

# CSE 493G1/599G1: Deep Learning

---

## Section 4: Backprop II & CNNs

Friday, January 26, 2024.

Welcome to section, we're glad you could make it!

### 1. Sigmoid Shenanigans

Consider the Sigmoid activation function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Draw a computational graph and work through the backpropagation. Then, fill in the Python function. If you finish early, work through the analytical derivation for Sigmoid.

As a hint, you could split Sigmoid into the following functions:

$$a(x) = -x$$

$$b(x) = e^x$$

$$c(x) = 1 + x$$

$$d(x) = \frac{1}{x}$$

Observe that chaining these operations gives us Sigmoid:  $d(c(b(a(x)))) = \sigma(x)$ .

Suppose  $x = 2$ . What would the gradient with respect to  $x$  be? Feel free to use a calculator on this part.

You should have gotten around 0.1. If the step size is 0.2, what would the value of  $x$  be after taking one gradient descent step? As a hint, remember that `parameters -= step_size * gradient`.

```
1  import numpy as np
2
3  # inputs:
4  # - a numpy array `x`
5  # outputs:
6  # - `out`: the result of the forward pass
7  # - `fx` : the result of the backwards pass
8  def sigmoid(x):
9      # provided: forward pass with cache
10     a = -x
11     b = np.exp(a)
12     c = 1 + b
13     f = 1/c
14     out = f
15
16     # TODO: backwards pass, "fx" represents df / dx
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40     return out, fx
```

Ignore the line numbers, they do NOT correspond to the number of lines you need to write.

## 2. Lost In The Sauce

Consider the following function:

$$f = \frac{\ln x \cdot \sigma(\sqrt{y})}{\sigma((x+y)^2)}$$

Break the function up into smaller parts, then draw a computational graph and finish the Python function.

For reference, the derivative of Sigmoid is  $\sigma(x) \cdot (1 - \sigma(x))$ .

The TA solution breaks  $f$  into equations  $a$  through  $i$ . Yours doesn't have to match this exactly.

*Python function printed on the following page.*

```

1  import numpy as np
2
3  # helper function
4  def sigmoid(x):
5      return 1/(1 +np.exp(-x))
6
7  # inputs: numpy arrays `x`, `y`
8  # outputs: forward pass in `out`, gradient for x in `fx`, gradient for y in `fy`
9  def complex_layer(x, y):
10
11      # forward pass
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29      # backwards pass
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58  return out, fx, fy

```

Ignore the line numbers, they do NOT correspond to the number of lines you need to write.

### 3. As Convoluting As Possible

Suppose that both your input and filter are squares. Let  $W$  be the width of your input,  $F$  be the width of your filter,  $P$  be the amount of zero-padding utilized, and  $S$  be the stride.

Recall that convolution layers utilize a dot product. For instance, when you apply a  $5 \times 5 \times 3$  conv filter, it reduces the corresponding 75 input values into 1 output value. One way to phrase this, borrowed from the cs231n notes, would be “every neuron in the conv layer has 180 connections to the input volume”.

What’s the formula for determining a conv layer’s output size, ignoring padding?

What’s the formula for determining a conv layer’s output size, taking padding into account?

Consider a layer that takes a  $32 \times 32 \times 3$  input and applies a  $5 \times 5 \times 3$  filter (with stride 1 and no padding). What would the output size be?

What would happen if we changed our stride to 2?

What about a stride of 3?

What about a stride of 3 with a padding of 3 (i.e., three zero-values go on *all* sides of our input)?

People will often leave out the channels dimension when writing out a filter (i.e., they might refer to a  $5 \times 5 \times 3$  filter as a  $5 \times 5$  filter). We will now adopt this shorthand as well.

Consider the first conv layer of AlexNet, which took an input of size  $227 \times 227 \times 3$ . The layer applied 96 separate  $11 \times 11$  filters with stride of 4 and no padding. You might sometimes see people write this as applying one  $11 \times 11$  filter with depth 96; they mean the same thing. What would our output dimensions be?

We can use  $K$  to refer to the number of filters we apply at a conv layer.

Also, let  $C$  be the number of channels in our input. For example,  $C = 3$  for a RGB image.

For any given conv layer, there will be  $KF^2C$  parameters for the “weights” of our filters, and  $K$  parameters for our biases. More simply, there are  $K(F^2C + 1)$  trainable parameters.

Consider the first layer of AlexNet mentioned above. How many trainable parameters are there?

Consider a conv layer which takes a  $31 \times 31 \times 5$  input and applies 25 filters of size  $3 \times 3 \times 5$  with stride 2 and padding 1. How many trainable parameters does it have?

A key takeaway here is that the values you choose to set  $K$  and  $F$  to (these are hyperparameters!) will have a significant impact on the number of parameters your model has. You must be careful not to add too many parameters to your model, especially since we have limited compute in class.

Recall your answer for the output of AlexNet’s first layer from above. The second layer applies a  $3 \times 3$  pool filter at stride 2, again with no padding. What will the output size be? How many trainable parameters will this layer introduce?

Did the pool layer change the number of channels? Does this pattern generalize to pool layers of all sizes?