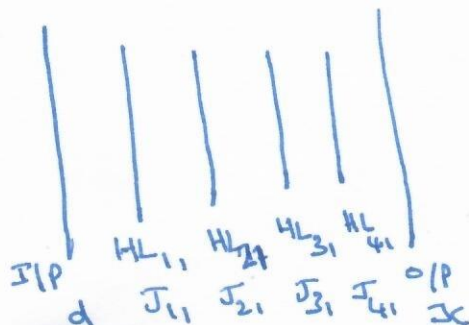


1. (a) Adam is expected to give the best convergence mainly because it combines the momentum based method and RMSProp. Numerator corresponds to the exponentially weighted average of the first moment of gradient. Denominator term corresponds to the exponentially weighted average of the second moment of gradient.

(b) Internal covariant shift: change in the distribution of output of a hidden node due to different mini-batches. Mean and variance used in normalization are computed using the output of a node for examples in a mini-batch. The scale and shift parameters are learnt during training.

2. (b) The AANN₁ is ~~used~~ trained using $\mathcal{D}_1 = \{x_n\}_{n=1}^N$. After AANN₁ is trained, every \bar{x}_n is passed through it to get \bar{z}_{1n} . The AANN₂ is trained using $\mathcal{D}_2 = \{\bar{z}_{1n}\}_{n=1}^N$. After AANN₂ is trained, every \bar{z}_{1n} is passed through it to get \bar{z}_{2n} . This process is continued to train AANN₃ and AANN₄.

(c)



HL_{i+1} is the first hidden layer of AANN_i, with $J_{i,1}$ nodes.
 $i = 1, 2, 3, 4$

3.

Gray level image.

Input: $W \times H \times 1$ CL : ~~conv~~ $F \times F$ kernel Tanh
3 feature maps $W_1 \times H_1 \times 3$

$$W_1 = \left\lfloor \frac{W - F + 2P}{S} \right\rfloor + 1$$

$$W_1 = W - F + 1$$

$$H_1 = H - F + 1$$

$$PL : W_2 = \frac{W_1}{2}$$

$$H_2 = \frac{H_1}{2}$$

$$W_2 \times H_2 \times 3$$

FC : L nodes~~Logistic~~ TanhO/P Layer : K nodes

Logistic

(a) No. of connections:

CL: No. of neurons in a feature map: $W_1 \times H_1$
No. of connections to a neuron: $F \times F$ Total no. of connections: $W_1 \times H_1 \times F \times F \times 3$ FC: No. of inputs to a neuron in FC layer: $W_2 \times H_2 \times 3$
No. of connections to FC layer: $W_2 \times H_2 \times 3 \times L$ O/P layer: No. of connections: $K \times L$ Total no. of connections with weights: $W_1 \times H_1 \times F \times F \times 3 + W_2 \times H_2 \times 3 \times L + K \times L$ No. of weight parameters in CL to be learnt: $F \times F \times 3$ Total no. of weight parameters: $F \times F \times 3 + W_2 \times H_2 \times 3 \times L + K \times L$

(b)

$$E_{tot} = \frac{1}{2} \sum_{k=1}^K (t_k - y_k)^2$$

$$\Delta w_{jk} = -\eta \frac{\partial E}{\partial w_{jk}}$$

$$= \eta (t_k - y_k) \cdot \frac{d \phi_0(a_k)}{da_k} \cdot \delta_L^{FC}$$

$$\Delta w_{jk} = \eta (t_k - y_k) \cdot \beta_0 \cdot \phi_0(a_k) (1 - \phi_0(a_k)) \cdot \delta_L^{FC}$$

$$\delta_k = (t_k - y_k) \cdot \beta_0 \cdot \phi_0(a_k) (1 - \phi_0(a_k))$$

(c)

$$\Delta w_{mijl} = \eta \delta_L \Delta m_{ij} \cdot r_{mij}$$

$$\delta_L = \left(\sum_{k=1}^K w_{Lk} \delta_k \right) \cdot \frac{d \phi_{FC}(a_L)}{da_L}$$

$$\delta_L = \left(\sum_{k=1}^K w_{Lk} \delta_k \right) \cdot \beta_{FC} (1 - \phi_{FC}^2(a_L))$$

Δm_{ij} : output of (i,j) th neuron in the m th FM

r_{mij} : 1, if Δm_{ij} is maximum of all the inputs to the neuron in PL to which the (i,j) th neuron is connected
 $= 0$, otherwise.

(d)

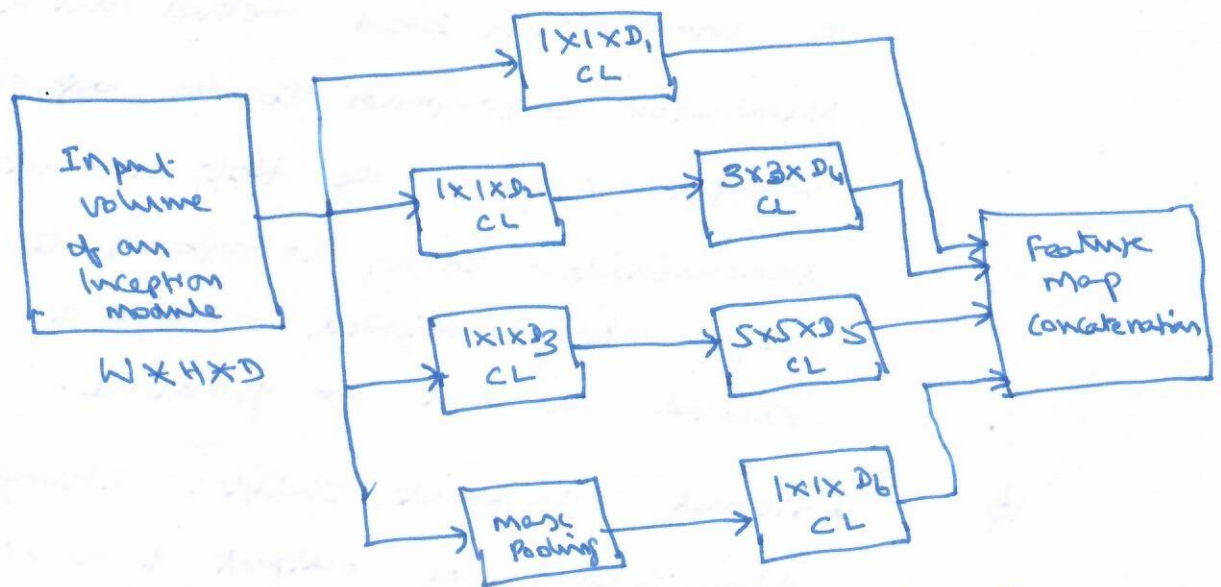
$$\Delta w_{pvi}^m = \frac{1}{W_p \times H_1} \sum_{i=1}^{W_1} \sum_{j=1}^{H_1} \Delta w_{pvij}^m$$

$$\Delta w_{pvij}^m = \eta \delta_{ij}^m x_{pvi}$$

$$\delta_{ij}^m = \left(\sum_{L=1}^L \delta_L \cdot w_{mijL} \cdot r_{mij} \right) \cdot \frac{d \phi_{CL}(a_{ij}^m)}{da_{ij}^m}$$

$$= \left(\sum_{L=1}^L \delta_L \cdot w_{mijL} \cdot r_{mij} \right) \cdot \beta_{CL} (1 - \phi_{CL}^2(a_{ij}^m))$$

4. Typical structure of inception module:



D_1, D_2, D_3, D_6 are smaller than D .

Main purpose:

Extract features with different resolution, using kernels of different sizes.

$1 \times 1 \times D_i$ CLs are used to reduce the depth of volume on which kernels of different sizes are applied.