

Matrix Operations I

- For efficient implementations, we should conduct convolutional operations by **matrix-matrix** and **matrix-vector** operations
- Let's collect images of all channels as the input

$$\begin{aligned} & Z^{\text{in},i} \\ = & \begin{bmatrix} z_{1,1,1}^i & z_{2,1,1}^i & \cdots & z_{a^{\text{in}},b^{\text{in}},1}^i \\ \vdots & \vdots & \ddots & \vdots \\ z_{1,1,d^{\text{in}}}^i & z_{2,1,d^{\text{in}}}^i & \cdots & z_{a^{\text{in}},b^{\text{in}},d^{\text{in}}}^i \end{bmatrix} \\ & \in \mathbb{R}^{d^{\text{in}} \times a^{\text{in}} b^{\text{in}}}. \end{aligned}$$

Matrix Operations II

- Let all filters

$$W = \begin{bmatrix} w_{1,1,1}^1 & w_{2,1,1}^1 & \cdots & w_{h,h,d^{in}}^1 \\ \vdots & \vdots & \ddots & \vdots \\ w_{1,1,1}^{d^{out}} & w_{2,1,1}^{d^{out}} & \cdots & w_{h,h,d^{in}}^{d^{out}} \end{bmatrix}$$
$$\in \mathbb{R}^{d^{out} \times hhd^{in}}$$

be variables (parameters) of the current layer

Matrix Operations III

- Usually a bias term is considered

$$\mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_{d^{\text{out}}} \end{bmatrix} \in R^{d^{\text{out}} \times 1}$$

- Operations at a layer

$$\begin{aligned} S^{\text{out},i} &= W\phi(Z^{\text{in},i}) + \mathbf{b}\mathbf{1}_{a^{\text{out}}b^{\text{out}}}^T \\ &\in R^{d^{\text{out}} \times a^{\text{out}}b^{\text{out}}}, \end{aligned} \tag{1}$$

Matrix Operations IV

where

$$\mathbb{1}_{a^{\text{out}} b^{\text{out}}} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \in R^{a^{\text{out}} b^{\text{out}} \times 1}.$$

- $\phi(Z^{\text{in},i})$ collects all sub-images in $Z^{\text{in},i}$ into a matrix.

Matrix Operations V

Specifically,

$$\phi(Z^{\text{in},i}) = \begin{bmatrix} z_{1,1,1}^i & z_{1+s,1,1}^i & & z_{1+(a^{\text{out}}-1)s,1+(b^{\text{out}}-1)s,1}^i \\ z_{2,1,1}^i & z_{2+s,1,1}^i & & z_{2+(a^{\text{out}}-1)s,1+(b^{\text{out}}-1)s,1}^i \\ \vdots & \vdots & \dots & \vdots \\ z_{h,h,1}^i & z_{h+s,h,1}^i & & z_{h+(a^{\text{out}}-1)s,h+(b^{\text{out}}-1)s,1}^i \\ \vdots & \vdots & & \vdots \\ z_{h,h,d^{\text{in}}}^i & z_{h+s,h,d^{\text{in}}}^i & & z_{h+(a^{\text{out}}-1)s,h+(b^{\text{out}}-1)s,d^{\text{in}}}^i \end{bmatrix}$$
$$\in \mathbb{R}^{hhd^{\text{in}} \times a^{\text{out}} b^{\text{out}}}$$

Activation Function I

- Next, an activation function scales each element of $S^{\text{out},i}$ to obtain the output matrix $Z^{\text{out},i}$.

$$Z^{\text{out},i} = \sigma(S^{\text{out},i}) \in R^{d^{\text{out}} \times a^{\text{out}} b^{\text{out}}}. \quad (2)$$

- For CNN, commonly the following RELU activation function

$$\sigma(x) = \max(x, 0) \quad (3)$$

is used

- Later we need that $\sigma(x)$ is differentiable, but the RELU function is not.

Activation Function II

- Past works such as Krizhevsky et al. (2012) assume

$$\sigma'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

The Function $\phi(Z^{\text{in},i})$ I

- In the matrix-matrix product

$$W\phi(Z^{\text{in},i}),$$

each element is the **inner product between a filter and a sub-image**

- Clearly ϕ is a **linear mapping**, so there exists a **0/1** matrix P_ϕ such that

$$\phi(Z^{\text{in},i}) \equiv \text{mat} (P_\phi \text{vec}(Z^{\text{in},i}))_{hhd^{\text{in}} \times a^{\text{out}} b^{\text{out}}}, \quad \forall i, \quad (4)$$

The Function $\phi(Z^{\text{in},i})$ II

- $\text{vec}(M)$: all M 's columns concatenated to a vector \mathbf{v}

$$\text{vec}(M) = \begin{bmatrix} M_{:,1} \\ \vdots \\ M_{:,b} \end{bmatrix} \in R^{ab \times 1}, \text{ where } M \in R^{a \times b}$$

- $\text{mat}(\mathbf{v})$ is the inverse of $\text{vec}(M)$

$$\text{mat}(\mathbf{v})_{a \times b} = \begin{bmatrix} v_1 & & v_{(b-1)a+1} \\ \vdots & \dots & \vdots \\ v_a & & v_{ba} \end{bmatrix} \in R^{a \times b}, \quad (5)$$

The Function $\phi(Z^{\text{in},i})$ III

where

$$\mathbf{v} \in R^{ab \times 1}.$$

- P_ϕ is a huge matrix:

$$P_\phi \in R^{hhd^{\text{in}} a^{\text{out}} b^{\text{out}} \times d^{\text{in}} a^{\text{in}} b^{\text{in}}}$$

and

$$\phi : R^{d^{\text{in}} \times a^{\text{in}} b^{\text{in}}} \rightarrow R^{hhd^{\text{in}} \times a^{\text{out}} b^{\text{out}}}$$

- The representation of using P_ϕ makes some subsequent derivations easier

The Function $\phi(Z^{\text{in},i})$ IV

- Past works using the form (4) include, for example, Vedaldi and Lenc (2015)

Optimization Problem I

- We collect all weights to a vector variable θ .

$$\theta = \begin{bmatrix} \text{vec}(W^1) \\ \mathbf{b}^1 \\ \vdots \\ \text{vec}(W^L) \\ \mathbf{b}^L \end{bmatrix} \in R^n, \quad n : \text{total \# variables}$$

- The output of the last layer L is a vector $\mathbf{z}^{L+1,i}(\theta)$.
- Consider any loss function such as the squared loss

$$\xi_i(\theta) = \|\mathbf{z}^{L+1,i}(\theta) - \mathbf{y}^i\|^2.$$

Optimization Problem II

- The optimization problem is

$$\min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}),$$

where

$$f(\boldsymbol{\theta}) = \frac{1}{2C} \boldsymbol{\theta}^T \boldsymbol{\theta} + \frac{1}{l} \sum_{i=1}^l \xi(\mathbf{z}^{L+1,i}(\boldsymbol{\theta}); \mathbf{y}^i, Z^{1,i})$$

C : regularization parameter.

- The formulation is almost the same as that for fully connected networks

Optimization Problem III

- Note that we divide the sum of training losses by the number of training data

Thus the second term becomes the average training loss

References I

- A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. 2012.
- A. Vedaldi and K. Lenc. MatConvNet: Convolutional neural networks for matlab. In *Proceedings of the 23rd ACM International Conference on Multimedia*, pages 689–692, 2015.