

CNN - MNIST Dataset

Overview of General CNN's

a) Forward Pass: Convolution

- Set of learnable filters (kernels) $\sim 3 \times 3$ or 5×5
- Slide convolutional operations
- Single input + stride of 1 (no padding) \Rightarrow 2D convolutional slide @ position (i, j) :

$$Y[i, j] = \sum_{p=0}^{F-1} \sum_{q=0}^{F-1} (X[i+p, j+q] \cdot W[p, q]) + b$$

X : input (for MNIST, it's a 28×28 matrix)
 W : Filter

Multiple channels (Deeper layers means colored images) \therefore

each kernel has shape $(F, F, C_{in}) \Rightarrow C_{in}$ # input channels for multiple feature maps.

$$Y_k[i, j] = \sum_{c=1}^{C_{in}} \sum_{p=0}^{F-1} \sum_{q=0}^{F-1} (X_c[i+p, j+q] \cdot W_{p,q,c}) + b_k$$

b) Backpropagation via Convolution: In training, we filter weights through PDEs (partial diff. eqs.) of loss fn. \mathcal{L} wrt for each weight:

$$\frac{\partial \mathcal{L}}{\partial W_{p,q,c,k}} = \sum_{i,j} \frac{\partial \mathcal{L}}{\partial Y_k[i,j]} \cdot \frac{\partial Y_k[i,j]}{\partial W_{p,q,c,k}}$$

bias term: $\frac{\partial \mathcal{L}}{\partial b_k} = \sum_{i,j} \frac{\partial \mathcal{L}}{\partial Y_k[i,j]}$ (Relativ Linear Unit)

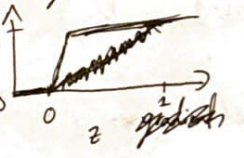
c) Activation Functions: Usually ~~ReLU~~ ReLU

$$ReLU(z) = \max(0, z)$$

Each element $Y[i, j, k]$ from convolution is passed through ReLU to obtain

$$A_k[i, j] = \max(0, Y_k[i, j])$$
$$\frac{\partial \mathcal{L}}{\partial Y_k[i, j]} = \begin{cases} \frac{\partial \mathcal{L}}{\partial A_k[i, j]} & \text{if } Y_k[i, j] \geq 0 \\ 0 & \text{if } Y_k[i, j] < 0 \end{cases}$$

ReLU(z) = piecewise linear \therefore gradient = 1 for $z \geq 0$ and 0 for $z < 0$.
(V similar to) $\left\{ \begin{matrix} ReLU(z) \\ \text{familiar x form} \end{matrix} \right.$



d) Pooling layers [Max Pooling]: w/ convolution + ReLU, CNNs use pooling layers to reduce spatial dimensions \therefore the most common form is a 2×2 pool size which is a moving box "striding" or square box which formulates a single output in the pooled map. This decreases variance when translating data.

$$\frac{\partial \mathcal{L}}{\partial X_c[i+p, j+q]} = \sum_{k} \frac{\partial \mathcal{L}}{\partial A_k[i', j']} \cdot \frac{\partial A_k[i', j']}{\partial X_c[i+p, j+q]}$$

Equation of pooling: $P_{i,j} = \max_{p,q} Y_{i+p,j+q}$

Used ReLU

~~$$A[i,j] = \max_{(p,q) \in \Omega} A[i \cdot 2 + p, j \cdot 2 + q]$$~~

$$A[i,j] = \max_{(p,q) \in \Omega} (A[i \cdot 2 + p, j \cdot 2 + q])$$

Ω = local 2x2 region in activation map A

Backpropagation through max pooling:

$$\frac{dL}{dA[i,j]} = \begin{cases} \frac{dL}{dA[i',j']} & \text{if } A[i,j] \text{ is max in } 2 \times 2 \text{ region} \\ 0 & \text{if else} \end{cases}$$

d) Fully Connected (Dense) Layers

$$z_m = \sum_{d=1}^D w_{m,d} h_d + b_m, \text{ for } m=1, \dots, M$$

↓
With a non-linear activation (ReLU) element wise:

$$a_m = \sigma(z_m)$$

For multiple fully connected layers, the output of one becomes input to the next:

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}, \quad a^{(l)} = \sigma(z^{(l)})$$

e) Output layer + Softmax:

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{j=0}^9 e^{z_j}}, \text{ for } k=0, \dots, 9$$

[MNIST digits 0 to 9]

f) Loss Function (Cross-Entropy):

$$L(\hat{y}, y) = - \sum_{j=0}^9 y_j \cdot \ln(\hat{y}_j)$$

g) Weight Updates (Gradient Descent):

$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

η = learning rate and $\frac{dL}{d\theta} \rightarrow$ computed in ~~for~~ backpropagation