# Deep Learning

Last Class:-

① Motivation for deep learning ⎰ scalability
                                  ⎱ compositional

② (a) CNN :- Convolutional Neural Networks

(b) RNN :- Recurrent Neural Networks

(c) GAN :- Generative Adversarial Networks

⎱ ⎰ - ~~self~~ Attention
  ⎱ - Dropouts

## Convolutional Neural Networks :-

Apply a filter

$K^2 d$

$x[i, j, e]$

$w[u, v, e]$

$k \times k$

$w$

$h$

Input :- 3-Dim

$\rightarrow$ A set of operations

$Z_o[i, o]$ output feature map.

filter
$\Uparrow$

$$Z[i, o] = \sum_{l=0}^{d-1} \sum_{u=0}^{k-1} \sum_{v=0}^{k-1} x[i+u, o+v, l] \cdot W[u, v, l] + b$$

tied across

Assuming stride $s=1$
Resulting feature
map

$0 \le i \le \left(\dfrac{w-k}{s}\right) + 1$

$0 \le o \le \left(\dfrac{h-k}{s} + 1\right)$

$[i, o]$
Zero Padding

Apply non-linearity to get the final output

$$O[i, \delta] = g(z[i, \delta]) \qquad \# i, \delta$$

$\quad \hookrightarrow$ sigmoid (activation fn.)

$\hookrightarrow$ Some intuition about how to backpropagate
gradients in this network :-

$\rightarrow$ think about vt

(A) $\quad \dfrac{\partial J}{\partial W[u,v,\ell]} \xrightarrow{\text{error metric}} = \left\{ \sum\limits_{i,\delta} \left[ \dfrac{\partial J}{\partial z[i,\delta]} \right] \cdot \left[ \dfrac{\partial z[i,\delta]}{\partial W[u,v,\ell]} \right] \right.$

$\hookrightarrow$ generalized chain rule

$0 \leq u, v \leq k$
$0 \leq \ell \leq d$

$$\boxed{\dfrac{\partial J}{\partial net_j}}$$

neural networks $[f(y_1 - - y_k)]$

$y_\ell (x_1 - - x_n)$

$\dfrac{\partial f}{\partial x_t} = \left[ \sum\limits_{\ell} \dfrac{\partial f}{\partial y_\ell} \dfrac{\partial y_\ell}{\partial x_t} \right]$

(B)

chain
rule $\leftarrow$ $\left[\dfrac{\partial J}{\partial Z_t[i,\delta]}\right]$ $=$ $\displaystyle\sum_{(i',\delta')} \dfrac{\partial J}{\partial Z_{t+1}[i',\delta']} \cdot \left[\dfrac{\partial Z_{t+1}[i',\delta']}{\partial Z_t[i,\delta]}\right]$

$\underbrace{(i',\delta')}_{Z_{t+1}[i',\delta']}$ layer $\overline{t+1}$

index $\leftarrow Z_t[i,\delta]$

$\dfrac{\partial \text{net}_\ell}{\partial \text{net}_j}$

$\ell \in \text{downNbr}(j)$

At layer t

(C) $\quad \dfrac{\partial Z_{t+1}[i',\delta']}{\partial Z_t[i,\delta]} \qquad \forall \ i,\delta, i',\delta' \ \Big]$

(D) $\dfrac{\partial O_t[i,\delta]}{\partial Z_t[i,\delta]} = \dfrac{\partial g(Z_t[i,\delta])}{\partial Z_t[i,\delta]} = O_t[i,\delta](1-O_t[i,\delta])$

Sigmoid

In summary:-    $Z_t[i,\delta]$                        $Z_{t+1}[i',\delta']$

(A) $\dfrac{\partial J}{\partial W_t[u,v,\ell]}$

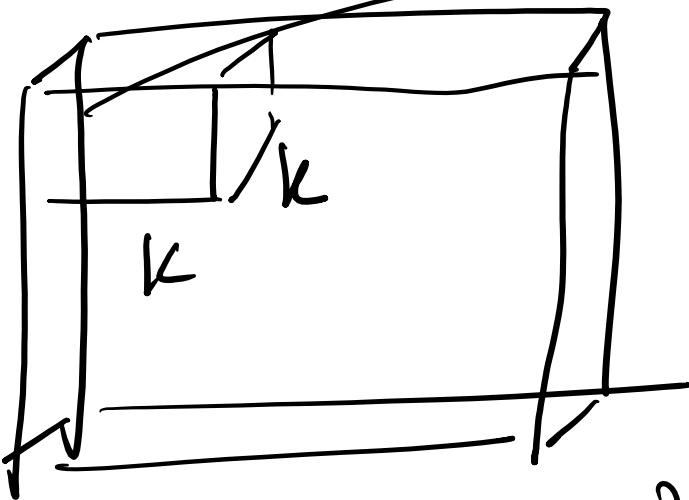[expressed in terms of

$\dfrac{\partial J}{\partial Z_t[i,\delta]}$

(B) $\dfrac{\partial J}{\partial Z_t[i,\delta]}$ expressed in terms of

$\partial J / \partial Z_{t+1}[i',\delta']$

(C) $\dfrac{\partial Z_{t+1}[i',\delta']}{\partial Z_t[i,\delta]}$ (reqi. in B)

(D) $\dfrac{\partial O_t[i,\delta]}{\partial Z_t[i,\delta]}$

"depth"

$x[i,j,\ell]$ (linear operation)

$Z[i,j]$

$0 \leq i < w$
$0 \leq j \leq h$
$0 \leq \ell < d$

$W[u,v,\ell]$

$k^2 d$ sized
kernel
(filter)

$\Rightarrow$ Applying multiple filters
(kernels)

#of kernels (feature map)

$g(Z[i,j])$
↳ activation
function

$W_t^s[u,v,\ell]$
$s \in \{1 --- d_{\ell+1}\}$  32,64, 128

$$\frac{\partial J}{\partial Z_t[i,j,l]} \Bigg]$$ uniform treatment

Arg [sum]
Max

$t$

Pooling Operation :—          depth 1

OP→$(Z_t[i,j])$

$g(Z_t[i,j])$          $:Z_t[i+1,j+1]$

Detering a
part

$b=2$

$Z_t[i,j]$

$0 \leq i' < w'$

$0 \leq j' < h'$

$w' = w/b, \; h' = h/b$

Convolution $g(z[i,j])$ $g(z[i+1,j])$

$g(z[i,j+1])$ $g(z[i+1,j+1])$

Sum/
Avg/max

$OP(g(z[i,j]),$
$g(z[i+1,j]),$
$g(z[i,j+1]),$
$g(z[i+1,j+1])$

→ Robust → less overfitting } → handle noise

→ scalability — $[w, h, d] → [w/b, h/b, d]$

Reduction in size by $b^2$

[AlexNet:- (2012) ] convolutional Neural Networks
Imagenet ] $\approx$ Lnullion (s) image
↳ Layers of [convolution/pooling] followed fully connected.