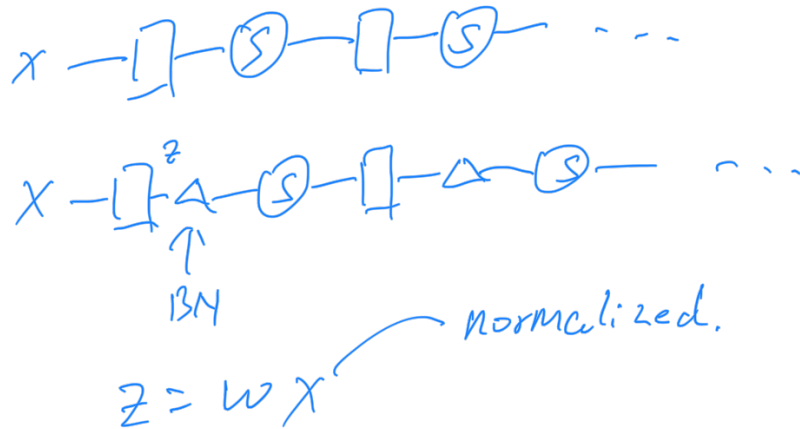


## Lecture 4

### 1. Batch normalization



$$z = ax + b$$

$$= \underbrace{\begin{bmatrix} a & b \end{bmatrix}}_w \underbrace{\begin{bmatrix} x \\ 1 \end{bmatrix}}_x$$

Batch size  $m$

$$z_1 = wX_1, \quad z_2 = wX_2, \quad \dots, \quad z_n = wX_n$$

$$\mu = \frac{1}{m} \sum_{i=1}^m z_i, \quad \tilde{z}_j = z_j - \mu$$

$$\sigma = \frac{1}{m} \sum_{i=1}^m \tilde{z}_i^2$$

$$\hat{\tilde{z}}_j = \frac{\tilde{z}_j}{\sqrt{\sigma + \epsilon}}$$

$$\dots \hat{\tilde{z}}_j + \gamma$$

$$H = \alpha z_j^T K$$

$\underbrace{\hspace{10em}}_{\text{learned}}$

## 2. Optimization

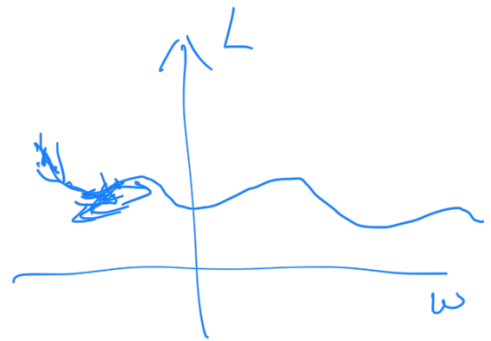
Gradient descent

$$\rightarrow w_{t+1} = w_t - \lambda \frac{\partial L}{\partial w_t}$$

$$L = \sum_{i=1}^m L_i$$

$m$  — full dataset

SAG  $m$  — batch size



Momentum

$$\text{GD: } w_{t+1} = w_t - \lambda g_t$$



momentum:  $\mu_{t+1} = \gamma \mu_t - \lambda g_t$

$$w_{t+1} = w_t + \mu_{t+1}$$

momentum example:  $\gamma = 0.9, \mu_0 = 0$

$$\mu_1 = 0.9 \cdot 0 - \lambda g_0 = -\lambda g_0$$

$$w_1 = w_0 - \lambda g_0$$

$$t=2 \quad \mu_2 = 0.9 \cdot \mu_1 - \lambda g_1$$

$$= -0.9\lambda g_0 - \lambda g_1 = -\lambda(0.9g_0 + g_1)$$

$$W_2 = W_1 - \lambda(0.9g_0 + g_1)$$

$\gamma$  change from 0.5 to 0.9  
over iterations.

Adam

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$V_t = \beta_2 V_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{V}_t = \frac{V_t}{1 - \beta_2^t}$$

$$U_t = -\lambda \frac{\hat{m}_t}{\sqrt{\hat{V}_t} + \epsilon}$$

$$W_{t+1} = W_t + U_t$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$



3. convolutional NN

~ convolution

one-D convolution

$$(X * W)(t) = \sum_{a=-\infty}^{\infty} x[a] w[t-a]$$

Two-D conv

$$(I * K)(i, j) = \sum_m \sum_n i[m, n] K[i-m, j-n]$$

$$= \sum_m \sum_n I[i-m, j-n] k[m, n]$$

cross-correlation

$$(I * K)(i, j) = \sum_m \sum_n I[i+m, j+n] K[m, n]$$

$$I = \begin{bmatrix} 1 & 2 & 3 & 2 & 1 \end{bmatrix}$$

$$K = \begin{bmatrix} 1 & 0 & 1 \end{bmatrix}$$

$$I * K = \begin{bmatrix} -2 & 0 & 2 \end{bmatrix}$$

$$I = \begin{bmatrix} 1 & 2 & 1 & 1 & 3 & 0 \\ 2 & 0 & 1 & 1 & 1 & 2 \\ 3 & 2 & 1 & 0 & 1 & \\ 1 & 1 & 2 & 1 & 1 & \end{bmatrix}$$

$$K = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$

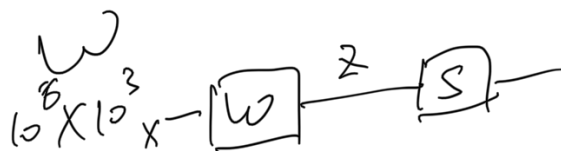
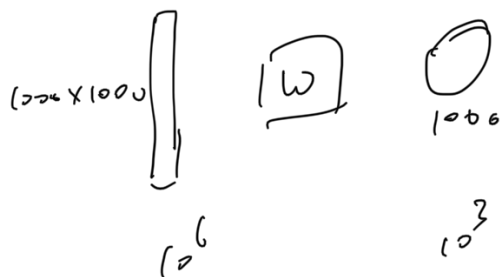
$$\sum \sum \begin{bmatrix} 1 & 2 & 1 \\ 2 & 0 & 1 \\ 3 & 2 & 1 \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = 0$$

$$\rightarrow \begin{bmatrix} 1 & 2 & 1 & 3 \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & 1 \end{bmatrix} = 4$$

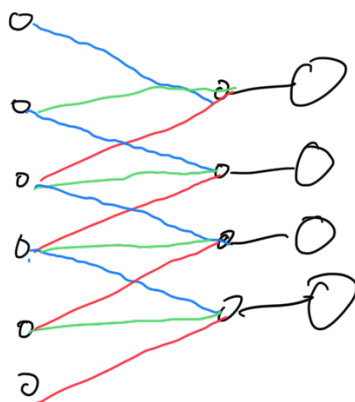
$$\Rightarrow \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}^{-1}$$

$$\Rightarrow \begin{bmatrix} 1 & 3 & 0 \\ 1 & 1 & 2 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix} = 2$$

$$= \begin{bmatrix} 0 & 4 & 2 \\ 6 & & \end{bmatrix}$$

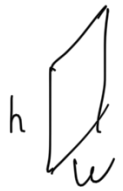


$w_1$   
 $w_2$   
 $w_3$

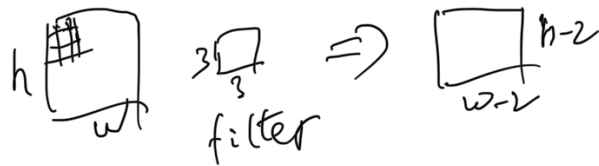



$$(X_{i+1} \dots X_{i+m}) \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix} = z_i$$

$$\frac{\partial z_i}{\partial w_j} = X_{i+j}$$


 Term  
 receptive field aka.  
 filter size  
 3x3, 5x5, 7x7,

padding size



  $3 \times 3$  padding #  $\approx 1$  for each side.

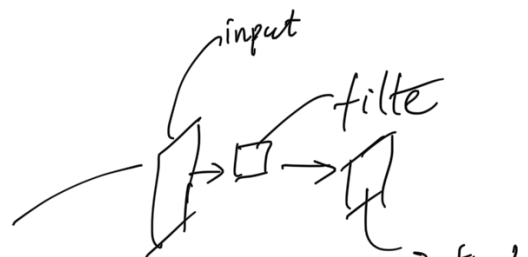
stride — step size,  $s$

filter size —  $F$

padding size —  $P$

input image  $w \times h$

size of output after conv



output

output size  $m \times n$

$$m = \frac{(W - F + 2P)}{S} + 1$$

$$n = \frac{(h - F + 2P)}{S} + 1$$