

$\frac{\partial g}{\partial x_p} = \frac{\partial g(x)}{\partial x_p}$

$g(x, y)$

$g(x, \dots, \dots)$

$g(x) = x^2$

$g'(x) = 2x$

$x \in f(x)$

$\frac{dg}{dx} = \lim_{\epsilon \rightarrow 0} \frac{g(x+\epsilon) - g(x-\epsilon)}{2\epsilon}$

Gradient Checking

$\frac{\partial g}{\partial x} > \frac{\partial g}{\partial y}$

$x \in g'(x)$ dated, close user limit

$x = 2$	$\frac{\partial g}{\partial x} = 4$	$\frac{\partial g}{\partial y} = 4$
$x = 2.01$	$\frac{\partial g}{\partial x} = 4.02$	$\frac{\partial g}{\partial y} = 4.02$
$x = 2.001$	$\frac{\partial g}{\partial x} = 4.002$	$\frac{\partial g}{\partial y} = 4.002$

$\lim_{\epsilon \rightarrow 0} \frac{g(x+\epsilon) - g(x-\epsilon)}{2\epsilon} = 4$

$f(x, y) = x + y$

$\frac{\partial f}{\partial x} = 1$

$f(x, y) = x * y$

$\frac{\partial f}{\partial x} = y$

$\frac{\partial f}{\partial y} = x$

$$\frac{\partial f}{\partial x} = 1$$

$f(x, y) = x/y = \frac{x}{y}$
 $\frac{\partial f}{\partial x} = \frac{1}{y} = \text{num} \cdot \frac{1}{\text{den}}$
 $\frac{\partial f}{\partial y} = -x \cdot \frac{1}{y^2} = -\text{num} \cdot \frac{1}{(\text{den})^2}$

$$f(x, y) = x - y$$

$\frac{\partial f}{\partial x} = 1$
 $\frac{\partial f}{\partial y} = -1$

$$f(x) = e^x$$

$\frac{df}{dx} = f$

$$f(x) = \log x$$

$\frac{df}{dx} = \frac{1}{x}$ /input

$$f(x) = \underline{x^2} = \underline{x * x}$$

$$\frac{df}{dx} = x$$

$$f(x) = \frac{\sin x}{\cos x}$$

$$f(x) = x^3 = x^2 * x = x * x * x$$

$$\frac{d}{dx} = x^2 + x^{-1} + x^3$$

$$= 3x^2$$

A hand-drawn diagram of a single neuron. It features a central circular node labeled "Sq". Two curved arrows point towards this node from the left, representing "2. input". A curved arrow points away from the right side of the node, representing the "x^2" output.

$$f(x) = \left(\frac{1-x}{1+e} \right)$$

$$\frac{\partial f}{\partial x} = 1 \times \cancel{1} \times L \times q \times (-1) = -\frac{e^{-x}}{(1+e^{-x})^2} = \frac{f(1-f)}{i}$$

$$\therefore 1 \times \left(-\frac{1}{e^x}\right) \times 1 \times q \times (-1) = \left(\frac{1}{1+e^x}\right) \left(1 - \frac{1}{1+e^x}\right)$$

The diagram illustrates a single neuron model with the following components:

- Forward Pass:** An input x is processed by weights $w_1 = -1$ and $w_2 = -x$. The bias is $b = 1$. The activation function is $f(x) = \frac{1}{1+e^{-x}}$.
- Backpropagation:** A target value $y_t = 1$ is compared with the output $y = f(x)$. The error is calculated as $\text{error} = y_t - y$. The derivative of the error with respect to the output is $\frac{\partial \text{error}}{\partial y}$.
- Gradients:** The gradient of the error with respect to the weights and bias is shown as $\frac{\partial \text{error}}{\partial w_1} = 1$ and $\frac{\partial \text{error}}{\partial b} = 1$.
- Sigmoid Function Plot:** A graph of the sigmoid function $f(x) = \frac{1}{1+e^{-x}}$ is plotted against x , showing its characteristic S-shape.

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

$$\lim_{x \rightarrow -\infty} f(x) = 0$$

$$\lim_{x \rightarrow \infty} f(x) = 1$$

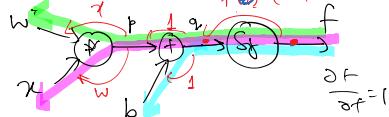
$$f'(x)|_{x=0} = \frac{1}{2}(1-\frac{1}{2}) = 0.25$$

$$f'(ax)|_{x=0} = ?$$

$$f'(ax) = \frac{a \cdot f(ax) \cdot (1-f(ax))}{1-a^2}$$

$$\frac{\partial f}{\partial x} = 1$$

$$f_{w,b}(x) = \frac{1}{1+e^{-(w^T x + b)}}$$



$$f = \text{Sig}(g)$$

$$\frac{\partial f}{\partial x} = 1$$

$$f'_b = 1 \times f \cdot (1-f) \times 1$$

$$f'_x = 1 \times f \cdot (1-f) \times 1 \times w$$

$$f'_w = 1 \times f \cdot (1-f) \times 1 \times x$$

$$f(w^T x + b) = \frac{1}{1+e^{-(w^T x + b)}} * \text{Sig}(w^T x)$$

$$g_1 = w_1 x + b$$

$$g_2 = w_2 x$$

$$(g_1 + g_2)x = (g_1 + g_2) \cdot w$$

$$\frac{\partial f}{\partial w} = \text{path}_3 + \text{path}_2$$

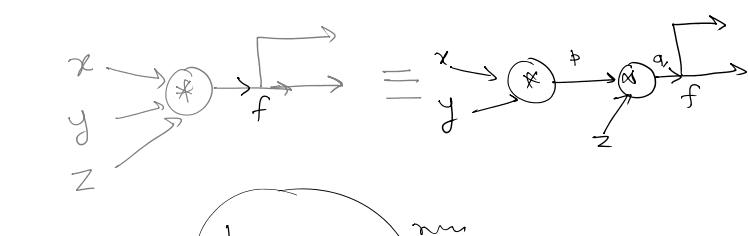
$$\frac{\partial f}{\partial x} = \text{path}_1 + \text{path}_4$$

$$\frac{\partial f}{\partial b} = \text{path}_5$$

$$\text{path 1: } 1 \times s \times r \times (a \circ (P)) \times w$$

$$\text{path 2: } 1 \times s \times r \times (a \circ (P)) \times x$$

$$\begin{cases} \text{path 3: } 1 \times s \times r \times (1-r) \times 1 \times x \\ \text{path 4: } 1 \times s \times r \times (1-r) \times 1 \times w \\ \text{path 5: } 1 \times s \times r \times (1-r) \times 1 \end{cases}$$



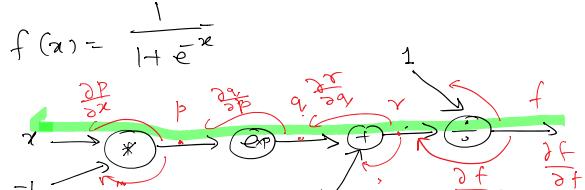
$$f(\omega^T x + b) = \frac{1}{1 + e^{-(\omega^T x + b)}}$$

$\frac{\partial f}{\partial \omega} = \frac{\partial f}{\partial \tau} \cdot \frac{\partial \tau}{\partial \omega}$

$\frac{\partial f}{\partial \tau} = \frac{\partial f}{\partial p} \cdot \frac{\partial p}{\partial \tau}$

$\frac{\partial p}{\partial \tau} = \frac{\partial p}{\partial x} \cdot \frac{\partial x}{\partial \tau}$

$\frac{\partial x}{\partial \tau} = r$



$$\frac{\partial f}{\partial x} = \underbrace{\frac{\partial f}{\partial t} \cdot \frac{\partial t}{\partial x}}_{\text{Chain rule}} \cdot \frac{\partial t}{\partial p} \cdot \frac{\partial p}{\partial x}$$

Class Test - 1: 3 + 2 points Mid Sem

$$f(x, y) = x^y \cdot \log(x+y)$$

$$\text{Compute } f_x'(x, y) = y \cdot x^{y-1} \log(x+y) + x^y \cdot \frac{1}{x+y}$$

$$f_y'(x, y) = x^y \cdot \ln(x) \cdot \log(x+y) + \frac{x^y}{x+y}$$

$f(x) = a^x$
$f'_x(x) = ?$

$a > 0$

Lim $a \rightarrow 0$

$a^x \cdot \ln(a)$ incorrect

Correlation

Sliding Vector
(filter / kernel)

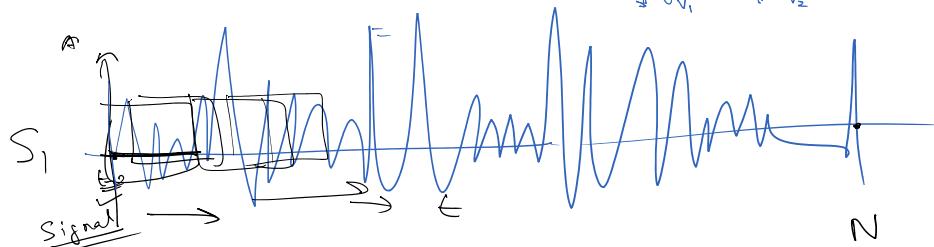
Base Vector Signal.

$$\|v_1\| = \|v_2\|$$

$$\text{Pearson Correl}(v_1, v_2) = \frac{1}{M} \frac{(v_1 - M_{v_1})}{\sqrt{\sigma_{v_1}^2}} \cdot \frac{(v_2 - M_{v_2})}{\sqrt{\sigma_{v_2}^2}} \frac{M_{v_1}}{\sqrt{\sigma_{v_1}^2}} \frac{M_{v_2}}{\sqrt{\sigma_{v_2}^2}}$$

Convolution

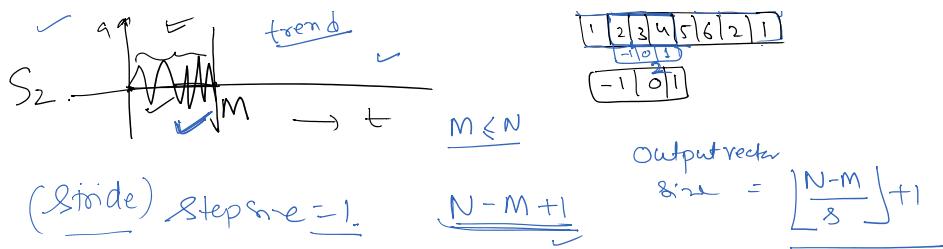
$$\begin{matrix} v_1 \\ v_2 \end{matrix} = \begin{matrix} 1 & 2 & 3 \\ 2 & 3 & 4 \\ 3 & 4 & 5 \\ 4 & 5 & 6 \\ 5 & 6 & 3 \\ 6 & 3 & 2 \end{matrix}$$



S_1 Signal

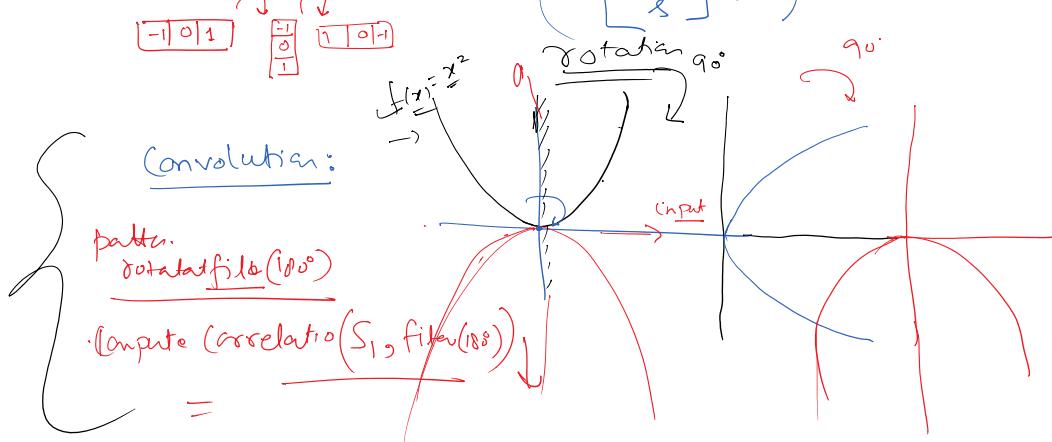
S_2 pattern trend match

2	-5
-1 0 1	-1 0 1
1 2 3 4 5 6 2 1	-1 0 1
-1 0 1	-1 0 1



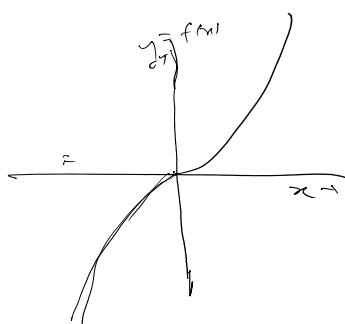
$$\text{Stride} = 2 \quad \left(\left\lfloor \frac{N-M}{2} \right\rfloor + 1 \right)$$

$$\text{Stride} = 8 \quad = \left(\left\lfloor \frac{N-M}{8} \right\rfloor + 1 \right)$$



mirror image ($x-y$ plane \rightarrow x -axis is the mirror)

$$y = f(x) \xrightarrow{\text{rotate } 180^\circ \text{ in } y} -f(x)$$



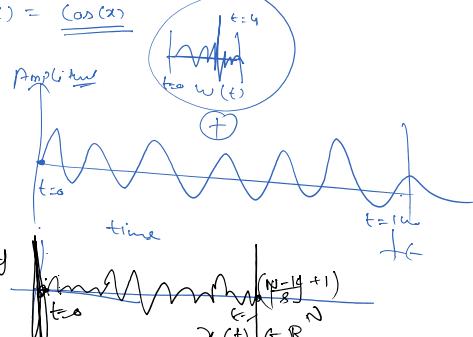
$$\begin{aligned}
 f(x) &= x^3 \\
 f(-x) &= -f(x) \\
 f(x) &\xrightarrow{\text{mirror } y \text{-axis}} f(-x) \\
 &= f(x) = x^2 \Rightarrow f(-x) = (-x)^2 \\
 f(x) &= \log(x) \Rightarrow f(-x) = \log(-x) / x
 \end{aligned}$$

$$\begin{aligned}
 f(x) &= \sin(x) : f(-x) = \sin(-x) = -\sin x \\
 f(x) &= \cos(x) : f(-x) = \cos(-x) = \cos x
 \end{aligned}$$

$$\xrightarrow{\text{Signal}} f(t) \xrightarrow{\text{Filter}} w(t)$$

$$N = 100, m = 5$$

$$\begin{aligned}
 \text{Stride} &= 1 \\
 (\text{Correl} \Rightarrow f(t) \cdot w(t)) & \\
 (\text{Conv} \Rightarrow f(t) * w(-t))
 \end{aligned}$$



$$\begin{aligned}
 \xrightarrow{\text{ID:}} x(t) &= \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 & \dots & x_N \end{bmatrix} \\
 f(t) &= \begin{bmatrix} f_1 & f_2 & f_3 & f_4 & f_5 & \dots & f_N \end{bmatrix}
 \end{aligned}$$

$$y = R^{N-m+1}$$

$$s = 1$$

$$y = \text{Gradient } f(t), f(t)$$

$$\begin{aligned}
 \frac{\partial y[0]}{\partial f_i} &= \frac{x_0 \ x_1 \ x_2 \ x_3 \ x_4}{x_0 \ x_1 \ x_2 \ x_3 \ x_4} y[0] = x \Big|_{t=0} T \Big|_{t=k-1} \cdot f(t) \Big|_{t=0} T \Big|_{t=k-1} \\
 &= \begin{bmatrix} x_1 & x_2 & x_3 & x_4 \end{bmatrix} = x_1 \cdot f_1 + x_2 \cdot f_2 + x_3 \cdot f_3 + x_4 \cdot f_4
 \end{aligned}$$

$$y[1] = x \Big|_{t=1} T \Big|_{t=k-1} f(t) \Big|_{t=0} T \Big|_{t=k-1} =$$

$$\frac{\partial y[i]}{\partial f} = \boxed{x_2 | x_3 | x_4 | x_5}$$

$$y[i] = x \left|_{t=1 \text{ to } t=k} \right. f(t) \Big|_{t=0 \text{ to } k-1} =$$

$$\boxed{x_2 | x_3 | x_4 | x_5} \quad \boxed{f_1 | f_2 | f_3 | f_4}$$

$$= x_2 \cdot f_1 + x_3 \cdot f_2 + x_4 \cdot f_3 + x_5 \cdot f_4$$

$$y[2] = x \left|_{t=2 \text{ to } t=k+1} \right. f(t) \Big|_{t=0 \text{ to } k-1}$$

$$\frac{\partial y[i^*]}{\partial f} \Big|_{t=0 \text{ to } k-1} = \boxed{x_{i+1} | x_{i+2} | \dots | x_{i+k-1}}$$

$$y[i^*] = x \left|_{t=i \text{ to } t=i+k-1} \right. f(t) \Big|_{t=0 \text{ to } k-1}$$

$$= \boxed{x_{i+1} | x_{i+2} | x_{i+3} | x_{i+4}} \quad \boxed{f_1 | f_2 | f_3 | f_4} =$$

$$y[1^*] = x_{i+1} \cdot f_1 + x_{i+2} \cdot f_2 + x_{i+3} \cdot f_3 + x_{i+4} \cdot f_4$$

$$y[1^*] = \sum_{j=1}^k x_{i+j} \cdot f_j$$

$$y = \text{corr}(x, f) \quad 0 \leq i \leq \lfloor \frac{N-k}{s} \rfloor + 1$$

$$f(t) = \boxed{f_1 | f_2 | f_3 | f_4} \Rightarrow f(-t) =$$

$$= \boxed{f_4 | f_3 | f_2 | f_1}$$

Convolution: $y[i^*] = \sum_{j=1}^k x_{i-j} \cdot f_j$

$$\text{Cov}(x(t) \cdot f(t)) = \text{corr}(x(t), f(-t))$$

Jacobian

$$\frac{\partial y}{\partial x} = ? \quad \frac{\partial y}{\partial f} = ?$$

$$\frac{\partial y}{\partial x} = \boxed{\frac{\partial y}{\partial x_1} \quad \dots \quad \frac{\partial y}{\partial x_{N-k+1}}} \quad + \quad \boxed{\begin{array}{c|ccccc} f_1 & x_1 & x_2 & x_3 & x_4 & f_4 \\ \hline x_1 & x_2 & x_3 & x_4 & x_5 & \\ x_2 & x_3 & x_4 & x_5 & -x_1 & \\ x_3 & x_4 & x_5 & -x_2 & & \\ \hline x_{N-k+1} & x_{N-k+2} & \dots & x_N & (N-k+1) \times k \end{array}} \quad \frac{\partial y}{\partial f} = ?$$

$$\frac{\partial y}{\partial x} = ? \quad \frac{\partial y}{\partial f} = ? \quad \boxed{(N-k+1) \times N}$$

$$\frac{\partial y[0]}{\partial x_1} = f_1 ; \quad \frac{\partial(y[0])}{\partial x_2} = f_2 ; \quad \frac{\partial(y[0])}{\partial x_u} = f_u = \frac{\partial(y[0])}{\partial x_s} = 0$$

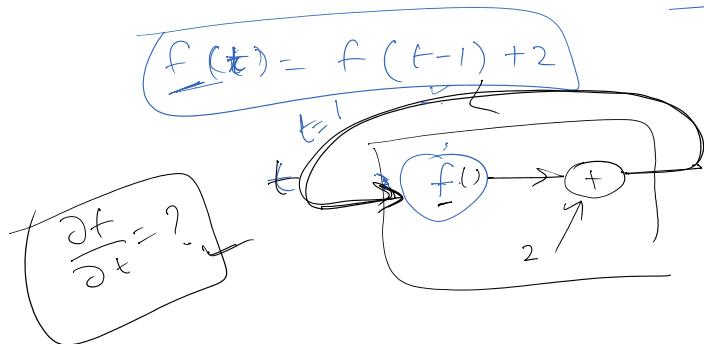
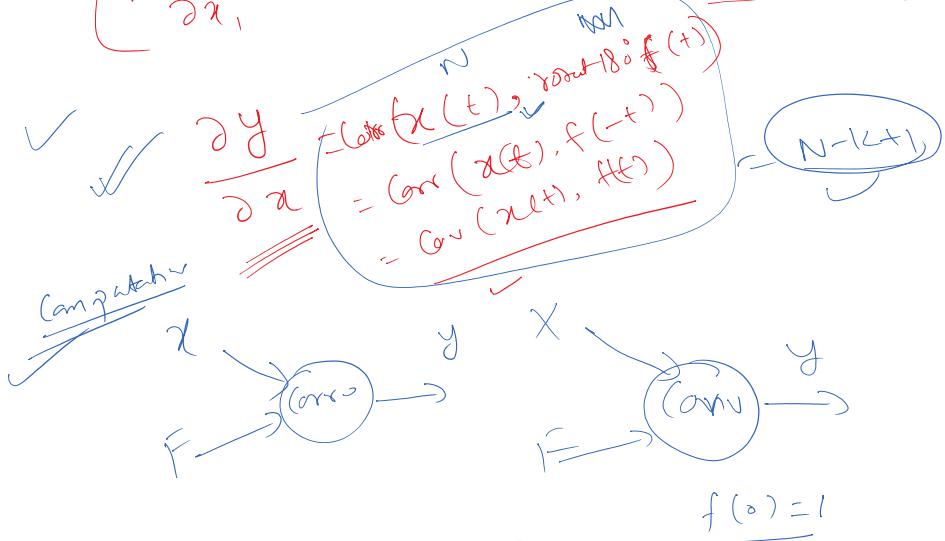
$$\frac{\partial(y[1])}{\partial x_1} = 0 ; \quad \frac{\partial(y[1])}{\partial x_2} = f_1 ; \quad \dots ; \quad \frac{\partial(y[1])}{\partial x_5} = f_4 ; \quad \frac{\partial(y[1])}{\partial x_6} = 0 ; \quad \frac{\partial(y[1])}{\partial x_N} = 0$$

$$\frac{\partial(y[2])}{\partial x_1} = 0 ; \quad \frac{\partial(y[2])}{\partial x_2} = 0 ; \quad \frac{\partial(y[2])}{\partial x_3} = f_1 ; \quad \dots ; \quad \frac{\partial(y[2])}{\partial x_N} = 0$$

$$\vdots$$

$$\frac{\partial(y[N-k+1])}{\partial x_1} = 0 ; \quad \dots ; \quad = 0 \quad \frac{\partial(y[N-k+1])}{\partial x_{N-k+1}} = 0 ; \quad \frac{\partial(y[N-k+1])}{\partial x_N} = f_4$$

$$\frac{\partial y_{(N-l+1)}}{\partial x_1} = 0 \dots = 0 \Rightarrow \frac{\partial y_{(N-k+1)}}{\partial x_{N+k+1}} = \frac{\partial y_{(N-l+1)}}{\partial x_N} = F_4$$



\rightarrow Convolution: (Deep learning)

$\frac{\partial y}{\partial x}$

$\frac{\partial y}{\partial w}$

$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w}$

$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial x}$

$x \in \mathbb{R}^N$

$w \in \mathbb{R}^{N-M+1}$

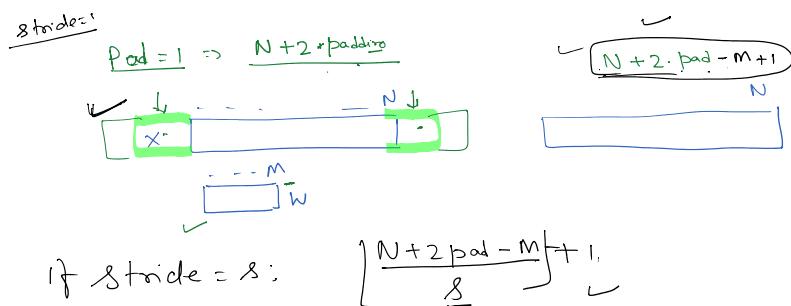
$y \in \mathbb{R}$

stride = 1

\checkmark Edges are the basic building

\checkmark Valid Convolution (Stride=1) = $N - M + 1$

Same Convolution \Rightarrow output vector should have same size as input signal vector.



Stride > 1

$N + 2 \cdot \text{pad} - M + 1 = N$

$\text{pad} = \lceil \frac{M-1}{2} \rceil$

: odd size
kernels

Stride = 8

$\lceil \frac{M-1}{2} \rceil + 1$

$r = 1, 2, \dots$ kernels

stride = 8

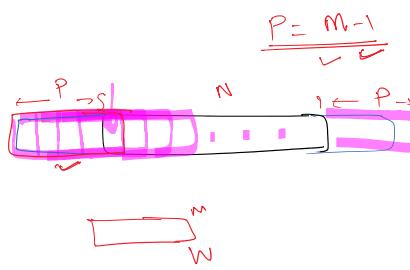
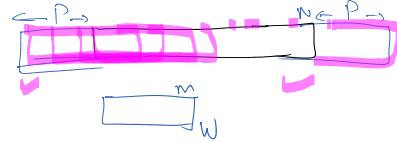
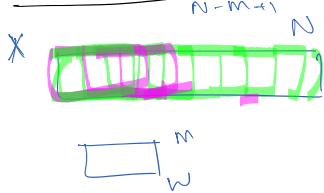
$$\left\lceil \frac{N+2\text{Pad}-M}{s} \right\rceil + 1 = N$$

$$\text{Pad} = \frac{s(N-1) + M - N}{2}$$

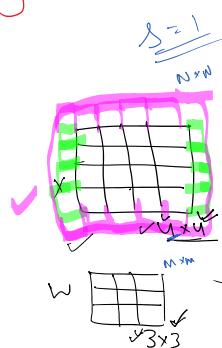
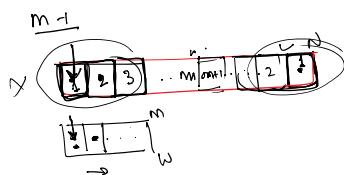
$$P = \frac{N(s-1) - s + M}{2}$$

$$P = \frac{M-1}{2}$$

Full convolution:



$$Y \in \mathbb{R}^{N+m-1}$$



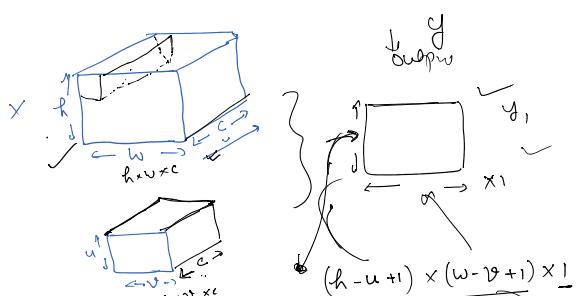
$$P = \frac{m-1}{2}, \frac{m+1}{2}$$

Valid Same Full

$$(u-3+1) \times (u-3+1) \quad 4 \times 4 \quad 6 \times 6$$

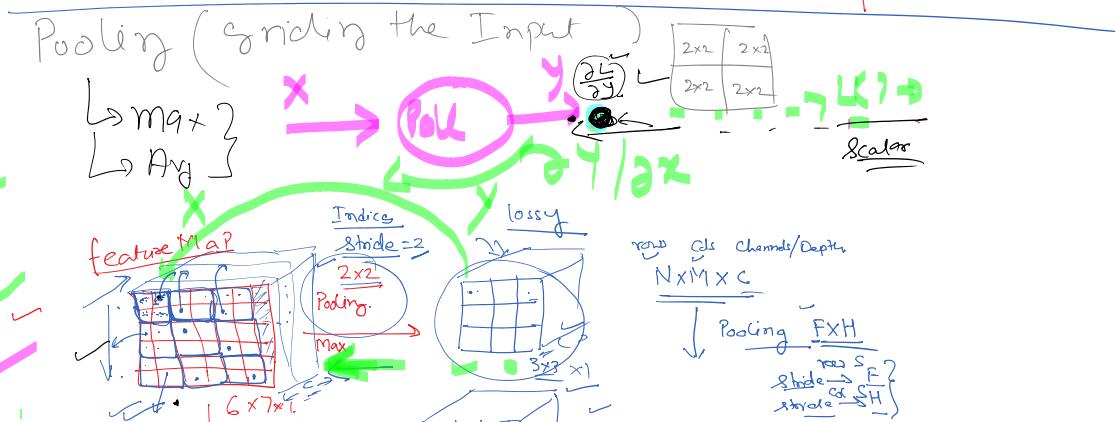
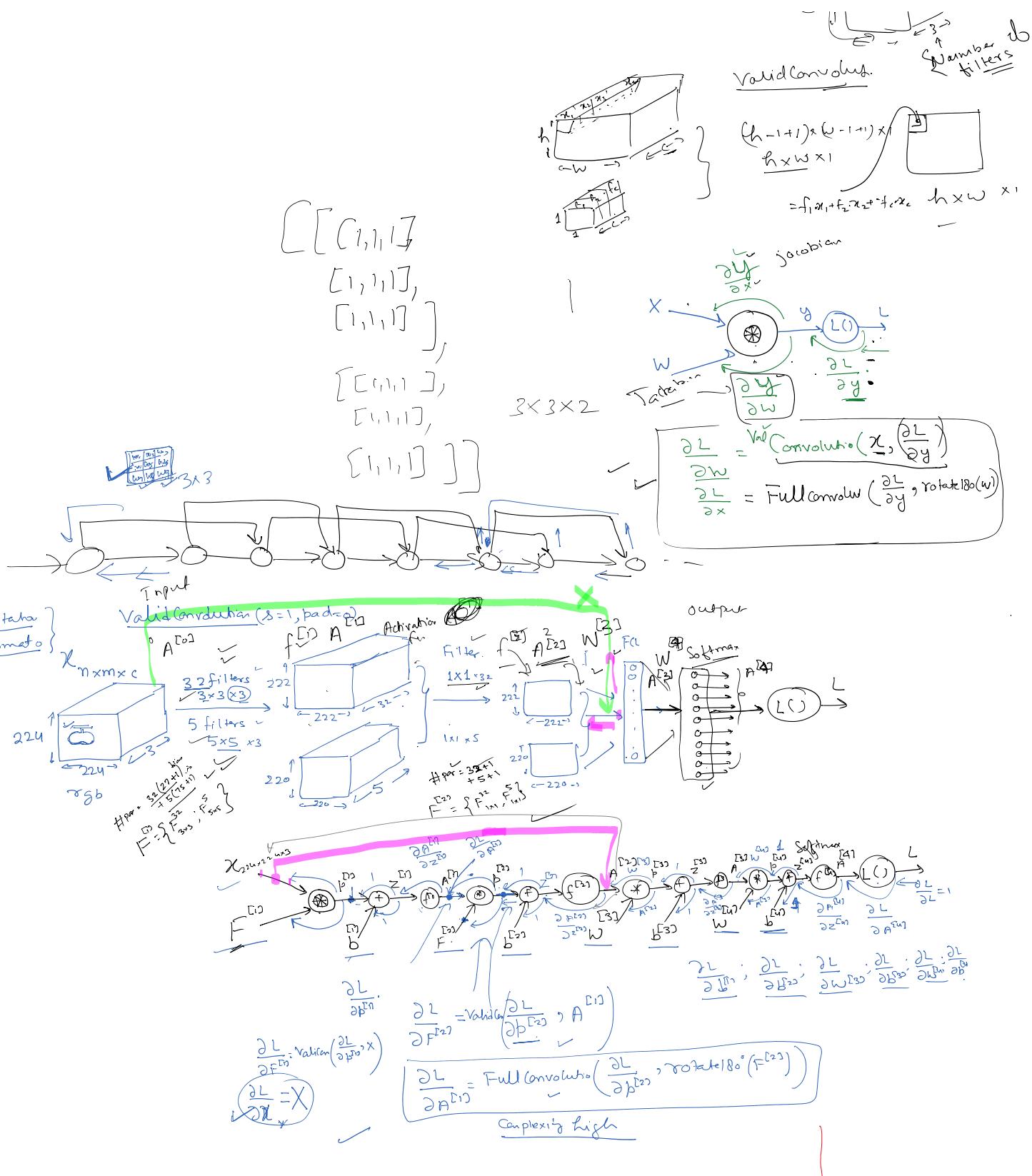
$$(N-m+1) \times (N-m+1) \quad N \times N \quad (N+m-1) \times (N+m-1)$$

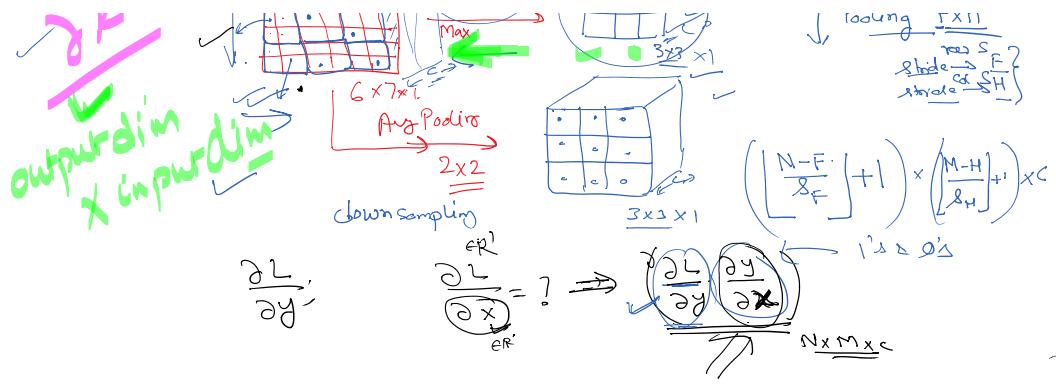
$$\frac{(h_x+2p)}{s} \times \frac{(w_x+2p)}{s} \quad h_x \times w_x \quad \frac{(h_x+h_f-1)}{(w_x-w_f+1)} \times \frac{(h_x+w_f-1)}{(w_x+w_f-1)}$$



$$N \approx M$$







ReLU

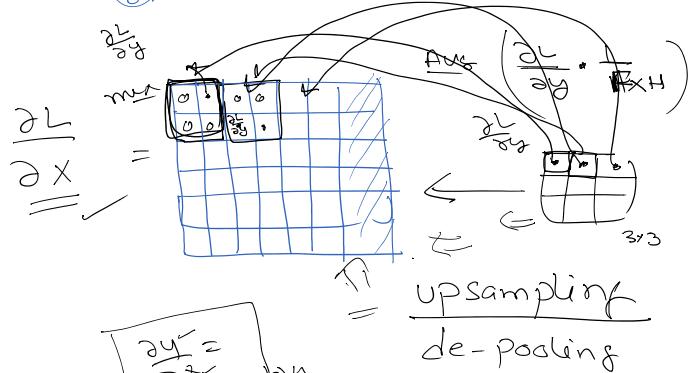
$$f(x) = \max(x, 0) \Rightarrow x \in \mathbb{R}$$

$$\frac{\partial f}{\partial x} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \\ \text{undefined if } x = 0 \end{cases}$$

$$f(x, y) = \max(x, y)$$

$$\frac{\partial f}{\partial x} = \begin{cases} 1 & \text{if } x > y \text{ otherwise } 0 \end{cases}$$

$$\frac{\partial f}{\partial y} = \begin{cases} 1 & \text{if } y > x \text{ otherwise } 0 \end{cases}$$



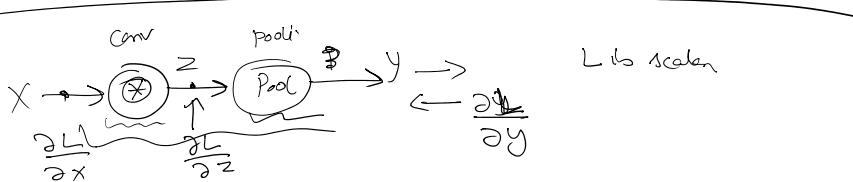
$$\frac{\partial y}{\partial z} = \frac{\partial y}{\partial t} \text{ Jacobian}$$

$$\frac{\partial L}{\partial x} = \dim \frac{\partial x}{\partial z}$$

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial y} \right) \cdot \left(\frac{\partial y}{\partial z} \right) \cdot \left(\frac{\partial z}{\partial x} \right)$$

$$\left(\frac{N-F}{S} + 1 \right) \times \left(\frac{M-H}{S} + 1 \right) \times \left(\frac{I \times \dim y}{\dim x} \right)$$

$$= \frac{I \times \dim y}{\dim x}$$



$\sum 1 \dots 10 \Rightarrow \frac{100 \times 101}{2} =$

$\sum 1 \dots 10 = 101$

$\sum 1 \dots 10 = 101$

$\sum 1 \dots 10 = 101$

$5 \times 10 \quad 1 \times 5$



$$f(x) = a \cdot f(x-1) + b \quad (\text{Recurrence Eqn})$$

Unrolling

x is discrete time $a \in \mathbb{R}$ $b \in \mathbb{R}$

$$\checkmark$$

$$f(x) = a \cdot (a \cdot f(x-2) + b) + b$$

$$= a \cdot (a \cdot (a \cdot f(x-3) + b) + b) + b$$

\vdots

$$= a^n \cdot f(x-n) + n \cdot b$$

$f(x-n)$ is independent \times

$$f'(x) = \frac{\partial f}{\partial x} = a^n \cdot f'(x-n)$$

$$f(x) = a \cdot f(x-1) + b$$

$$\checkmark$$

$$f'(x) = a \cdot \left(\frac{\partial f(x-1)}{\partial x} \right) \stackrel{x \rightarrow 0}{\rightarrow}$$

$$= a \cdot a \left(\frac{\partial f(x-2)}{\partial x} \right)$$

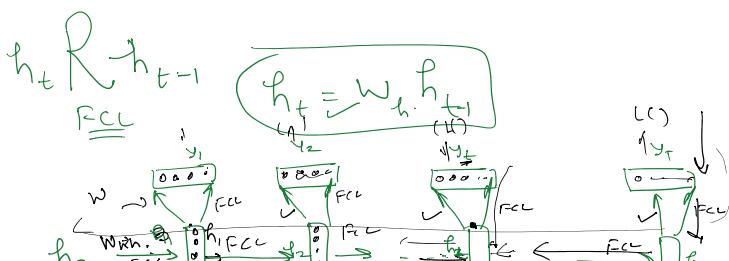
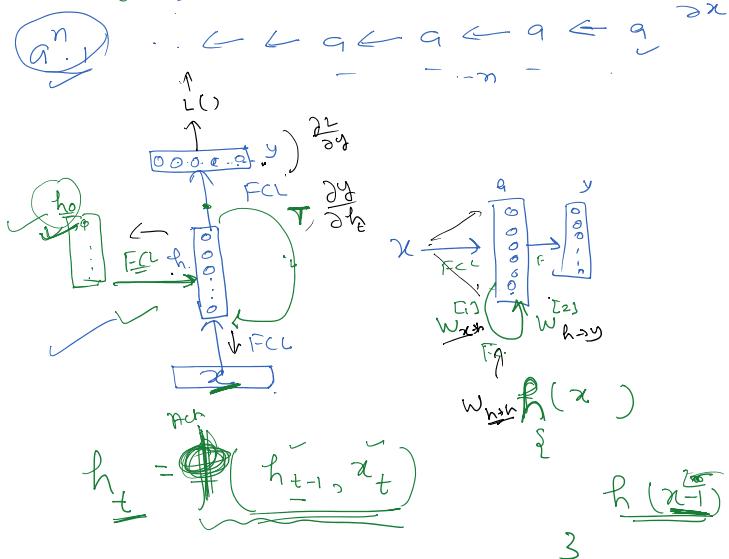
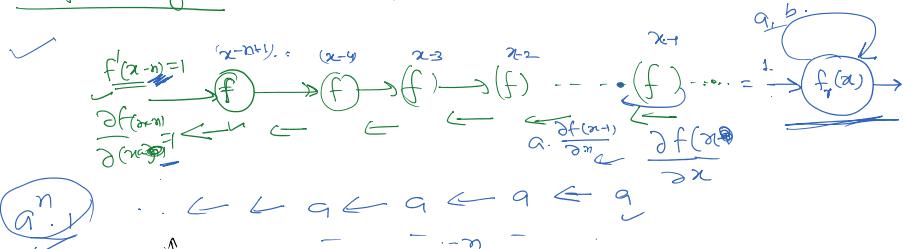
$n \downarrow$

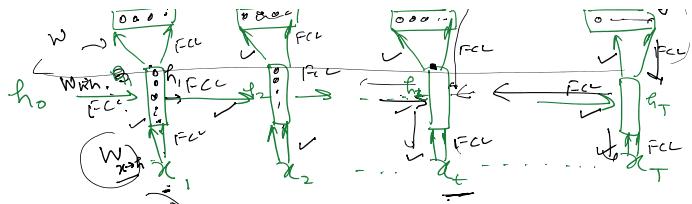
$$= a^n$$

$$f'(x-n) = 1$$

$$\left(\frac{\partial f(x)}{\partial f(x-1)} \right)$$

Unfolding Recurrence Eqn





$$f_1, f_2 = \text{(sigmoid, tanh)}$$

$$h_t = f_1(w_{x \rightarrow h} \cdot x_t + w_{h \rightarrow h} \cdot h_{t-1} + b_h)$$

$$h_t = f_2(w_{h \rightarrow h} \cdot h_{t-1} + b_h)$$

\Rightarrow Forward Computation Graph

$$\checkmark \frac{\partial h_t}{\partial w_{h \rightarrow h}} = f_1'(z_t) \cdot \frac{\partial z_t}{\partial w_{h \rightarrow h}}$$

=

$$z_t = w_{x \rightarrow h} \cdot x_t + w_{h \rightarrow h} \cdot h_{t-1} + b_h$$

$$\frac{\partial z_t}{\partial w_{h \rightarrow h}} = \frac{\partial (w_{x \rightarrow h} \cdot x_t)}{\partial w_{h \rightarrow h}} + \frac{\partial (w_{h \rightarrow h} \cdot h_{t-1})}{\partial w_{h \rightarrow h}} + \frac{\partial b_h}{\partial w_{h \rightarrow h}}$$

$$= \left(\frac{\partial w_{x \rightarrow h} \cdot x_t}{\partial w_{h \rightarrow h}} + \frac{\partial w_{h \rightarrow h} \cdot h_{t-1}}{\partial w_{h \rightarrow h}} \right) + \left(w_{h \rightarrow h} \frac{\partial h_{t-1}}{\partial w_{h \rightarrow h}} + \frac{\partial b_h}{\partial w_{h \rightarrow h}} \right)$$

$$= 0 + 0 + w_{h \rightarrow h} \cdot \frac{\partial h_{t-1}}{\partial w_{h \rightarrow h}} + 1 \cdot h_{t-1}$$

(Note: $h_{t-1} = f_1(w_{x \rightarrow h} \cdot x_{t-1} + w_{h \rightarrow h} \cdot h_{t-2} + b_h)$)

$$\frac{\partial z_t}{\partial w_{h \rightarrow h}} = h_{t-1} + w_{h \rightarrow h} \cdot \left(\frac{\partial h_{t-1}}{\partial w_{h \rightarrow h}} \right)$$

$$\frac{\partial z_{t-1}}{\partial w_{h \rightarrow h}} = h_{t-2} + w_{h \rightarrow h} \frac{\partial h_{t-2}}{\partial w_{h \rightarrow h}}$$

$$h_1 = f_1(w_{x \rightarrow h} \cdot x_1 + w_{h \rightarrow h} \cdot h_0 + b_h)$$

$\underline{z_1}$

$\frac{\partial z_1}{\partial w_{h \rightarrow h}} = h_0$

$f(z) = z$
 $\frac{\partial f}{\partial z} = 1$ linear ≤ 1

$$\frac{\partial h_t}{\partial w_{n \rightarrow h}} = f_1'(z_t) \cdot \left(h_{t-1} + w_{n \rightarrow h} \left(f_1'(z_{t-1}) \cdot \left(h_{t-2} + w_{n \rightarrow h} \left(f_1'(z_{t-2}) \cdot \left(h_{t-3} + w_{n \rightarrow h} \right) \right) \right) \right) \right)$$

⋮

$$= \dots - f_1'(z_2) \left(h_1 + w_{n \rightarrow h} \left(f_1'(z_1) \cdot (h_0) \right) \right) \dots$$

↙

$$\frac{\partial h_t}{\partial w_{n \rightarrow h}} = (\dots ? \dots)$$

$$\frac{\partial y_t}{\partial h_t} = f_2'(\quad) \cdot w_{h \rightarrow y}$$

$$\frac{\partial y_t}{\partial w_{h \rightarrow y}} = f_2'(\quad) \cdot h_t$$

↙



Computational Graph

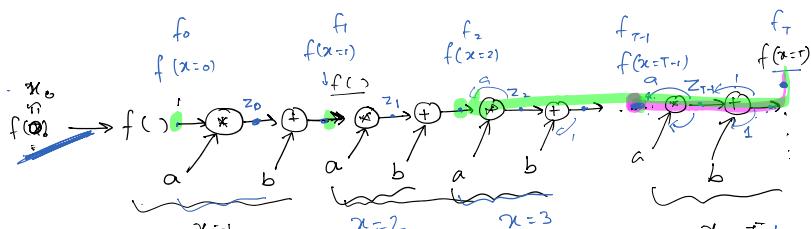
$$\underbrace{P(y_T | x_1, \dots, x_T)}_{RNN}$$

$$f(x) = \dots$$

$$f(x) = a \cdot f(x-1) + b$$

$x=3$:

$$\begin{aligned} f(3) &= (a^3 \cdot x_0 + 3b) \\ a \cdot f(2) + b &= a^2 \cdot x_0 + 2b \\ a \cdot f(1) + b &= a \cdot x_0 + b \\ a \cdot & \end{aligned}$$

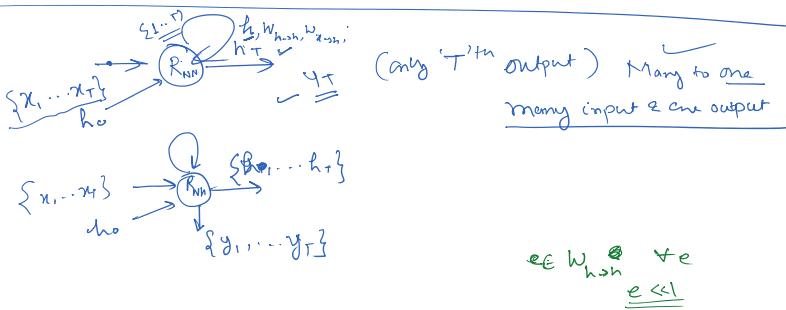


$$\frac{\partial f_T}{\partial z_{T-1}} = 1 \cdot \frac{\partial z_{T-1}}{\partial f_{T-1}} = a \cdot \frac{\partial f_{T-1}}{\partial f_{T-1}} = 1 \cdot a$$

$$\frac{\partial f_T}{\partial f_2} = (1 \cdot a) \times (1 \cdot a) \times (1 \cdot a) \cdots \underset{(T-2) \text{ times}}{\dots} = \underline{\underline{a^{T-2}}}$$

$$\frac{\partial f_T}{\partial f_1} = \underline{\underline{a^{T-1}}}$$

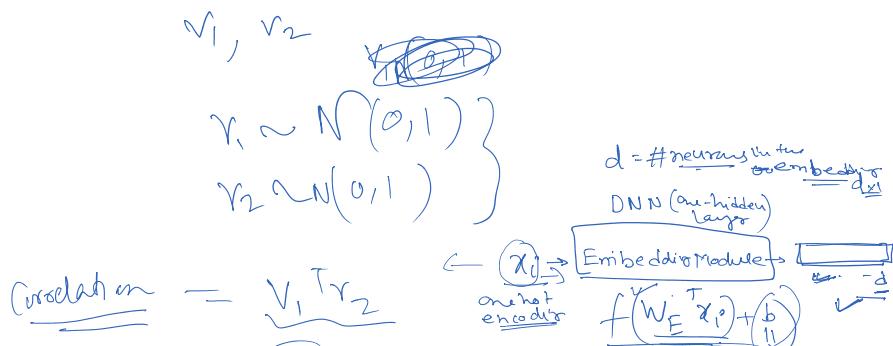
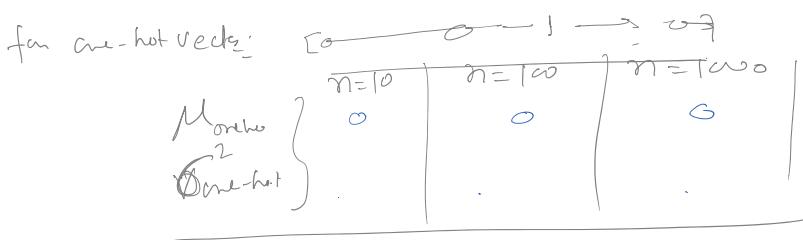
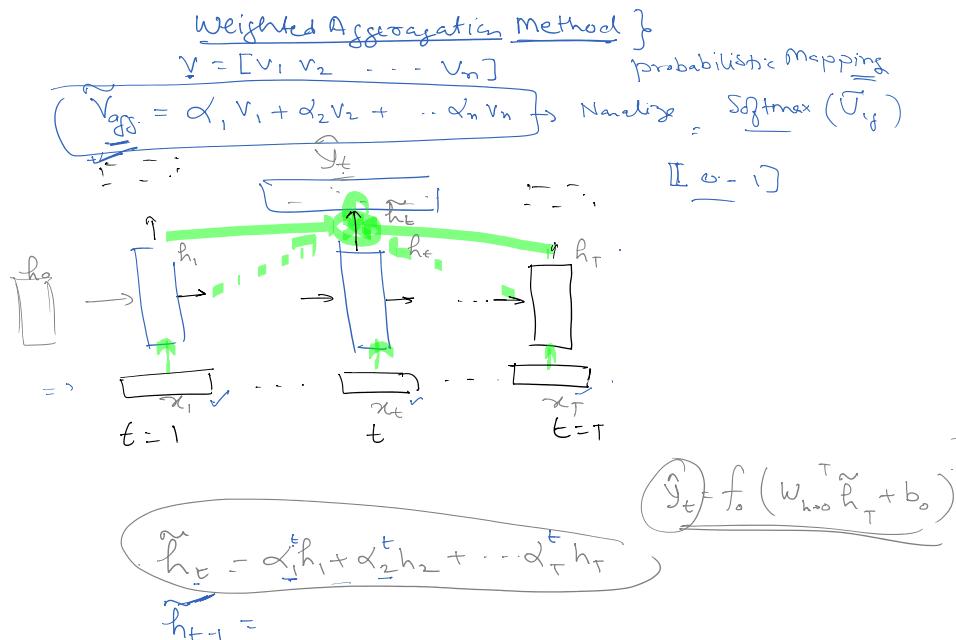
$$\frac{\partial f_T}{\partial f_0} = \underline{\underline{a^T}} \cdot \underline{\underline{?}}$$



$$\begin{aligned}\frac{\partial h_t}{\partial \mathbb{W}_{h \rightarrow h}} &= h_{t-1} + \underbrace{W_{h \rightarrow h}(h_{t-2} + W_{h \rightarrow h}(h_{t-3} + \dots))}_{W_{h \rightarrow h}(h_0)} \\ &= h_{t-1} + W_{h \rightarrow h} \cdot h_{t-2} + \underbrace{(W_{h \rightarrow h}) W_{h \rightarrow h} \cdot h_{t-3} + \dots}_{(W_{h \rightarrow h})^T h_0}\end{aligned}$$

→ RNN, LSTM, GRU
Bi-directional ; //

✓ Attention Mechanism ↴



$$\text{Correlation} = \sqrt{\chi_1^T \chi_2}$$

one-hot
 encoder $\xrightarrow{\chi_i \rightarrow \text{Embedding vectors}}$ $f(\chi_i^T \chi_i + b)$
 $\forall i=1 \dots T$

$\chi_i = \text{Embedding Vector}$

$$S_1, S_2, \dots, S_N \leftarrow \langle x_0, x_1, x_2, x_3, x_4, \dots, x_T \rangle$$

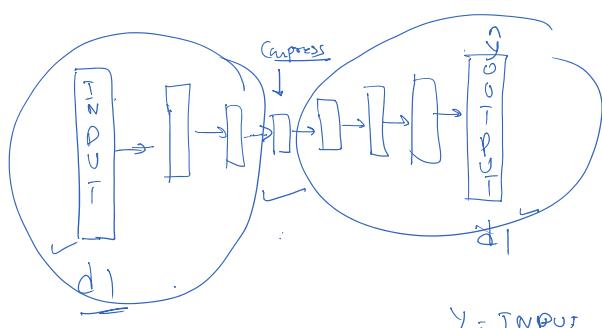
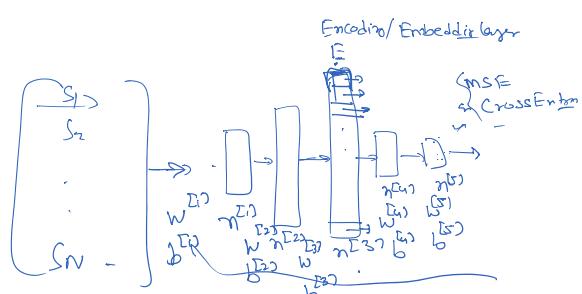
$x_i = [x_{i0}, x_{i1}, x_{i2}, x_{i3}, x_{i4}]$

$\rightarrow \chi_i = \text{Embedding Vector}$

$$\text{Similarity} = \chi_i^T \chi_j$$

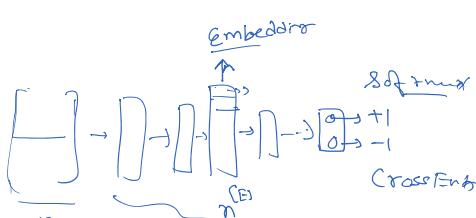
$$S_1 \rightarrow \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \end{bmatrix}_{n \times 30}$$

$$S_2 \rightarrow \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \end{bmatrix}_{n \times 30}$$



$$S_1 \leftarrow \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \end{bmatrix}_{n \times 2^n}$$

$\checkmark \langle x_3, x_{100} \rangle + 1$
 $\checkmark \langle x_3, x_4 \rangle + 1$
 $\checkmark \langle x_3, x_{300} \rangle - 1$
 $\checkmark \langle x_3, x_{700} \rangle - 1$



$\rightarrow \int_{\text{tokens}} (\text{Word in text, Value at certain time steps, frames, } \dots)$

\rightarrow token (word in text, value at certain time steps, frames, ...)

token \rightarrow Embedding Vector: $\underbrace{[\quad]}_{1 \times d}$ Embedding dimension

Positional Encoding + Attention

text: $\langle x_1, x_2, \dots, x_t, \dots, x_T \rangle$

$x_1 \rightarrow E_{x_1} [\quad]_{1 \times d}$? Position b_1 $\stackrel{\text{binary}}{=} [\quad]_{d+1}$ $\stackrel{\text{Sinusoids}}{=} [\quad]_{1 \times d}$

$x_2 \rightarrow E_{x_2} [\quad]_{1 \times d}$? b_2 $[\quad]_{d+1}$ $[\quad]_{1 \times d}$

$x_t \rightarrow E_{x_t} [\quad]_{1 \times d}$? b_t $[\quad]_{d+1}$ $[\quad]_{1 \times d}$

position in Position Vector $\Rightarrow c = 0 \text{ to } \frac{d}{2}$

Even i : $2i$ $2i+1$ $\sin\left(\frac{\text{position} \cdot \text{token}}{S^{2i}/\text{embedding dim of the token}}\right)$

Odd i : $2i+1$ $2i$ $\cos\left(\frac{\text{position} \cdot \text{token}}{S^{2i}/\text{embedding dim of the token}}\right)$

$S = 10000$

$x_1 = [\quad]_{1 \times 4}$

$x_2 = [\quad]_{1 \times 4}$

$x_3 = [\quad]_{1 \times 4}$

$x_6 = [\quad]_{1 \times 4}$

Position Emb: $\begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 0 & 1 \end{bmatrix}$

$i=0, \text{pos}=0$
 $\sin\left(\frac{0}{10000}\right) = 0$
 $\cos\left(\frac{0}{10000}\right) = 1$

$i=1, \text{pos}=0$
 $\sin\left(\frac{0}{10000}\right) = 0$
 $\cos\left(\frac{0}{10000}\right) = 1$

$i=0, \text{pos}=1$
 $\sin\left(\frac{1}{10000}\right) = \underline{\underline{\sin(1)}}$
 $\cos\left(\frac{1}{10000}\right) = \underline{\underline{\cos(1)}}$

$i=1, \text{pos}=1$
 $\sin\left(\frac{1}{10000^2}\right) = \sin\left(\frac{1}{10000}\right)$
 $\cos\left(\frac{1}{10000^2}\right) = \cos\left(\frac{1}{10000}\right)$

Assignment:

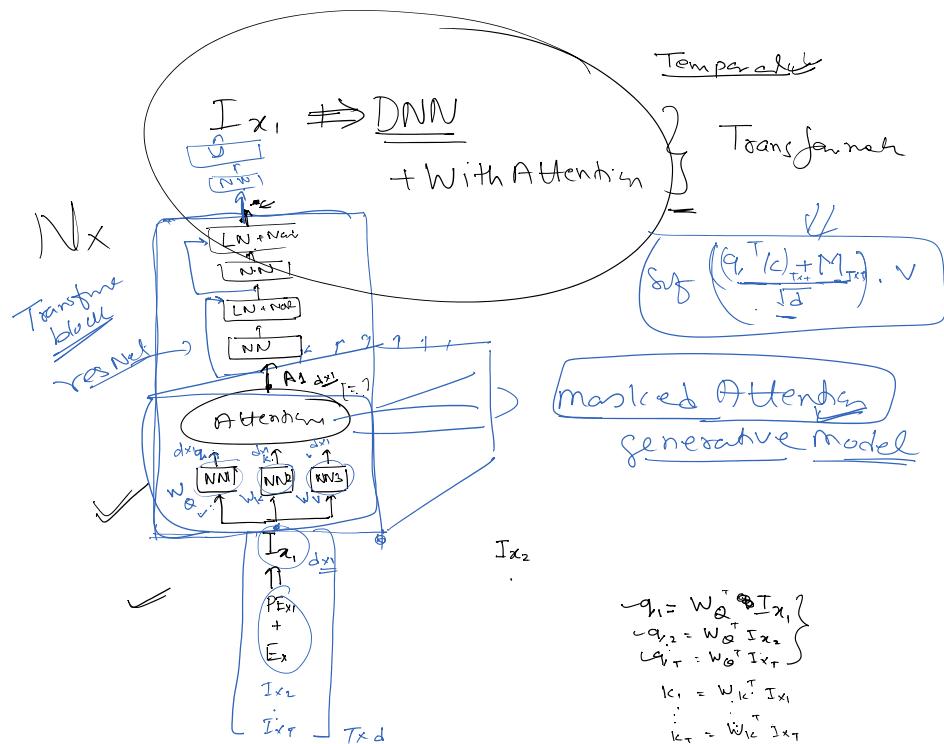
for a fixed Embedding dim d)

What is the maximum position that can be uniquely encoded in $\underline{\underline{P}}$?

Visualize by varying $\underline{\underline{d}}$

$$I_{x_1} = E_{x_1} + PE_{x_1} \rightarrow \text{Semantic} + \text{Position}$$

\rightarrow Simple Classification (Sentimental)
DNN



Self semantic (or Self Attention)

$$q_i = W_Q^T I_{x_1}$$

$$k_i = W_K^T I_{x_1}$$

$$v_i = W_V^T I_{x_1}$$

$$q_{i,2} = W_Q^T I_{x_2}$$

$$k_{i,2} = W_K^T I_{x_2}$$

$$v_{i,2} = W_V^T I_{x_2}$$

$$q_{i,T} = W_Q^T I_{x_T}$$

$$k_{i,T} = W_K^T I_{x_T}$$

$$v_{i,T} = W_V^T I_{x_T}$$

Q	K	V
$q_{i,1} = [1 \ 0 \ 0 \ \dots]$ $q_{i,2} = [0 \ 1 \ 0 \dots]$ $q_{i,3} = [0 \ 0 \ 1 \dots]$ $q_{i,n} = [0 \ 0 \ 0 \dots]$	$k_{i,1} = [1 \ 1 \ 0 \ 0 \ 0 \dots]$ $k_{i,2} = [0 \ 1 \ 0 \ 0 \ 0 \dots]$ $k_{i,3} = [0 \ 0 \ 1 \ 0 \ 0 \dots]$ $k_{i,n} = [0 \ 0 \ 0 \ 1 \ 0 \dots]$	$v_{i,1} = [1 \ 0 \ 1 \ 0 \ 0 \dots]$ $v_{i,2} = 0 \ 1 \ 0 \ 0 \ 0 \dots$ $v_{i,3} = 0 \ 0 \ 1 \ 0 \ 0 \dots$ $v_{i,n} = 0 \ 0 \ 0 \ 1 \ 0 \dots$

First token (I_{x_1}) = atten vector ($\underline{\alpha}$)

$$\underline{\alpha} = q_{i,1} \Rightarrow q_{i,1}^T k_1 \rightarrow 1 \times v_1 \Rightarrow v_1$$

$$q_{i,1}^T k_2 \rightarrow 0$$

$$q_{i,1}^T k_3 \rightarrow 0$$

$$q_{i,1}^T k_n \rightarrow 0$$

$$\underline{\alpha} : q_{i,2} \rightarrow q_{i,2}^T k_1 \rightarrow 0$$

$$q_{i,2}^T k_2 \rightarrow 1 \rightarrow v_2 \dots$$

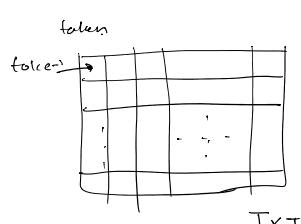
$$q_{i,2}^T k_3 \rightarrow 0$$

$$q_{i,2}^T k_n \rightarrow 0$$

$$A_{i,i} = \text{softmax}(q_{i,1}^T k_i) \cdot v_i$$

$$A_{i,j} = \text{softmax}(q_{i,1}^T k_j) \cdot v_i$$

$$A_{i,i} = \text{softmax}(q_{i,1}^T k_i) \cdot v_i$$



atten = $\text{softmax}(q_{i,1}^T k_i) \cdot v_i$

$$S1 < t_1 + \tau$$

$$S1 = \langle \overline{t_1}, \overline{t_2}, \overline{t_3} \rangle$$

Cross Attention



