

Supervised Learning VII – Neural Networks (cont.)

1. Direction of Gradient

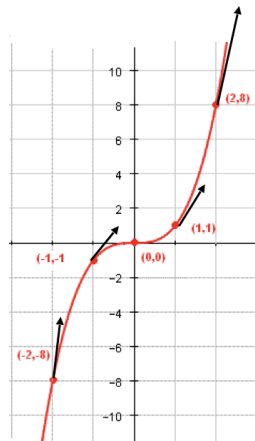
- a. Lets say we have a (differentiable) function ‘f’ we want to minimize
- b. The derivative of ‘f’ is a direction
 - i. Points to local maxima (including saddle points)
- c. We can use this to “descend!”
 - i. Move in opposite direction of derivative!
- d. Basic algorithm (in 1-D):

while not converged and $t < t_{max}$ do:

compute $\frac{df}{dx} \big|_{x^{(t)}}$

update $x^{(t+1)} = x^{(t)} - \eta \frac{df}{dx} \big|_{x^{(t)}}$

i. $t = t + 1$



e.

2. Gradient Descent in Parameter Space

- a. Remember, we want to find θ that optimizes a loss function
- b. So, we need to compute

$$\nabla_{\theta} L(\mathbf{D}, \theta) = \left(\frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \dots, \frac{\partial L}{\partial \theta_n} \right)^T$$

- c. Basic algorithm (in parameter space)

while not converged and $t < t_{max}$ do:

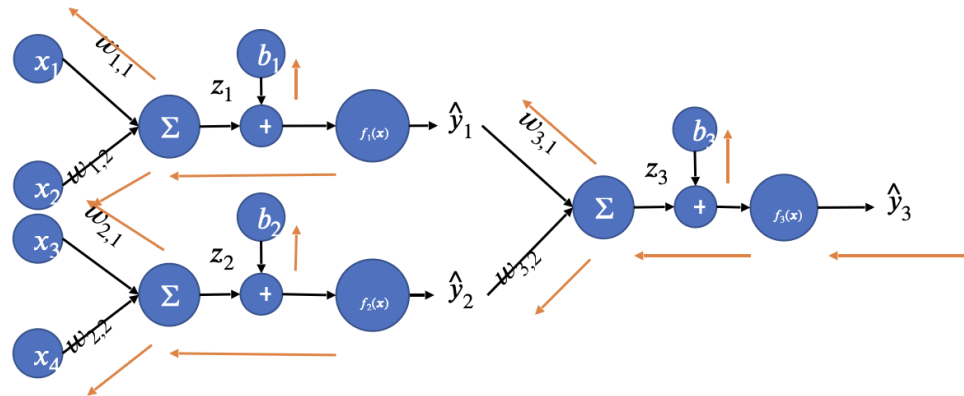
compute $\nabla_{\theta} L(\mathbf{D}, \theta) \big|_{\theta^{(t)}}$

update $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} L(\mathbf{D}, \theta) \big|_{\theta^{(t)}}$

$t = t + 1$

i.

3. Computing the Gradient



a.

b. Gradient “flows” backward (from output to input)

c. Lets do a small example, see if we notice any patterns

Look at all the redundant calculations!!!

$$\nabla_{\theta} L = \begin{pmatrix} \frac{\partial L}{\partial b_3} \\ \frac{\partial L}{\partial w_{3,1}} \\ \frac{\partial L}{\partial w_{3,2}} \\ \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_{1,1}} \\ \frac{\partial L}{\partial w_{1,2}} \\ \dots \end{pmatrix} = \begin{pmatrix} \frac{\partial L}{\partial \hat{y}_3} & \frac{\partial \hat{y}_3}{\partial z_3} & \frac{\partial z_3}{\partial w_{3,1}} & \frac{\partial z_3}{\partial w_{3,2}} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial z_1}{\partial w_{1,1}} & \frac{\partial z_1}{\partial w_{1,2}} & \dots \\ \frac{\partial L}{\partial \hat{y}_3} & \frac{\partial \hat{y}_3}{\partial z_3} & \frac{\partial z_3}{\partial w_{3,1}} & \frac{\partial z_3}{\partial w_{3,2}} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial z_1}{\partial w_{1,1}} & \frac{\partial z_1}{\partial w_{1,2}} & \dots \\ \frac{\partial L}{\partial \hat{y}_3} & \frac{\partial \hat{y}_3}{\partial z_3} & \frac{\partial z_3}{\partial w_{3,1}} & \frac{\partial z_3}{\partial w_{3,2}} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial z_1}{\partial w_{1,1}} & \frac{\partial z_1}{\partial w_{1,2}} & \dots \\ \frac{\partial L}{\partial \hat{y}_3} & \frac{\partial \hat{y}_3}{\partial z_3} & \frac{\partial z_3}{\partial w_{3,1}} & \frac{\partial z_3}{\partial w_{3,2}} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial z_1}{\partial w_{1,1}} & \frac{\partial z_1}{\partial w_{1,2}} & \dots \\ \frac{\partial L}{\partial \hat{y}_3} & \frac{\partial \hat{y}_3}{\partial z_3} & \frac{\partial z_3}{\partial w_{3,1}} & \frac{\partial z_3}{\partial w_{3,2}} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial z_1}{\partial w_{1,1}} & \frac{\partial z_1}{\partial w_{1,2}} & \dots \\ \frac{\partial L}{\partial \hat{y}_3} & \frac{\partial \hat{y}_3}{\partial z_3} & \frac{\partial z_3}{\partial w_{3,1}} & \frac{\partial z_3}{\partial w_{3,2}} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial \hat{y}_1}{\partial z_1} & \frac{\partial z_1}{\partial w_{1,1}} & \frac{\partial z_1}{\partial w_{1,2}} & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

i.

ii. Now, we cache the backpropagation values and use it again for later (to remove redundancy)

4. Backpropagation

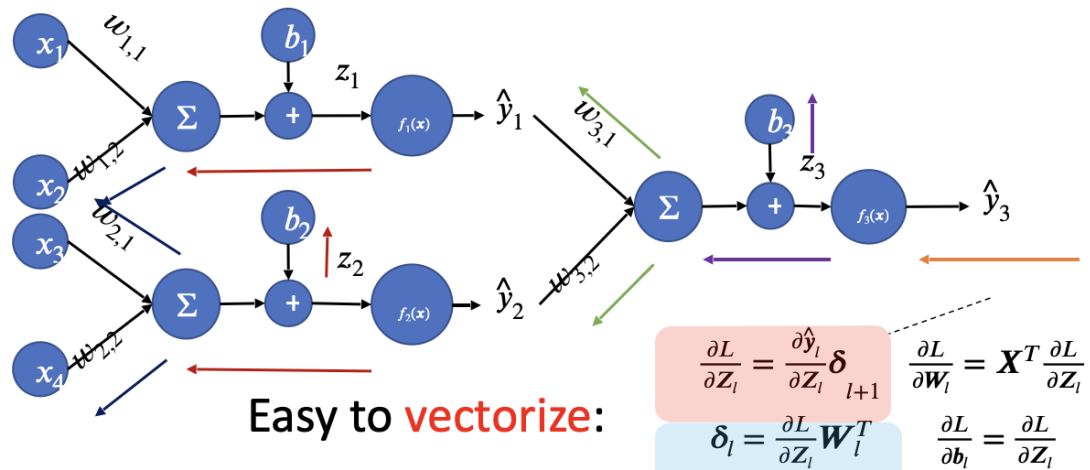
- NN computation is always a DAG
- Cache error as a function of layer
 - Remove redundant computation (from previous slide)

$$\begin{aligned}\delta_{\hat{y}_3} &= \frac{\partial L}{\partial \hat{y}_3} \\ \delta z_3 &= \delta_{\hat{y}_3} \frac{\partial \hat{y}_3}{\partial z_3} \\ \delta_{\hat{y}_1} &= \delta z_3 \frac{\partial z_3}{\partial \hat{y}_1} \\ \delta z_1 &= \delta_{\hat{y}_1} \frac{\partial \hat{y}_1}{\partial z_1}\end{aligned}$$

$$\nabla_{\theta} L = \begin{array}{ccccccc} \frac{\partial L}{\partial b_3} & \delta_{\hat{y}_3} & \delta z_3 & & & & \\ \frac{\partial L}{\partial w_{3,1}} & \delta_{\hat{y}_3} & \delta z_3 & \frac{\partial z_3}{\partial w_{3,1}} & & & \\ \frac{\partial L}{\partial w_{3,2}} & \delta_{\hat{y}_3} & \delta z_3 & \frac{\partial z_3}{\partial w_{3,2}} & & & \\ \frac{\partial L}{\partial b_1} & \delta_{\hat{y}_3} & \delta z_3 & & \delta_{\hat{y}_1} & \delta z_1 & \\ \frac{\partial L}{\partial w_{1,1}} & \delta_{\hat{y}_3} & \delta z_3 & & \delta_{\hat{y}_1} & \delta z_1 & \frac{\partial z_1}{\partial w_{1,1}} \\ \frac{\partial L}{\partial w_{1,2}} & \delta_{\hat{y}_3} & \delta z_3 & & \delta_{\hat{y}_1} & \delta z_1 & \frac{\partial z_1}{\partial w_{1,2}} \\ \dots & \delta_{\hat{y}_3} & \dots & & \delta_{\hat{y}_1} & \dots & \dots \\ & & & & \dots & \dots & \dots \end{array}$$

Error from activation func

c.



d.

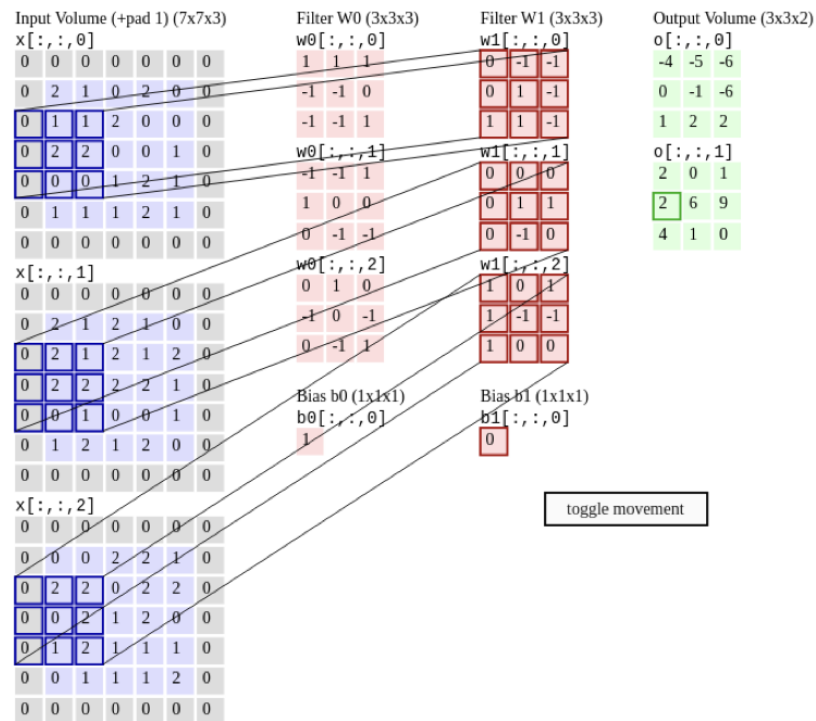
Error flowing out of layer

- Calculate once and reuse it (cache the matrix that has the values)
- Data inputs have to be vectors

5. Images and Your Eyes

- How do your eyes parse visual data?
- How do we localize different patterns in images?
 - Answer: Convolution!

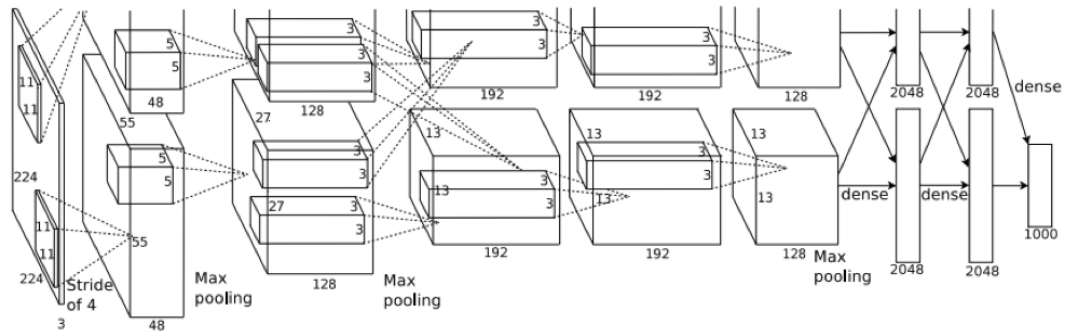
No statistical difference between eye response and convolution with gabor (wavelet) filters¹



- High level: convolution “rubs” one function (called the filter or kernel) over a signal (image)
 - Response signal is proportional to how similar geometry in signal is to filter
 - High response \rightarrow local geometry of signal is similar to geometry of filter!
- Can we learn the filters?

6. Convolutional NNs

- Implement convolution operator (over images) using learnable kernels
- Generate kernels (like gabor wavelets!)
- Successive layers \rightarrow more abstract features



d.

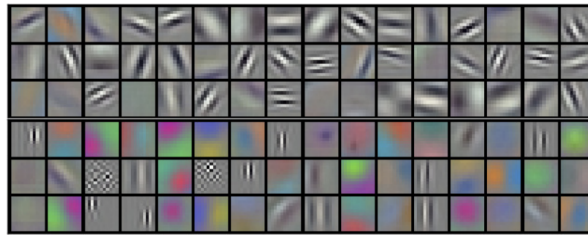


Figure 3: 96 convolutional kernels of size $11 \times 11 \times 3$ learned by the first convolutional layer on the $224 \times 224 \times 3$ input images. The top 48 kernels were learned on GPU 1 while the bottom 48 kernels were learned on GPU 2. See Section 6.1 for details.

e.