actual set of inputs

Internal Covariance shift

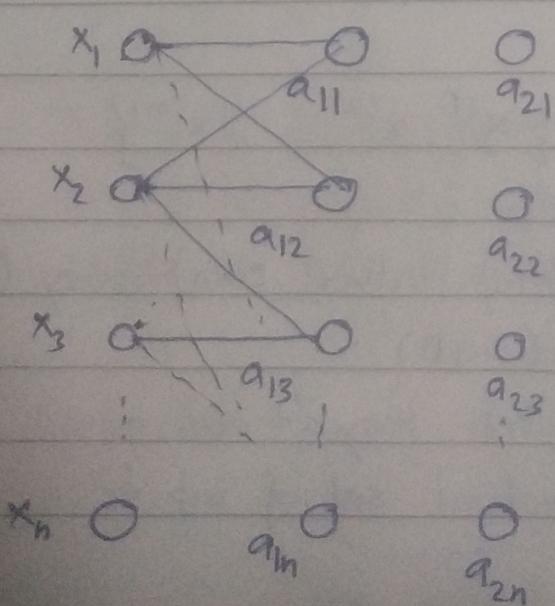
- when the distribution of input features varies across mini-batches
 - ↳ (in other words) mini batches fail to represent actual input effectively
 - ↳ may lead to slow convergence

solⁿ - (proposed by Ioffe & Szegedy)

- normalization may lead to stable & representative mini-batches
- normalization at each layer of NN may help in better propagation of info

Batch Normalization (BN)

- would help if the pre-activations at each layer are normalized
- ensured by standardizing the pre-activations at each unit



$$\hat{q}_{ik} = \frac{q_{ik} - E[q_{ik}]}{SE[q_{ik}]}$$

$E \rightarrow$ expected value

$SE \rightarrow$ standard error

q_{ik}	BN	$z_{ik} = \text{activation func}(q_{ik})$
$q_{11}, q_{12}, q_{13}, \dots, q_{1n}$		

- * local pixels may convey more meaningful info
- ** objects in motion appear different from those at rest
- *** an object is built up from smaller sub-objects or compositions
- **** stride is the no of steps CM is strided on input matrix

a_{ik} has been shifted and scaled which needs to be rescaled & retransformed for getting an output.

$$y^{(k)} = \gamma^{(k)} \hat{a}_{ik} + \beta_b^{(k)}$$

$\gamma^{(k)}$ & $\beta^{(k)}$ need to be learnt & aren't hyperparameters.

DL - 20/03

Image properties

- Locality *
- Stationarity **
- Hierarchy ***

$CM \equiv \text{kernel}$

CNN

↳ convolution layer

$N \times N$ input matrix, $F \times F$ predefined convolution matrix
 It reduces $N \times N \rightarrow (N-F+1) \times (N-F+1)$ (CM)
 for 1 stride. If stride = s, $\left(\frac{N-F}{s}+1\right) \times \left(\frac{N-F}{s}+1\right)$

is the resultant matrix

$$\cos \theta = \frac{\vec{u} \cdot \vec{v}}{|\vec{u}| \cdot |\vec{v}|}$$

\vec{u} & \vec{v} are similar when $\cos \theta \uparrow$ or $\theta \rightarrow 0$



Underlying principle of cosine similarity

further insights paper ref.,

Group equivariance in CNNs

so the convolution operation which is a dot product of kernel with the image, basically tries to identify parts of image similar to kernel. Multiple kernels can be used. Trained kernels are better than predefined ones.
(Learned)

o_bi

The no of units in layer i affecting a unit in layer $i+1$ is the receptive field of unit in layer $(i+1)$.

Due to sparse connectivity in CNN, receptive field goes vulnerable. To accommodate, the network must be deep, otherwise info will be lost.

Invariance

→ If I is an object, $g(I)$ is translation of I for I , then f is said to be translationally invariant if $f(g(I)) = f(I)$

Equivariance*

→ If I is an object, $g(I)$ is translational of I for I , then f is said to be translationally equivariant if $f(g(I)) = g(f(I))$

Convolution is partially invariant as well as partially equivariant.

DL-23/03

Equivariance - shift in object position in image doesn't affect the result

Invariance - slight perturbations in image doesn't affect the result

ReLU is an effective activation funcⁿ as tanh & sigmoid tend to have saturated gradient values for larger values.

Pooling → dimensionality redⁿ
↳ additionally to maintain invariance property

Maxpool vs Avgpool → maxpool better
↳ better invariance

Pooling comes after convolution as convolution weighs the important parts of the image. If pooling is done earlier, it would just identify the more intense parts of the image which may not be significant at all.

padding → borders are under-represented in convolution
↳ to maintain actual dimensions

Let p be the padding size. Stride value = s , input size n and filter size f .

$$\text{Output size} = (n + 2p - f + 1) \times (n + 2p - f + 1)$$

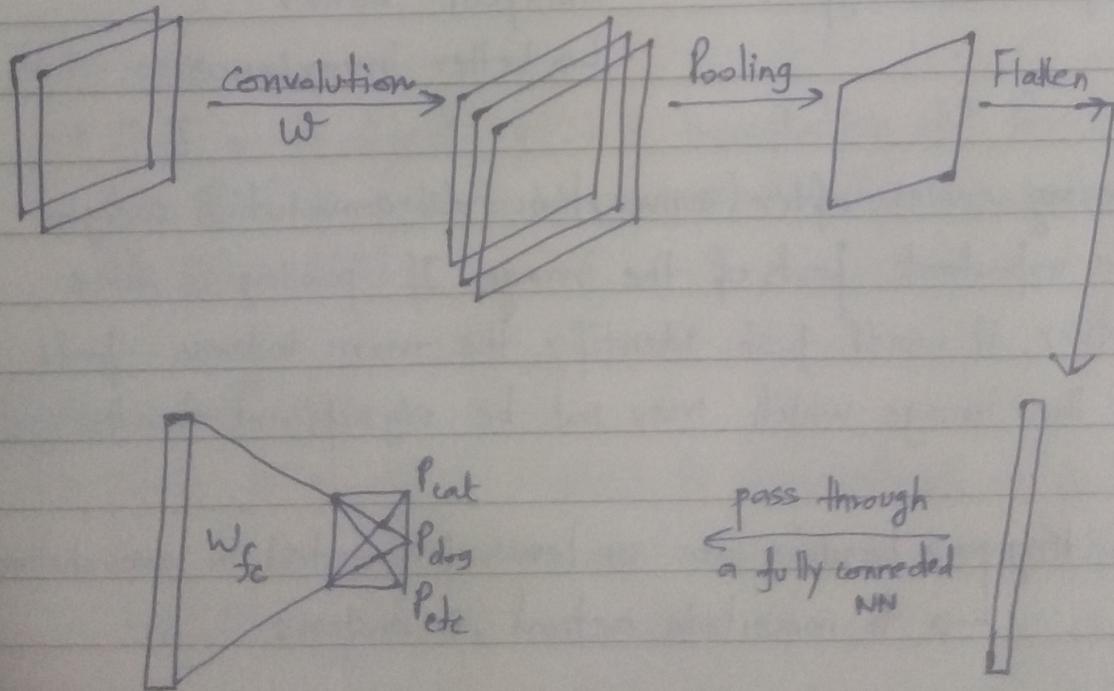
If output size is to be same as input size,

$$\Rightarrow \boxed{p = \frac{f-1}{2}}$$

$$\frac{n+2p-f+1}{s} = n$$

With stride = s,

$$\Rightarrow \boxed{p = \frac{ns - (n+s) + f}{2}}$$



$$L = \max(x_1, x_2) \rightarrow (\text{maxpool})$$

$$\frac{\partial L}{\partial x_1} = \begin{cases} 1 & \text{if } x_1 > x_2 \\ 0 & \text{otherwise} \end{cases}$$

DL - 24/03

AlexNet Architecture

Input: $227 \times 227 \times 3$

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

$$\text{Output size: } \frac{227 - 11}{4} + 1 = 55 \times 55 \times 96$$

Parameters: $(11 \times 11 \times 3) \times 96$

Max pool: 3x3 filters of stride 2

$$\text{Output size: } \frac{55 - 3}{2} + 1 = 27 \times 27 \times 96$$

Parameters: 0

Normalization: $\text{-----}^{(x-H)/6}$ Padding = 2

Conv. 2: 256 filters of dimension $5 \times 5 \times 96$, stride = 1

$$\text{Output size: } \frac{27 - 5 + 1}{1} + 1 = 27 \times 27 \times 256$$

Parameters: ~~($5 \times 5 \times 96$)~~ $(5 \times 5 \times 96) \times 256$

Max pool: 3x3 filters of stride 2

$$\text{Output size: } \frac{27 - 3}{2} + 1 = 13 \times 13 \times 256$$

Parameters: 0

Normalization: $\text{-----}^{(x-H)/6}$ pad = 1

Conv. 3: 384 filters of dimension $3 \times 3 \times 256$, stride = 1

$$\text{Output size: } 13 + 2 - 3 + 1 = 13 \times 13 \times 384$$

Parameters: $(3 \times 3 \times 256) \times 384$ pad = 1

Conv. 4: 384 filters of dimension $3 \times 3 \times 384$, stride = 1

$$\text{Output size: } 13 \times 13 \times 384$$

Parameters: $(3 \times 3 \times 384) \times 384$

* 1000 classes of ~~dataset~~ images in dataset

Conv. 5: 256 filters of dimension $3 \times 3 \times 384$, stride=1
 $\text{Pad} = 1$

Output size: $13 \times 13 \times 256$

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

Max pool: 3×3 filters of stride 2

Output size: $\frac{13-3}{2} + 1 = 6 \times 6 \times 256$

Parameters: 0

Fully-Connected Layer 1: 4096 neurons

Output size: 4096 Parameters: 4096×4096

Fully-Connected Layer 2: 4096 neurons

Output size: 4096 Parameters: 4096×4096

Fully-Connected Layer 3: 1000 neurons

Output Size: 1000 * Parameters: 4096×1000

Total 60M parameters, no of parameters ↑↑

ReLU activation is used after every convolution.

Challenges - performance

Improvement - VggNet

- made the network deeper to boost accuracy & performance (kernel size ↓)
- parameters learned at each layer is less, but as depth ↑ total no of parameters ↑ (Parameters = 138M)
- receptive field ↑ for ~~deeper layers~~ whole network, helps extracting complex features

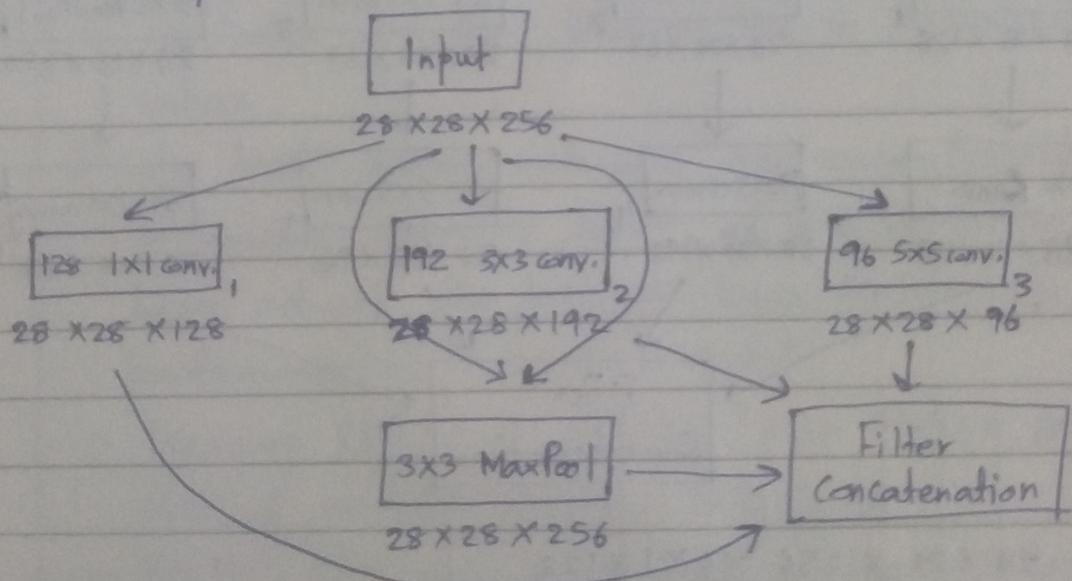
* Perception & observation aren't complete.
We perceive an image at different scales
at once & try to augment our observations
DL - 25/03 to learn what the image tries to convey.

Complexity of algorithms like VggNet & AlexNet is measured by no of Floating-Point Operations (FLOPs).

VggNet FLOPs ↑

GoogleNet objective → to reduce FLOPs & achieve comparable performance as that based on the idea of Inception

Naive Inception*



$$\begin{aligned} \text{Output size: } & 28 \times 28 \times (128 + 192 + 96 + 256) \\ & = 28 \times 28 \times 672 \end{aligned}$$

Considering only MULTIPLICATIONS, total no of operations;
(MULTs) (FLOPs)

$28 \times 28 \times 256 \times 1 \times 1 \times 128$ for conv. 1

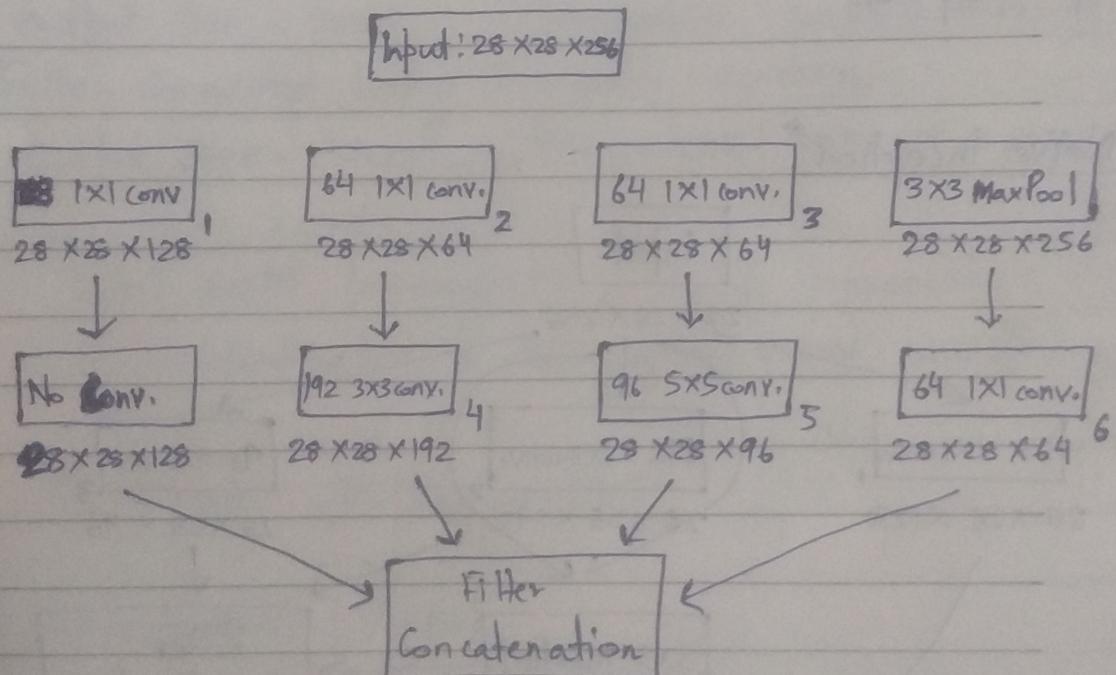
$28 \times 28 \times 256 \times 3 \times 3 \times 192$ for conv. 2

$28 \times 28 \times 256 \times 5 \times 5 \times 96$ for conv. 3

Max Pool doesn't require significant no. of MULTs.

Total FLOPs ≈ 854 M which is very high.

GoogleNet introduced a bottleneck layer to naive inception to reduce FLOPs.



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

Total FLOPs ≈ 358 M

- * paper on ResNet by Microsoft scientists, reference
- ** ResNet-1202 had error % of 7.93, more than ResNet-110 due to overfitting.

GoogleNet

- tries to visualise an image the same way a human does
- reduces no. of Flops than naive inception using bottleneck layer

ResNet *

- assumption of a deeper network performing better is contradicted
- 56-layer network performed worse than a 20-layer network experimentally & it was not due to overfitting as training error was also worse for 56-layer network
- conclusion was that problem lies in initialisation values
- also that deeper networks are uncomfortable in learning identity funcⁿ, as if in the 56-layer network, the last 36 learn $f(x) = x$, the performance of both would match
- learning $F(x) = 0$ is feasible than $F(x) = x$, idea of residual learning introduced
- ResNet idea was implemented & the results were now as expected. ResNet-34 performed better than ResNet-18 on both validation & ~~test cases~~ generalization
- ResNet performed better & error went down with increasing depth, ResNet-110 is minⁿ with 1.7M parameters & 6.43% error **

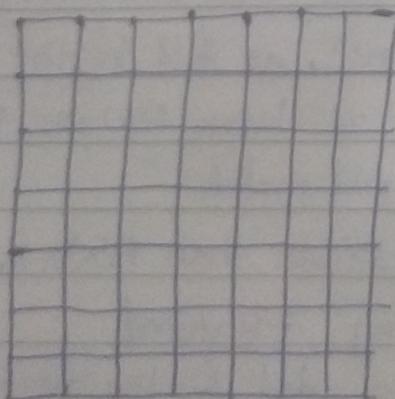
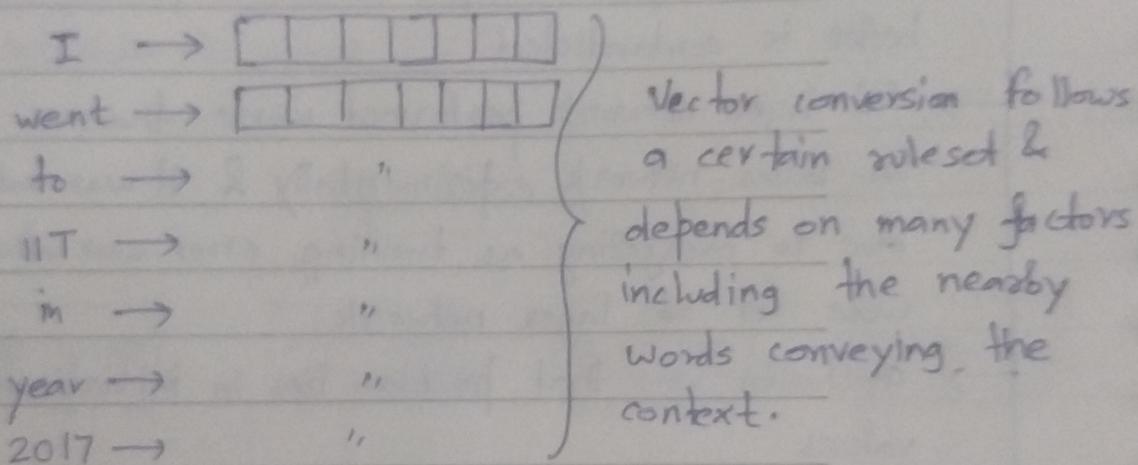
* 'I went to IIT in year 2017' &

DL - 27/03 * 'In 2017, I went to IIT' are same but CNN would fail in this model to identify them similar.

Recurrent NN (RNN)

Text comprehension,

Model → convert each word to a fixed-length vector & concatenate to form an image, apply CNN on the image for perception

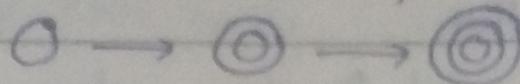


→ apply CNN

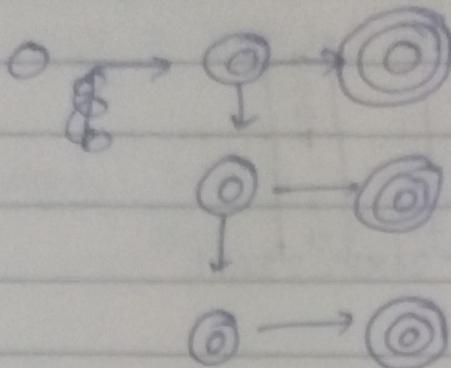
Problem

- Texts can be indefinitely long arbitrarily
- there may be non-local strong dependencies
- grammatically & orientationally different texts may convey same meaning *

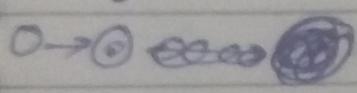
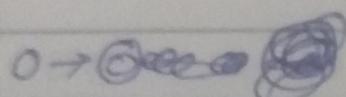
Common architectures,



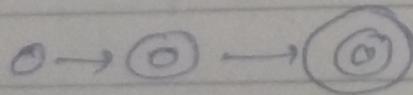
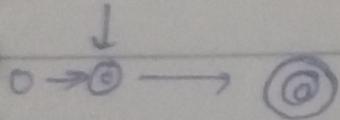
Vector conversion of word



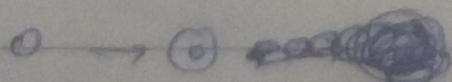
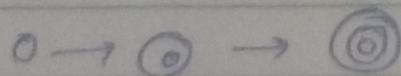
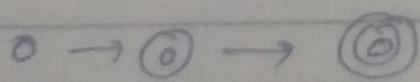
Text generation



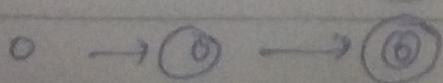
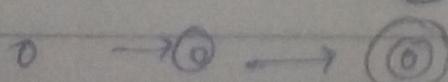
Sentiment classification



Speech to text



Machine translation



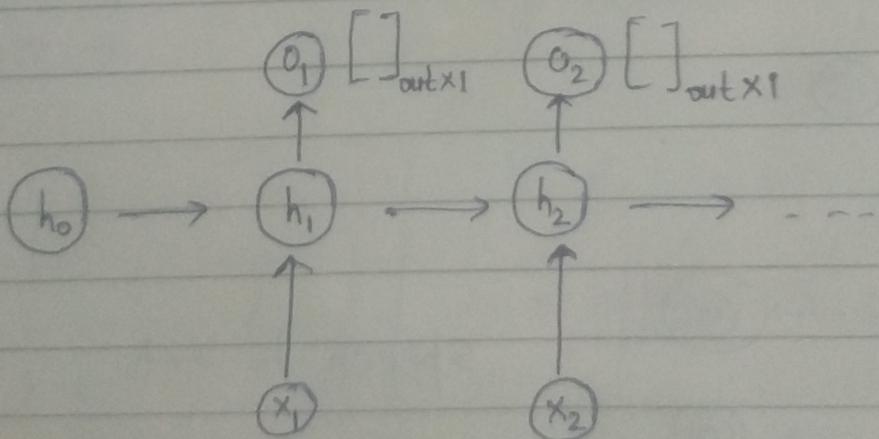
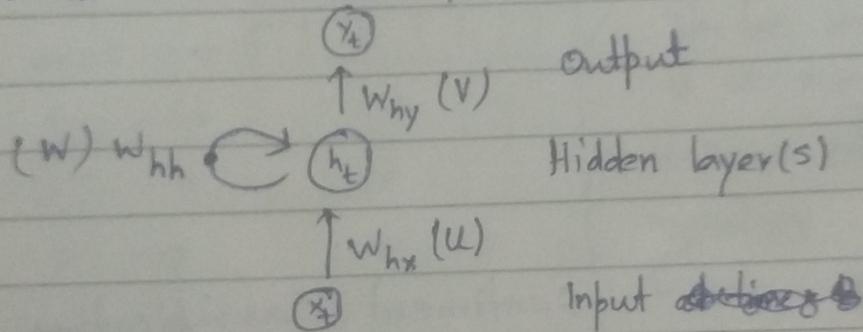
One-hot encoding / representation

↳ representing a word using vectors s.t. only 1 of the elements is 1 & rest 0

ex. bad →

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{26 \times 1}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{26 \times 1}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}_{26 \times 1}$$

Basic RNN architecture,



$$E_t = \frac{1}{2} \sum_{k=1}^{out} (\hat{y}_k - y_k)^2 \quad \text{for 1 layer}$$

$$E = \frac{1}{2} \sum_{t=1}^T \sum_{k=1}^{out} (\hat{y}_{tk} - y_{tk})^2$$



Basic RNN eqns,

$$h_t = Ux_t + W\phi(h_{t-1})$$

$$y_t = V\phi(h_t)$$

$$\text{Gradient, } \frac{\partial E}{\partial W} = \sum_{t=1}^T \frac{\partial E_t}{\partial W} = \sum_{t=1}^T \sum_{k=1}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$

$$\frac{\partial h_t}{\partial h_k} = \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdots \frac{\partial h_{k+1}}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}}$$

DL - 30/03

$$L = ab \Rightarrow \frac{\partial L}{\partial z} = a^T \frac{\partial b}{\partial z} + b^T \frac{\partial a}{\partial z}$$

$$h_i = Ux_t + W\phi(h_{i-1})$$

$$\Rightarrow \frac{\partial h_i}{\partial h_{i-1}} = W^T \frac{\partial \phi(h_{i-1})}{\partial h_{i-1}}$$

$$\phi(h_{i-1}) = \begin{bmatrix} \phi^{(1)}_{h_{i-1}} \\ \phi^{(2)}_{h_{i-1}} \\ \vdots \\ \phi^{(k)}_{h_{i-1}} \end{bmatrix}$$

$$h_i = \begin{bmatrix} h_i^{(1)} \\ h_i^{(2)} \\ \vdots \\ h_i^{(k)} \end{bmatrix}$$

$$\frac{\partial \phi(h_{i-1})}{\partial h_{i-1}} = \begin{bmatrix} \frac{\partial \phi^{(1)}_{h_{i-1}}}{\partial h_{i-1}^{(1)}} & \frac{\partial \phi^{(1)}_{h_{i-1}}}{\partial h_{i-1}^{(2)}} \\ \frac{\partial \phi^{(2)}_{h_{i-1}}}{\partial h_{i-1}^{(1)}} & \frac{\partial \phi^{(2)}_{h_{i-1}}}{\partial h_{i-1}^{(2)}} \\ \vdots & \vdots \\ \frac{\partial \phi^{(k)}_{h_{i-1}}}{\partial h_{i-1}^{(1)}} & \frac{\partial \phi^{(k)}_{h_{i-1}}}{\partial h_{i-1}^{(2)}} \end{bmatrix}, \quad = \begin{bmatrix} \phi^{(1)}_{h_{i-1}} & 0 \\ 0 & \phi^{(2)}_{h_{i-1}} \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix}$$

$$\text{So, } \frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=k+1}^t W^T \text{diag}[\phi'(h_{i-1})]$$

If $\frac{\partial h_i}{\partial h_{i-1}} < 1 \quad \forall i, \frac{\partial h_t}{\partial h_k} \rightarrow 0.$

- ↳ Problem of Vanishing Gradients
 - RNN stops learning at a point
 - New weights aren't updated, non-optimal weights remain

If $\frac{\partial h_i}{\partial h_{i-1}} > 1 \quad \forall i, \frac{\partial h_t}{\partial h_k} \rightarrow \infty.$

- ↳ Problem of Exploding Gradients
 - RNN weights updated drastically
 - can be handled by clipping gradients by a threshold, called gradient clipping

Solutions to Vanishing Gradients

- ↳ do backpropagation over smaller no of steps
- ↳ problem - error is also propagated.

$$h_t = Ux_t + W\phi(h_{t-1})$$

$m \times 1$ $m \times 1$

new input $I = \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \in \mathbb{R}^{2m}$ $W_a^\phi = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ w & & & & u \\ 1 & \dots & 1 & \dots & 1 \end{bmatrix}$

\rightarrow activation funcⁿ

$$h_t = \phi(W_a [h_{t-1}, x_t] + b_a)$$

↓
bias

Applications of RNN - Language Generation and Sequence Models

Let y_1, y_2, \dots, y_n be the words of a language
($n \rightarrow 10^6$ or even 10^{80} for rich language like English)

Goal of a language model

- ~~try to find~~ derive the probability of a given sequence of words $P(y_1, y_2, \dots, y_k)$
- helpful in tasks like speech recognition

Ex: The man — a cart. (Pooled / pulled)

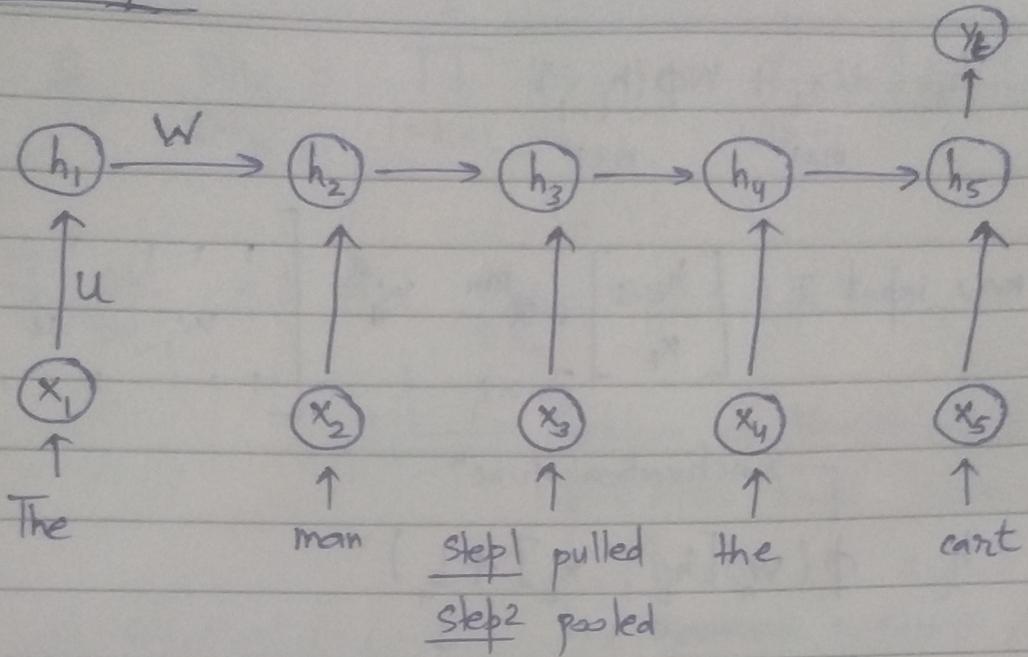
Find $P(\text{The man pooled a cart})$ & $P(\text{The man pulled a cart})$

$$P(E_1)$$

$$P(E_2)$$

* The model must establish $P(E_1) < P(E_2)$.

DL - 31/03



$x_t \rightarrow$ prob. score, whether input sentence is correct

Input with higher prob. score is used to fill up the blank.

Training

→ use a set of sentences

→ known sentences can be altered by changing / removing words so that incorrect sentences are also in the set so as to make the model learn

Model

$$h_t = \sigma_h (Ux_t + Wh_t + b_h)$$

$$\hat{y}_t = \sigma_y (Vh_t + b_y)$$

$b_h, b_y \rightarrow$ bias

$\sigma_h, \sigma_y \rightarrow$ activation function

Improvement

→ use extra hidden layers

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

Problem extended Now the position of missing word is also unknown.

Ex. The man the cart.

→ To identify the appropriate position & appropriate word.

Approach → Pick random words from corpus, put at each position in the sentence & check probability scores. 1 word, 5 positions in ex.

- Problem
- Sentences can be arbitrarily long, so the problem of vanishing gradients may come up
 - carrying info over large portions of text is difficult as knowledge will fade away due to vanishing gradients (Memory problem)

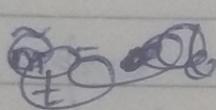
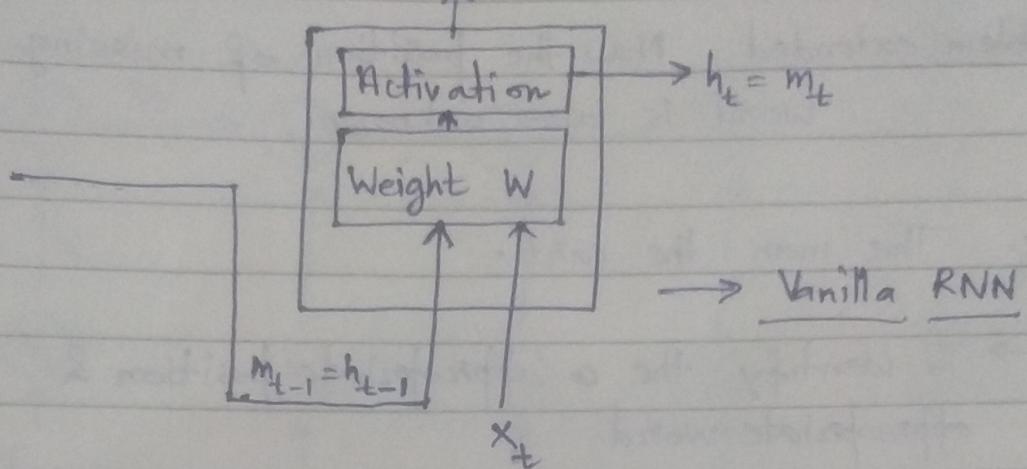
Solⁿ → Gated Recurrent Unit (GRU)

- maintains a memory cell $m_t = h_t$ that maintains contains relevant info like whether noun encountered was singular/plural & so on
- at every step, GRU needs to decide whether ~~copy~~ m_t needs to be replaced by an alternate/candidate value \tilde{m}_t

* $\odot \rightarrow$ Hadamard product

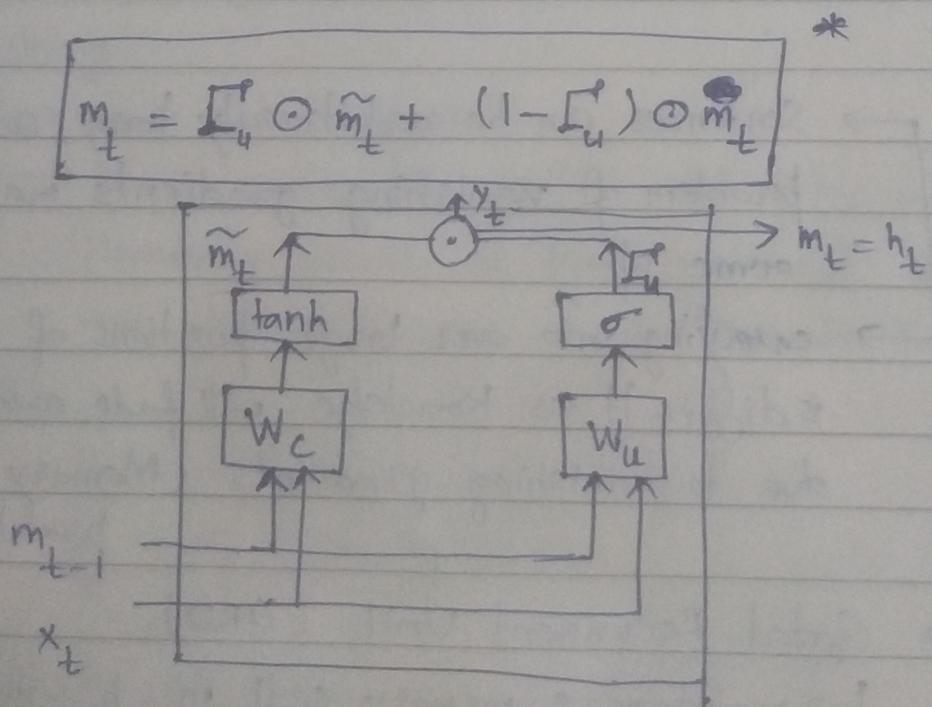
$$a = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \quad y_t$$

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



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

To make the decision, an update gate $\Gamma_u = \sigma(W_u[m_{t-1}, x_t] + b_u)$ is used.

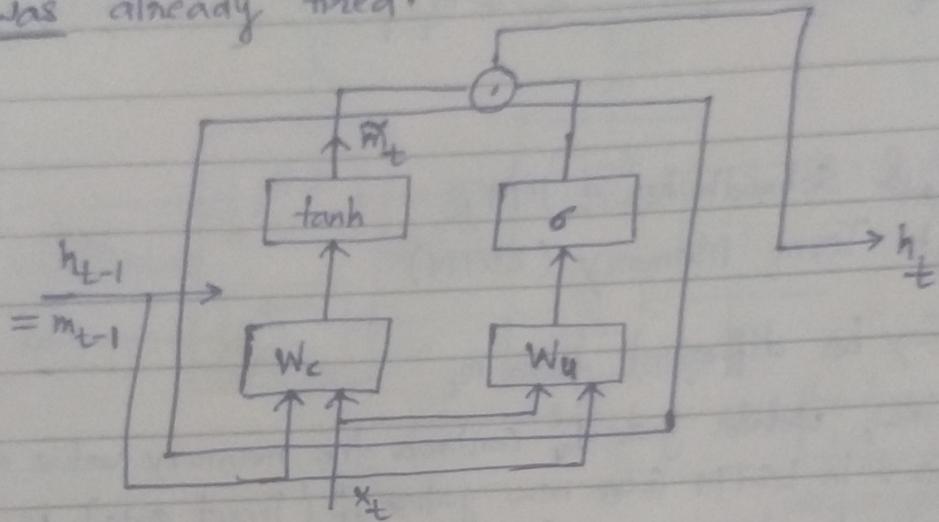


Simplified GRU

DL - 06/04

In arbitrary-long sentence, GRU on memory retention over a long sequence is necessary for noun-verb agreement.

Ex:
The man who pulled the cart when the king arrived
was already tired.



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

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

$$m_t = \tilde{f}_u \odot \tilde{m}_t + (1 - \tilde{f}_u) \odot m_{t-1}$$

even if $\tilde{f}_u \rightarrow 0$ due to vanishing gradients, m_t won't be 0 due to this. So memory retention is improved.

A full GRU tries to further improve over simplified GRU. The idea is to consider only data relevant to assigned task by using a relevance gate.

* GRU architecture ref.,

Cho et al. 2014, on the properties of Neural Machine Translation, encoder & decoder approaches

Full GRU \leftrightarrow

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

$$l_u^t = \sigma(W_u [m_{t-1}, x_t] + b_u) \quad (\text{Update Gate})$$

$$l_r^t = \sigma(W_r [m_{t-1}, x_t] + b_r) \quad (\text{Relevance Gate})$$

$$m_t = l_u^t \odot \tilde{m}_t + (1 - l_u^t) \odot m_{t-1}$$

④

Hochreiter & Schmidhuber - 1997

Long Short Term Memory (LSTM)

→ h_t may be different from m_t

→ In GRU, update gate l_u^t controls the memory value m_t , whereas in LSTM ~~as~~ new gates l_f^t (forget gate) is introduced along with l_o^t (output gate)

→ l_f^t controls m_t & l_o^t controls h_t

$$l_f^t = \sigma(W_f [h_{t-1}, x_t] + b_f)$$

$$m_t = l_u^t \odot \tilde{m}_t + l_f^t \odot m_{t-1}$$

$$h_t = l_o^t \odot m_t$$

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

$$\tilde{m}_t = \tanh(W_m [h_{t-1}, x_t] + b_m) \quad (W_m = W_c \text{ in GRU})$$

$$l_u^t = \sigma(W_u [h_{t-1}, x_t] + b_u)$$

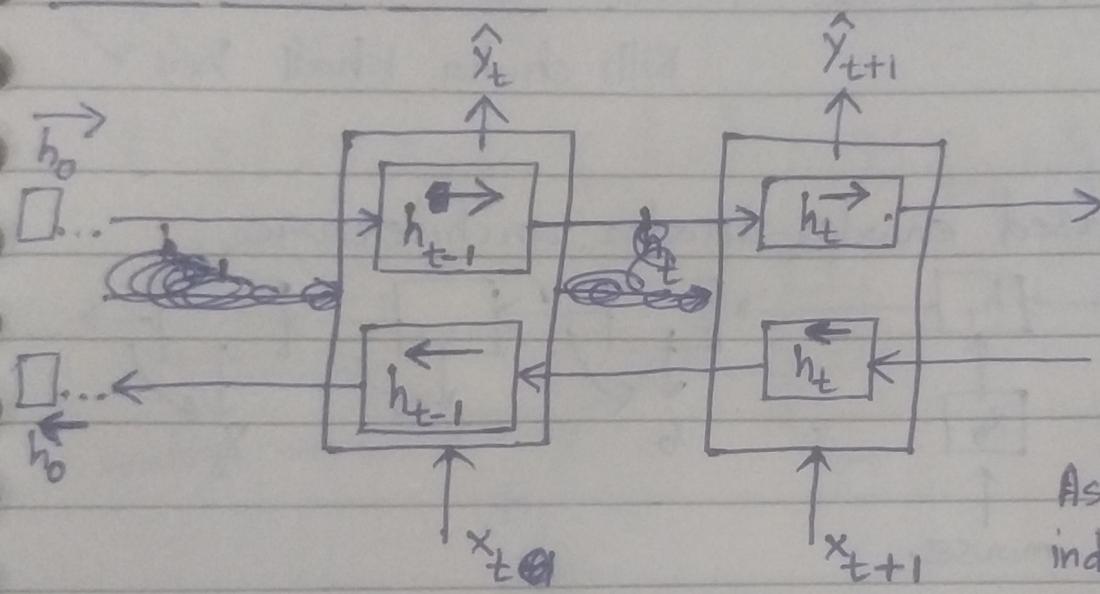
{ He said, " Teddy bears are on sale." }
He said, " Teddy Roosevelt was a great president" }

→ significance of looking forward in a sentence

Understanding forward parts is also important as illustrated.

↳ Bidirectional networks!

Bi-directional RNN



Forward & backward can ~~not~~ be done independently or in parallel.

As they may be independent, GRU ~~is~~ can also be used here.

$$h_t = [\vec{h}_t \quad \overleftarrow{h}_t]$$

$$\vec{h}_t = \sigma(W_h [h_{t-1}, x_t] + b_{h_t})$$

$E[(y - \hat{y})^2]$ is to be minimized.

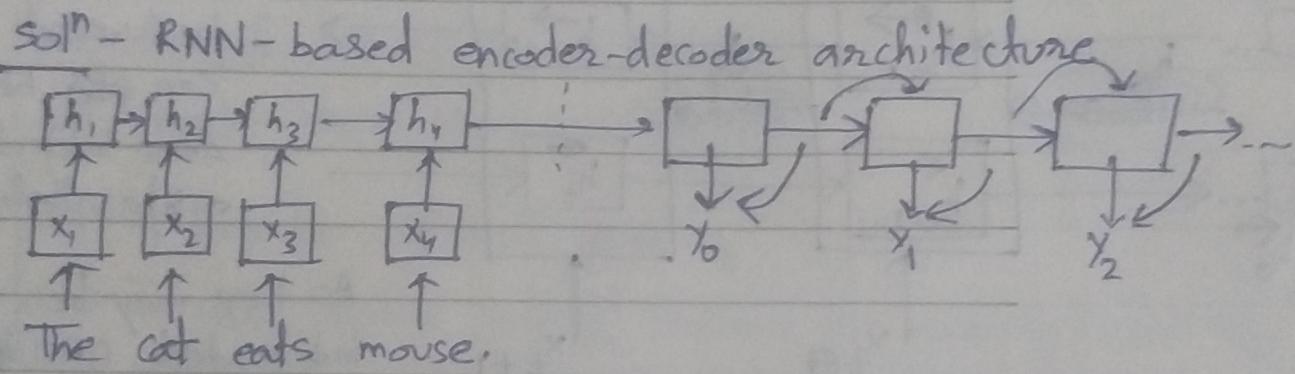
* Ref. Neural Machine Translation by jointly learning to align & translate
— Bahdanau et al., ICLR 2015

DL-10/04

Attention & Transformers*

Word to word translation ^{may be} is meaningless.

Ex. The cat eats mouse. → Wo billi khati chuha X
Billi chuha khati hai ✓



Problem — Sentences can be arbitrarily long. Capturing the context of the whole sentence in a lower dimension embedding may produce inefficient results.

Solⁿ — Create a context vector. Feed it to the decoder so as the decoder can choose the best possible output given the input sentence & the context.

RNN-based encoders

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

Encoder

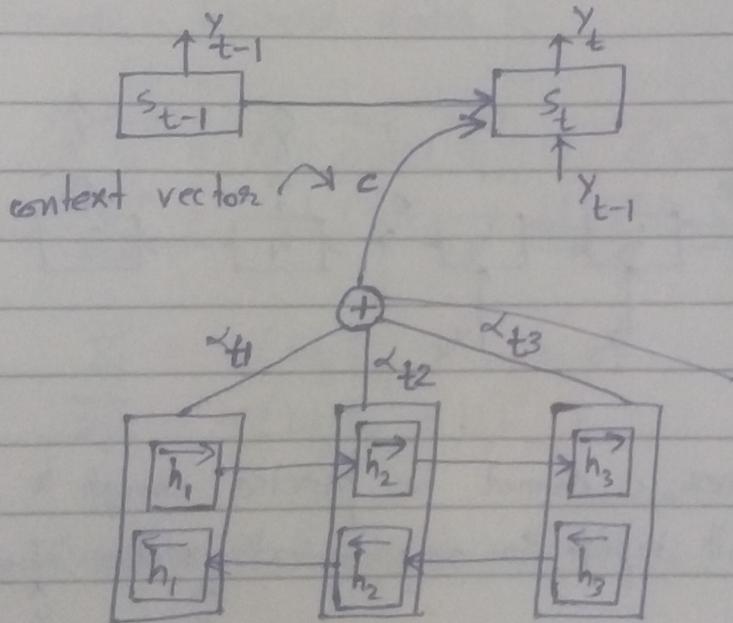
$$h_t = f(h_{t-1}, x_t)$$

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

$$p(y) = \prod_{t=1}^{T_x} p(y_t | \{y_1, y_2, \dots, y_{t-1}\}, c)$$

Decoder

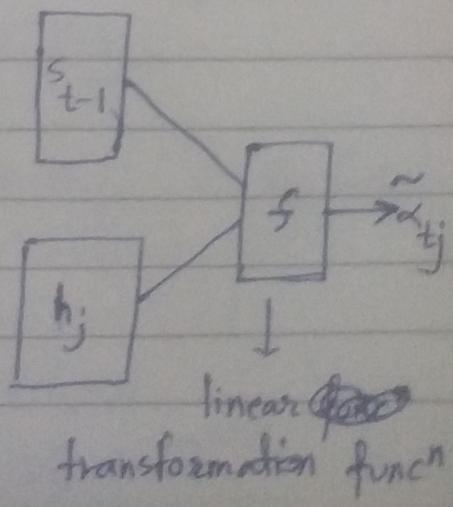
$$p(y_t | \{y_1, y_2, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$$



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

$$\tilde{\alpha}_{tj} = f(s^{(t-1)}, h^{(j)})$$

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



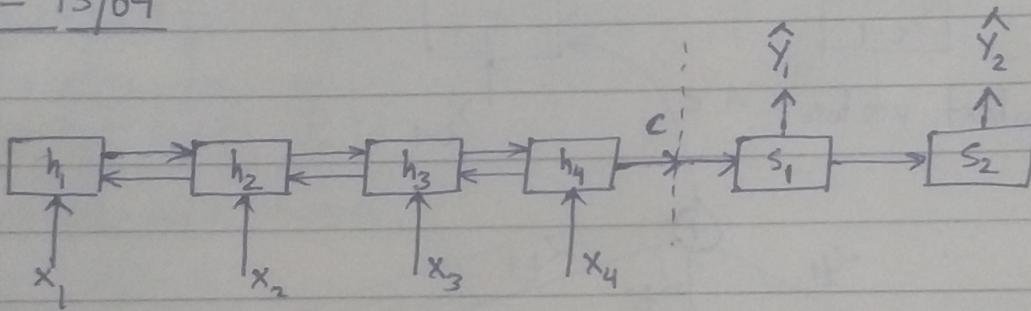
Learnable parameters are weights,

- (1) Weights used in f on linearly transforming $s^{(t-1)}$ & $h_j^{(t)}$
- (2) Weights used to generate \vec{h}_t & \hat{h}_t from given input

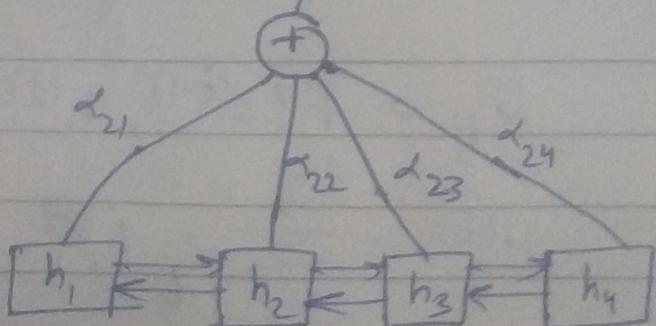
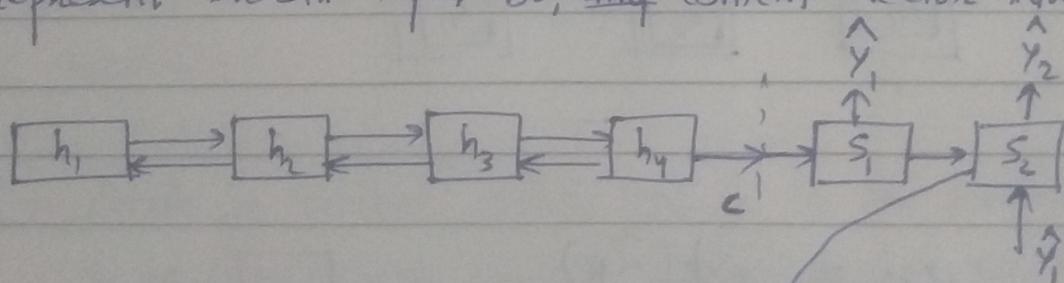
$$y_1 \ y_2 \ y_3 \dots y_n \quad p(\hat{y}_1) p(\hat{y}_2) \dots p(\hat{y}_n)$$

$$\alpha = \prod_{y_i \in T_y} p(\hat{y}_i) \rightarrow \text{optimization by maximizing the likelihood}$$

DL - 13/04



For long sentences, c cannot be precise enough to represent overall input. So, ~~isn't~~ context vector introduced.



* z_i for encoder & decoder aren't the same

$$** \quad \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^K \exp(e_{ik})} \quad | \quad e_{ij} = v_a^T \tanh(u_h s_{i-1} + u_b h_j)$$

Encoder :- ~~GRU~~ Bidirectional GRU

Source sentence \rightarrow word tokens \rightarrow 1 of K vectors

Construct a vocabulary of K_x words.

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

Input in $X = \{x_1, x_2, \dots, x_T\}$ where $x_i \in \mathbb{R}^{K_x}$

↑
1 of K vector

Translation output $Y = \{y_1, y_2, \dots, y_T\}$ where $y_i \in \mathbb{R}^{K_y}$

Encoder: forward states $\overrightarrow{h_i}$

$$\begin{aligned}\overrightarrow{h_i} &= (1 - \overline{z}_i) \odot \overrightarrow{h}_{i-1} + \overline{z}_i \odot \overrightarrow{h}_i \\ \overrightarrow{h_i} &= \tanh(W_E x_i + U [z_i \odot \overrightarrow{h}_{i-1}]) \quad (E_{x_i} : \text{embedding of } x_i)\end{aligned}$$

$$\overline{z}_i = \sigma(W_z E_{x_i} + U_z \overrightarrow{h}_{i-1}) \quad (\text{Update Gate})$$

$$\overline{z}_i = \sigma(W_r E_{x_i} + U_r \overrightarrow{h}_{i-1}) \quad (\text{Relevance gate or Reset})$$

Same equ" for backward states.

~~Dec~~ Decoder :- GRU

w, u, c are weights

$$s_i = (1 - z_i^d) \odot s_{i-1} + z_i^d \odot \underline{s}_{i-1} *$$

$$\underline{s}_i = \tanh(W_E y_{i-1} + U [z_i \odot s_{i-1}] + C c_i) \quad c_i = \sum_{j=1}^K \alpha_{ij} h_j **$$

$$z_i^d = \sigma(W_z E_{y_{i-1}} + U_z s_{i-1} + C_z c_i)$$

$$z_i^d = \sigma(W_z E_{y_{i-1}} + U_z s_{i-1} + C_z c_i)$$

$$s_i = \tanh(W_s \overleftarrow{h}_1) \quad W_s \in \mathbb{R}^{n \times n}$$

Probability of output,

$$P(y_i | s_i, y_{i-1}, c_i) \propto \exp(y_i^T w + t_i)$$

$$t_i \sim U_0 s_{i-1} + V_0 E y_{i-1} + C_0 c_i$$

→ maximize likelihood to minimize loss funcⁿ &
learn the weights

Performance Measures in Machine Translation

↳ Bleu score

Consider 2 sentences,

(1) The ball is blue (Predicted)

(2) The ball has a blue color (Translated)

Attempt - consider the ~~no~~^{overlap} of words appearing in
predicted sentence & translated sentence

Failure - "The ball the ball" overlaps 100%
with (2) but is very wrong

Solⁿ: Tokenise into N-grams, find precision

* If all weights are $w_n = \lambda_N$, then

$$P_{\text{gavg}} = \exp\left(\frac{N}{n=1} \log p_n^{\lambda_N}\right) = \exp(\log(p_1 p_2 \dots p_n)^{\lambda_N})$$

$$\Rightarrow P_{\text{gavg}} = (p_1 p_2 \dots p_n)^{\lambda_N}$$

For (1)

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

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

3-gram : "The ball is", "ball is blue"

$$\text{precision} = \frac{\# \text{ correct predicted words}}{\text{total predicted words}}$$

$$\text{In (1), \& (2), precision} = \frac{3}{4}$$

But there can be repetitions, "the the the ball"
or multiple target sentences

So, use clipped precision, if a word appears > 1 times
in predicted sentence, or if it is there in multiple
sentences (target), count its value as 1 only.

To find bleu score,

find precision for 1-gram - p_1

" " " 2-gram - p_2

" " " 3-gram - p_3

- - - - -

" " " n-gram - p_n

(GAP)

Find geometric avg precision score, P_{gavg} ,

$$P_{\text{gavg}} = \exp\left(\frac{N}{n=1} w_n \log p_n\right)$$

w_n = weight for
n-gram precision

γ = reference text length

find brevity ~~length~~ penalty, (BP)

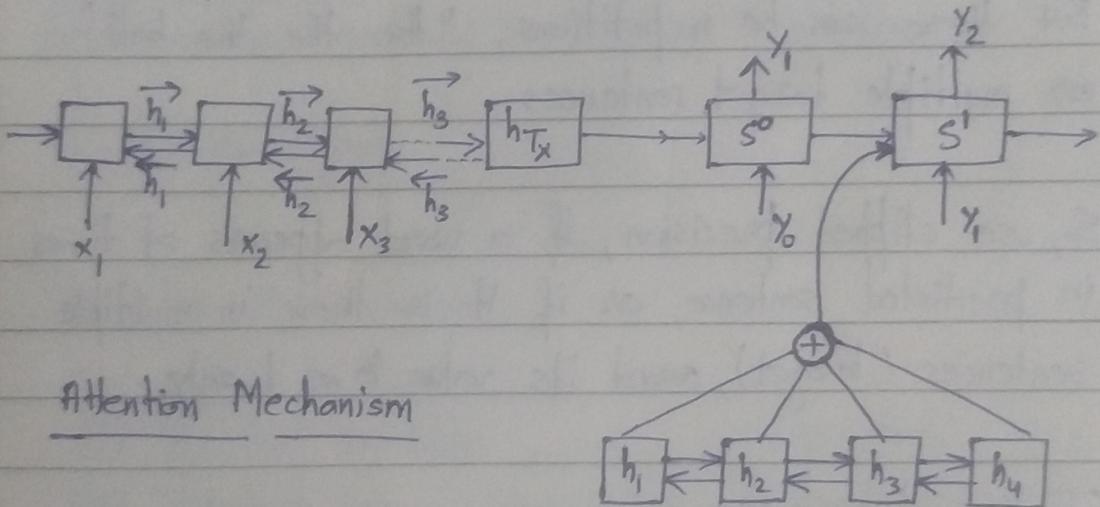
(penalises short translations)

$$BP = \begin{cases} 1 & ; \text{ if length} > \gamma^* \\ e^{(1-\frac{\gamma}{\text{len}})} & ; \text{ else} \end{cases}$$

↳ predicted length

Bleu Score = GAP. BP

DL - 15/04



Attention Mechanism

For arbitrarily long sentences, there are recurrence units in the model which tend to slow down the training significantly.

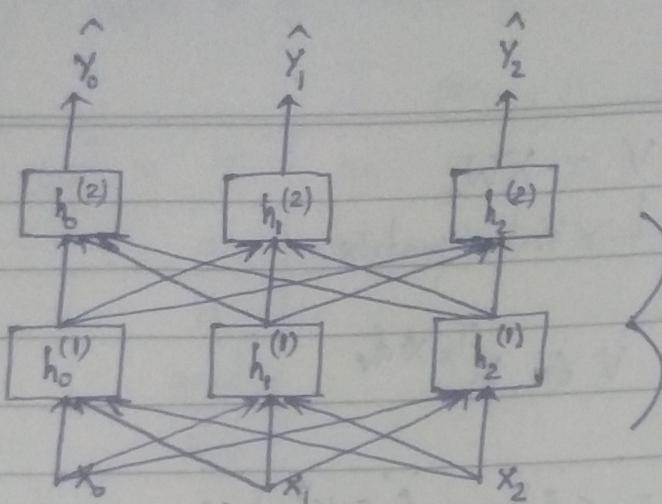
Improvement

Attention is all you need

- Vaswani et al., NIPS 2017

Sequential RNN-based architecture can't be parallelized to optimize computation complexity. Transformations introduced.

Intuition for Transformer architecture



Attention of any unit
is upon all units in
previous layer

Advantage Sequential dependency & recurrences have been eliminated. Parallelization can be done to optimize computational complexity.

Ex. Given $x_0, x_1 \& x_2$, $h_0^{(1)}, h_1^{(1)} \& h_2^{(1)}$ can be evaluated ~~in~~ in parallel, & so on for other layers.

Consider a word sequence, "I am a good boy."

$$X = \begin{bmatrix} I & .75 & .66 & .21 & \dots \\ am & .5 & .2 & .34 & \dots \\ a & .6 & .4 & .32 & \dots \\ good & .1 & .25 & .35 & \dots \\ boy & .5 & .26 & .5 & \dots \end{bmatrix} \in \mathbb{R}^{T_x \times d_k}, T_x = 5$$

transformation follows some rule

$$q_i = W_q x_i ; W_q \in \mathbb{R}^{d_k \times d_k}$$

$$Q = XW_q \in \mathbb{R}^{T_x \times d_k}$$

↳ Query matrix formed by a transformation of X by W_q , a set of learnable parameters, similar for $K \& V$

$K = XW_k$ & $V = XW_v$
 ↳ Key matrix ↳ Value matrix

$$K \in \mathbb{R}^{T_x \times d_k} \quad V \in \mathbb{R}^{T_x \times d_v}$$

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

$$Q = XW_q = I \begin{bmatrix} .65 & .62 & \dots \\ .31 & .67 & \dots \\ .45 & .47 & \dots \\ .32 & .62 & \dots \\ .47 & .68 & \dots \end{bmatrix} = \text{Some linear transformation of } X$$

$$QK^T = I \begin{bmatrix} \text{am} \\ \text{a} \\ \text{good} \\ \text{boy} \end{bmatrix} \begin{bmatrix} \text{I} \\ \text{am} \\ \text{a} \\ \text{good} \\ \text{boy} \end{bmatrix}^T$$

• when evaluating QK^T , row of "I" is multiplied by K^T column of "I", so result is some kind of a pairwise attention of "I" over "I", same for other words.

Also since words are aligned vertically in Q & horizontally in K^T , a pairwise dependency of each word upon another is also evaluated.

* d_k may also be same as d_v . The difference just helps in a smooth generalisation.

I am a good boy

$$QK^T = I \begin{bmatrix} .75 & .45 & .36 & \dots \\ .86 & .86 & .75 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \text{good} & \text{boy} & \text{good} & \text{boy} \end{bmatrix}$$

~~as~~ as $Q \in \mathbb{R}^{Tx \times dk}$ & $K^T \in \mathbb{R}^{dk \times Tx}$
 $QK^T \in \mathbb{R}^{Tx \times Tx}$

as d_k is dimension of Q & K^T , it ~~does~~ does a good job at scaling the values in QK^T . Also known as scaled normalization.

The softmax over $\frac{QK^T}{\sqrt{dk}}$ now changes the matrix into a

probability distribution. In other words, in the row of "I", the column of word with highest value (probability) is the most relevant in deciding "I" & so on. Also known as scaled dot product normalization.

Now when multiplied by V , it finally gives how much attention is to be paid to a word in deciding another word.

$$\text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V = I \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ \text{good} & \text{boy} & \text{good} & \text{boy} \end{bmatrix} \begin{bmatrix} .5 & .7 & \dots \\ .8 & .6 & \dots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

(w_q, w_k, w_v)

↳ In multi-head attention, use different values for these 3 learnable parameters

In multi-head attention, different matrices are obtained which are then normalized after adding over entire input set.

$$\begin{array}{c}
 b_1 \quad b_2 \quad \dots \quad b_m \\
 \downarrow \qquad \downarrow \qquad \dots \qquad \downarrow \\
 x_1 \quad | \quad x'_1 \quad | \quad \dots \quad | \quad x'_m \\
 x_2 \quad | \quad x'_2 \quad | \quad \dots \quad | \quad x'_m \\
 \vdots \quad | \quad \vdots \quad | \quad \dots \quad | \quad \vdots \\
 x_n \quad | \quad x'_n \quad | \quad \dots \quad | \quad x'_m
 \end{array}
 \quad
 x'_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

normalize normalize

↳ Layer-wise batch normalization

As there are separate batches, the whole values are evenly distributed.

DL - 20/04

Ex I am a good boy.

$$X = \begin{bmatrix} I & .7 & .5 & \dots \\ am & .2 & .8 & \dots \\ a & & & \dots \\ good & & & \dots \\ boy & & & \dots \end{bmatrix}$$

$$\begin{aligned}
 q_i &= w_q x_i, & Q &= X w_q \\
 k_i &= w_k x_i, & K &= X w_k \\
 v_i &= w_v x_i, & V &= X w_v
 \end{aligned}$$

Self-attention (Q, K, V) = softmax

$$X \in \mathbb{R}^{T_x \times d_{\text{model}}}$$

$$w_q \in \mathbb{R}^{d_{\text{model}} \times d_k}$$

$$x w_q = Q \in \mathbb{R}^{T_x \times d_k}$$

$$\left(\frac{Q K^T}{\sqrt{d_k}} \right) V$$

$$QK^T \in \mathbb{R}^{T_x \times T_x}$$

↳ kind-of self-attention for each word on each word

I am a good boy

$$Q = I \quad \left[\begin{array}{c} \text{am} \\ \text{a} \\ \text{good} \\ \text{boy} \end{array} \right] \quad QK^T = I \quad \left[\begin{array}{c} \text{am} \\ \text{a} \\ \text{good} \\ \text{boy} \end{array} \right] \quad \left[\begin{array}{c} \text{I} \\ \text{am} \\ \text{a} \\ \text{good} \\ \text{boy} \end{array} \right]$$

$$QK^T = I \quad \left[\begin{array}{ccccc} \text{I} & \text{am} & \text{a} & \text{good} & \text{boy} \\ .7 & .5 & .3 & .5 & .2 \end{array} \right] \quad \rightarrow \text{self-attention of each word on each word}$$

Problem :- values can be arbitrarily large

- loss funcⁿ doesn't change much wrt weights
- gradient doesn't change much
- inefficient learning

Solution - scale ($\sqrt{d_k}$) & apply softmax

→ gets converted to probability scores

→ "how important words are wrt some word"

* $d_V = d_K$ in paper. For general case, they may be taken to be different.

Multi-head attention, there are different sets of weights of w_q, w_k & w_v which can be stacked or concatenated.

$$\text{let } z_i = \underset{\text{softmax}}{\max} \left(\frac{(Q^{(i)})^T k(i)}{\sqrt{d_k}} \right) v(i) \in \mathbb{R}^{T_x \times d_v}$$

$$z_1 = T_x \begin{bmatrix} \\ \vdots \\ \\ \end{bmatrix} \quad \dots \quad z_n = T_x \begin{bmatrix} \\ \vdots \\ \\ \end{bmatrix}$$

$$z_{n+1} = T_x \begin{bmatrix} \\ \vdots \\ \\ \end{bmatrix}$$

Layer Normalization ($X + \text{sublayer}(X)$)

$$Z = [z_1 \ z_2 \ \dots \ z_n] W_o$$

\rightarrow transformation matrix to bring Z to dimensions same as X ($\mathbb{R}^{T_x \times d_{\text{model}}}$)

I-1 I-2 ... I-n

$$\begin{bmatrix} x_{-1} \\ x_{-2} \\ x_{-3} \\ \vdots \\ x_{-m} \end{bmatrix} \xrightarrow{\text{Column-wise Normalization}} \begin{bmatrix} \text{Difference of mean across batches is normalized} \end{bmatrix}$$

Row-wise Normalization

Batch Normalization (BN)

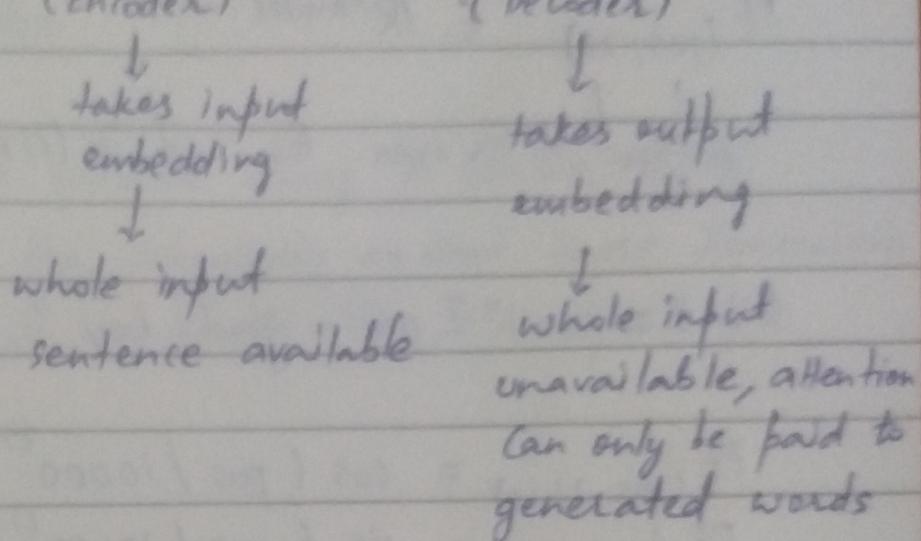
\downarrow
Layerwise Normalization (LN)

$$LN \quad \mu = \frac{\sum x_i}{f}$$

$$x'_i = \frac{x_i - \mu}{SE[x_i]}$$

Residual networks help in learning identity functions.

Difference b/w Multi-head & Masked multi-head attention,



Stage I am a good boy → Main ek achcha baccha han

$$\begin{matrix} \text{Main} \\ \text{ek} \\ \text{achcha} \end{matrix} \left[\begin{matrix} \dots \\ \dots \\ \dots \end{matrix} \right] + \left[\begin{matrix} 1 & -\infty & -\infty & -\infty \\ 1 & 1 & -\infty & -\infty \\ 1 & 1 & 1 & -\infty \\ \dots & \dots & \dots & \dots \end{matrix} \right] \text{Mask}$$

Idea is to deplete attention of non-output words to $-\infty$ so that softmax is 0, this is masked multi-head attention. Attention/Importance of a word w.r.t ~~an~~ look-ahead word gets to 0.

Positional Encoding.

- Position of a word in a sentence is important, or I am a good boy = Am a I boy good
- Positional info is feeder of words to the model

$$PE_{(pos, 2i)} \in \mathbb{R}^{d_{\text{model}}}$$

position ↓ for even. for odd is $(2i+1)$

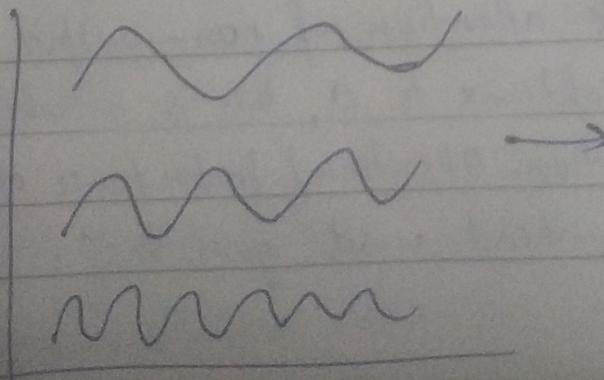
for I in "I am a good boy", pos = 1

$$PE_{(pos, 2i+1)} = \cos \left(pos / 10000^{2i+1/d_{\text{model}}} \right)$$

$$PE_{(pos, 2i)} = \sin \left(pos / 10000^{2i/d_{\text{model}}} \right)$$

$$(2i, 2i+1) \rightarrow (1, d_{\text{model}})$$

$$PE_{(1)} \text{ for } I = \begin{bmatrix} PE_{(1,1)} \\ PE_{(1,2)} \\ \vdots \\ PE_{(1,d_{\text{model}})} \end{bmatrix}$$



sin, cos wavy curves
encode positional
info into the model

* Ref. Bert : Pre-training of Deep Bi-directional
Transformers for Language Understanding

- Jacob Devlin et al.

DL-21/04

BERT : Bi-directional Encoder Representations using Transformers*

- 1) Masked Word Prediction
- 2) Next line Prediction

Autoencoders - slides

PCA aims to vanish variance along a certain no of dimensions so that the dependency on ~~that~~ those dimensions of the output ~~is~~ is eradicated. Thus, PCA reduces the dimensionality.