

Homework 7: Computation Graphs, Back-propagation, and Neural Networks

Due: Thursday, May 6th, 2021 at 11:59PM EST

Instructions: Your answers to the questions below, including plots and mathematical work, should be submitted as a single PDF file. It's preferred that you write your answers using software that typesets mathematics (e.g. LaTeX, LyX, or MathJax via iPython), though if you need to you may scan handwritten work. You may find the minted package convenient for including source code in your LaTeX document. If you are using LyX, then the listings package tends to work better.

1 Introduction

There is no doubt that neural networks are a very important class of machine learning models. Given the sheer number of people who are achieving impressive results with neural networks, one might think that it's relatively easy to get them working. This is a partly an illusion. One reason so many people have success is that, thanks to GitHub, they can copy the exact settings that others have used to achieve success. In fact, in most cases they can start with “pre-trained” models that already work for a similar problem, and “fine-tune” them for their own purposes. It's far easier to tweak and improve a working system than to get one working from scratch. If you create a new model, you're kind of on your own to figure out how to get it working: there's not much theory to guide you and the rules of thumb do not always work. Understanding even the most basic questions, such as the preferred variant of SGD to use for optimization, is still a very active area of research.

One thing is clear, however: If you do need to start from scratch, or debug a neural network model that doesn't seem to be learning, it can be immensely helpful to understand the low-level details of how your neural network works – specifically, back-propagation. With this assignment, you'll have the opportunity to linger on these low-level implementation details. Every major neural network type (RNNs, CNNs, Resnets, etc.) can be implemented using the basic framework we'll develop in this assignment.

To help things along, Philipp Meerkamp, Pierre Garapon, and David Rosenberg have designed a minimalist framework for computation graphs and put together some support code. The intent is for you to read, or at least skim, every line of code provided, so that you'll know you understand all the crucial components and could, in theory, create your own from scratch. In fact, creating your own computation graph framework from scratch is highly encouraged – you'll learn a lot.

2 Computation Graph Framework

To get started, please read the tutorial on the computation graph framework we'll be working with. (Note that it renders better if you view it locally.) The use of computation graphs is not specific to machine learning or neural networks. Computation graphs are just a way to represent a function that facilitates efficient computation of the function's values and its gradients with respect to inputs. The tutorial takes this perspective, and there is very little in it about machine learning, per se.

To see how the framework can be used for machine learning tasks, we've provided a full implementation of linear regression. You should start by working your way through the `__init__` of the `LinearRegression` class in `linear_regression.py`. From there, you'll want to review the node class definitions in `nodes.py`, and finally the class `ComputationGraphFunction` in `graph.py`. `ComputationGraphFunction` is where we repackage a raw computation graph into something that's more friendly to work with for machine learning. The rest of `linear_regression.py` is fairly routine, but it illustrates how to interact with the `ComputationGraphFunction`.

As we've noted earlier in the course, getting gradient calculations correct can be difficult. To help things along, we've provided two functions that can be used to test the backward method of a node and the overall gradient calculation of a `ComputationGraphFunction`. The functions are in `test_utils.py`, and it's recommended that you review the tests provided for the linear regression implementation in `linear_regression.t.py`. (You can run these tests from the command line with `python3 linear_regression.t.py`.) The functions actually doing the testing, `test_node_backward` and `test_ComputationGraphFunction`, may seem a bit intricate, but they're implementing the exact same `gradient_checker` logic we saw in the second homework assignment.

Once you've understood how linear regression works in our framework, you're ready to start implementing your own algorithms. To help you get started, please make sure you are able to execute the following commands:

- `cd /path/to/hw7`
- `python3 linear_regression.py`
- `python3 linear_regression.t.py`

3 Ridge Regression

When moving to a new system, it's always good to start with something familiar. But that's not the only reason we're doing ridge regression in this homework. In ridge regression the parameter vector is "shared", in the sense that it's used twice in the objective function. In the computation graph, this can be seen in the fact that the node for the parameter vector has two outgoing edges. This sharing is common many popular neural networks (RNNs and CNNs), where it is often referred to as *parameter tying*.

`ridge_regression.py` provides a skeleton code and `ridge_regression.t.py` is a test code, which you should eventually be able to pass.

1. Complete the class `L2NormPenaltyNode` in `nodes.py`. If your code is correct, you should be able to pass `test_L2NormPenaltyNode` in `ridge_regression.t.py`. Please attach a screenshot that shows the test results for this question.

```
class L2NormPenaltyNode(object):
    """ Node computing  $l2\_reg * ||w||^2$  for scalars  $l2\_reg$  and vector  $w$  """
    def __init__(self, l2_reg, w, node_name):
        """
        Parameters:
        l2_reg: a numpy scalar array (e.g. np.array(.01)) (not a node)
        w: a node for which w.out is a numpy vector
        node_name: node's name (a string)
        """
        self.node_name = node_name
        self.out = None
        self.d_out = None
        self.l2_reg = np.array(l2_reg)
        self.w = w

    def forward(self):
        self.out = self.l2_reg * np.sum(self.w.out**2)
        self.d_out = np.zeros(self.out.shape)
        return self.out

    def backward(self):
        self.w.d_out += self.d_out * self.l2_reg * 2 * self.w.out

    def get_predecessors(self):
        return [self.w]
```

Test results for all 3 tests (as the .t.py file is not to be modified, running all 3 tests at once):

```
DEBUG: (Node l2 norm node) Max rel error for partial deriv w.r.t. w is 1.2987169065072546e-09.  
.DEBUG: (Node sum node) Max rel error for partial deriv w.r.t. a is 1.6365788917502879e-09.  
DEBUG: (Node sum node) Max rel error for partial deriv w.r.t. b is 5.838672524325393e-10.  
.DEBUG: (Parameter w) Max rel error for partial deriv 3.180347658756531e-10.  
DEBUG: (Parameter b) Max rel error for partial deriv 4.379296446951102e-10.  
.  
-----  
Ran 3 tests in 0.002s  
  
OK  
  
Process finished with exit code 0
```

2. Complete the class `SumNode` in `nodes.py`. If your code is correct, you should be able to pass `test_SumNode` in `ridge_regression.t.py`. Please attach a screenshot that shows the test results for this question.

```
class SumNode(object):
    """ Node computing  $a + b$ , for numpy arrays  $a$  and  $b$  """
    def __init__(self, a, b, node_name):
        """
        Parameters:
        a: node for which a.out is a numpy array
        b: node for which b.out is a numpy array of the same shape as a
        node_name: node's name (a string)
        """
        self.node_name = node_name
        self.out = None
        self.d_out = None
        self.b = b
        self.a = a

    def forward(self):
        self.out = self.a.out + self.b.out
        self.d_out = np.zeros(self.out.shape)
        return self.out

    def backward(self):
        d_a = self.d_out * 1
        d_b = self.d_out * 1
        self.a.d_out += d_a
        self.b.d_out += d_b

    def get_predecessors(self):
        return [self.a, self.b]
```

Test results for all 3 tests (as the .t.py file is not to be modified, running all 3 tests at once):

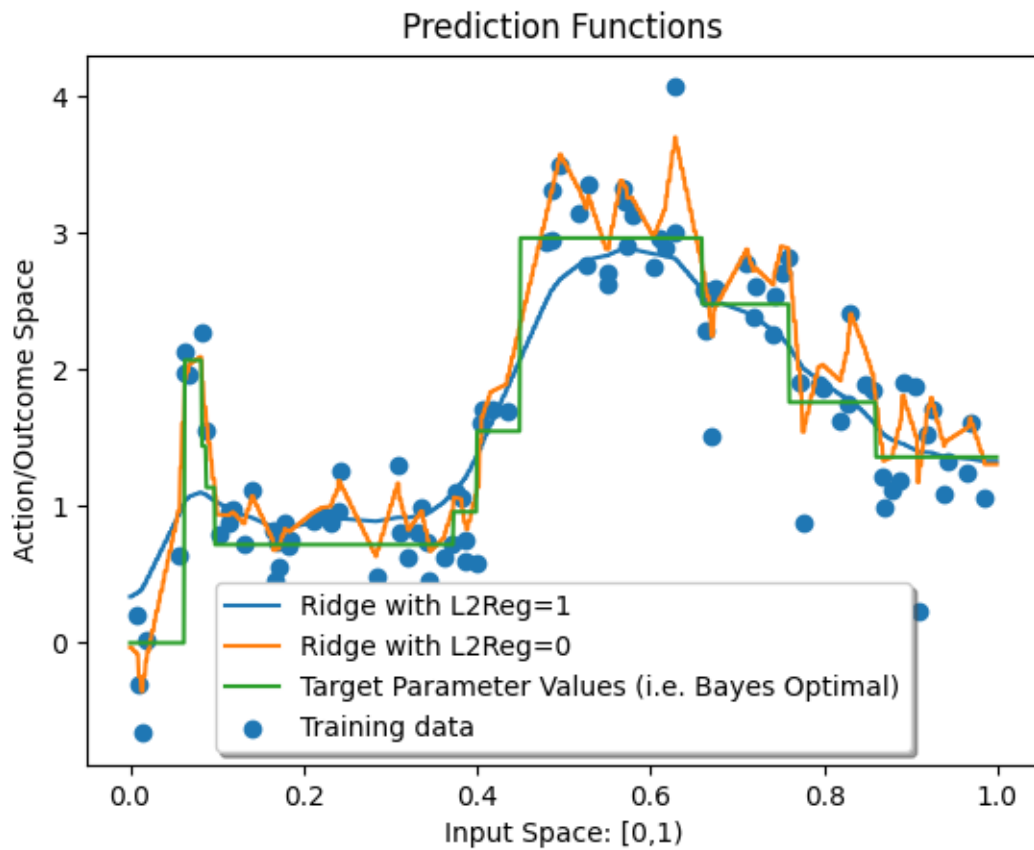
```
DEBUG: (Node l2 norm node) Max rel error for partial deriv w.r.t. w is 1.2987169065072546e-09.  
.DEBUG: (Node sum node) Max rel error for partial deriv w.r.t. a is 1.6365788917502879e-09.  
DEBUG: (Node sum node) Max rel error for partial deriv w.r.t. b is 5.838672524325393e-10.  
.DEBUG: (Parameter w) Max rel error for partial deriv 3.180347658756531e-10.  
DEBUG: (Parameter b) Max rel error for partial deriv 4.379296446951102e-10.  
.  
-----  
Ran 3 tests in 0.002s  
  
OK  
  
Process finished with exit code 0
```


Test results for all 3 tests (as the .t.py file is not to be modified, running all 3 tests at once):

```
DEBUG: (Node l2 norm node) Max rel error for partial deriv w.r.t. w is 1.2987169065072546e-09.
.DEBUG: (Node sum node) Max rel error for partial deriv w.r.t. a is 1.6365788917502879e-09.
DEBUG: (Node sum node) Max rel error for partial deriv w.r.t. b is 5.838672524325393e-10.
.DEBUG: (Parameter w) Max rel error for partial deriv 3.180347658756531e-10.
DEBUG: (Parameter b) Max rel error for partial deriv 4.379296446951102e-10.
.
-----
Ran 3 tests in 0.002s

OK

Process finished with exit code 0
```



Final average training square objective and loss for the 2 settings are: (made a minor change in the code to print when epoch+1 is divisible by 50, so that last epochs loss will be printed out)

```
Epoch 2000 : Ave objective= 0.30478222133992705  Ave training loss: 0.19979320133887904
```

```
Epoch 500 : Ave objective= 0.04630343290301937  Ave training loss: 0.06352717526036394
```

4 Multilayer Perceptron

Let's now turn to a multilayer perceptron (MLP) with a single hidden layer and a square loss. We'll implement the computation graph illustrated below:

The crucial new piece here is the nonlinear **hidden layer**, which is what makes the multilayer perceptron a significantly larger hypothesis space than linear prediction functions.

4.1 The standard non-linear layer

The multilayer perceptron consists of a sequence of “layers” implementing the following non-linear function

$$h(x) = \sigma(Wx + b),$$

where $x \in \mathbb{R}^d$, $W \in \mathbb{R}^{m \times d}$, and $b \in \mathbb{R}^m$, and where m is often referred to as the number of **hidden units** or **hidden nodes**. σ is some non-linear function, typically tanh or ReLU, applied element-wise to the argument of σ . Referring to the computation graph illustration above, we will implement this nonlinear layer with two nodes, one implementing the affine transform $L = W_1x + b_1$, and the other implementing the nonlinear function $h = \tanh(L)$. In this problem, we'll work out how to implement the backward method for each of these nodes.

The Affine Transformation

In a general neural network, there may be quite a lot of computation between any given affine transformation $Wx + b$ and the final objective function value J . We will capture all of that in a function $f: \mathbb{R}^m \rightarrow \mathbb{R}$, for which $J = f(Wx + b)$. Our goal is to find the partial derivative of J with respect to each element of W , namely $\partial J / \partial W_{ij}$, as well as the partials $\partial J / \partial b_i$, for each element of b . For convenience, let $y = Wx + b$, so we can write $J = f(y)$. Suppose we have already computed the partial derivatives of J with respect to the entries of $y = (y_1, \dots, y_m)^T$, namely $\frac{\partial J}{\partial y_i}$ for $i = 1, \dots, m$. Then by the chain rule, we have

$$\frac{\partial J}{\partial W_{ij}} = \sum_{r=1}^m \frac{\partial J}{\partial y_r} \frac{\partial y_r}{\partial W_{ij}}.$$

4. Show that $\frac{\partial J}{\partial W_{ij}} = \frac{\partial J}{\partial y_i} x_j$, where $x = (x_1, \dots, x_d)^T$. [Hint: Although not necessary, you might find it helpful to use the notation $\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases}$. So, for examples, $\partial_{x_j} (\sum_{i=1}^n x_i^2) = 2x_j$.]

Solution:

$$y = Wx + b$$

Consider only 1 output y_i and let d be the number of features

$$y_i = \sum_{j=1}^d W_{ij}x_j + b_i$$

Differentiate this wrt W_{ij}

$$\frac{\partial y_i}{\partial W_{ij}} = x_j$$

Now, for all outputs y_r other than y_i ($r \neq i$), the differentiation of y wrt W_{ij} is 0. This is because, those outputs are a feature of W_{rj} instead

So,

$$\frac{\partial y_r}{\partial W_{ij}} = \delta_{rj}x_j$$

Thus, when $r = i$, $\delta_{rj} = 1$, and so the differentiation is non-zero, otherwise it is 0

Substituting this value in the given equation

$$\frac{\partial J}{\partial W_{ij}} = \sum_{r=1}^m \frac{\partial J}{\partial y_r} \frac{\partial y_r}{\partial W_{ij}}$$

$$\frac{\partial J}{\partial W_{ij}} = \sum_{r=1}^m \frac{\partial J}{\partial y_r} \delta_{rj}x_j$$

Now again, when $r \neq i$, this expression is 0, and so the summation can be removed

$$\frac{\partial J}{\partial W_{ij}} = \frac{\partial J}{\partial y_i} x_j$$

Thus proved

5. Now let's vectorize this. Let's write $\frac{\partial J}{\partial y} \in \mathbb{R}^{m \times 1}$ for the column vector whose i th entry is $\frac{\partial J}{\partial y_i}$. Let's also define the matrix $\frac{\partial J}{\partial W} \in \mathbb{R}^{m \times d}$, whose ij 'th entry is $\frac{\partial J}{\partial W_{ij}}$. Generally speaking, we'll always take $\frac{\partial J}{\partial A}$ to be an array of the same size ("shape" in numpy) as A . Give a vectorized expression for $\frac{\partial J}{\partial W}$ in terms of the column vectors $\frac{\partial J}{\partial y}$ and x . [Hint: Outer product.]

Solution:

$\frac{\partial J}{\partial y}$ is a column vector, so has shape $(m,1)$

Similarly, x is a column vector, so has shape $(d,1)$ where d is the number of features

So, we take outer product of the 2 vectors by transposing x

$$\frac{\partial J}{\partial W} = \frac{\partial J}{\partial y} x^T$$

6. In the usual way, define $\frac{\partial J}{\partial x} \in \mathbb{R}^d$, whose i 'th entry is $\frac{\partial J}{\partial x_i}$. Show that

$$\frac{\partial J}{\partial x} = W^T \left(\frac{\partial J}{\partial y} \right)$$

[Note, if x is just data, technically we won't need this derivative. However, in a multilayer perceptron, x may actually be the output of a previous hidden layer, in which case we will need to propagate the derivative through x as well.]

Solution:

Consider only 1 value in the vector, using chain rule

$$\frac{\partial J}{\partial x_i} = \sum_{r=1}^m \frac{\partial J}{\partial y_r} \frac{\partial y_r}{\partial x_i} \dots (i)$$

$$y_r = \sum_{i=1}^d W_{ri} x_i + b_r$$

Differentiate wrt x_i

$$\frac{\partial y_r}{\partial x_i} = W_{ri}$$

Substitute this in equation (i)

$$\frac{\partial J}{\partial x_i} = \sum_{r=1}^m W_{ri} \frac{\partial J}{\partial y_r}$$

Now, this is the i^{th} row for the LHS

We find the i^{th} row for the RHS next, using rules of matrix multiplication. The matrix multiplication involves inner product of i^{th} row of first matrix and j^{th} column of the second matrix. But here, the second matrix is a column vector, so has only 1 column and so we get

$$W^T \left(\frac{\partial J}{\partial y} \right) = \left(\cdot \cdot \cdot \sum_{r=1}^m W_{ri} \frac{\partial J}{\partial y_r} \cdot \cdot \cdot \right)^T$$

Thus, the two sides LHS and RHS are equal, as the i^{th} row of both the column vectors is equal, and by definition, there is only 1 column, so the entire matrices are equal

Hence proved

7. Show that $\frac{\partial J}{\partial b} = \frac{\partial J}{\partial y}$, where $\frac{\partial J}{\partial b}$ is defined in the usual way.

Solution:

By chain rule, (consider only 1 row of the column vector, which is row i)

$$\frac{\partial J}{\partial b_i} = \sum_{r=1}^m \frac{\partial J}{\partial y_r} \frac{\partial y_r}{\partial b_i}$$

Also, $y_i = \sum_{j=1}^d W_{ij}x_j + b_i$ and so

$$\frac{\partial y_i}{\partial b_i} = 0 + 1 = 1$$

For all values of r other than i , this differentiation will be 0, and so

$$\frac{\partial J}{\partial b_i} = \sum_{r=1}^m \frac{\partial J}{\partial y_r} \delta_{ir}$$

$$\frac{\partial J}{\partial b_i} = \frac{\partial J}{\partial y_i}$$

Now we vectorize back this equation (as both LHS and RHS have same variable i as the index),

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial y}$$

Hence, proved

Element-wise Transformers

Our nonlinear activation function nodes take an array (e.g. a vector, matrix, higher-order tensor, etc), and apply the same nonlinear transformation $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ to every element of the array. Let's abuse notation a bit, as is usually done in this context, and write $\sigma(A)$ for the array that results from applying $\sigma(\cdot)$ to each element of A . If σ is differentiable at $x \in \mathbb{R}$, then we'll write $\sigma'(x)$ for the derivative of σ at x , with $\sigma'(A)$ defined analogously to $\sigma(A)$.

Suppose the objective function value J is written as $J = f(\sigma(A))$, for some function $f : S \mapsto \mathbb{R}$, where S is an array of the same dimensions as $\sigma(A)$ and A . As before, we want to find the array $\frac{\partial J}{\partial A}$ for any A . Suppose for some A we have already computed the array $\frac{\partial J}{\partial S} = \frac{\partial f(S)}{\partial S}$ for $S = \sigma(A)$. At this point, we'll want to use the chain rule to figure out $\frac{\partial J}{\partial A}$. However, because we're dealing with arrays of arbitrary shapes, it can be tricky to write down the chain rule. Appropriately, we'll use a tricky convention: We'll assume all entries of an array A are indexed by a single variable. So, for example, to sum over all entries of an array A , we'll just write $\sum_i A_i$.

8. Show that $\frac{\partial J}{\partial A} = \frac{\partial J}{\partial S} \odot \sigma'(A)$, where we're using \odot to represent the **Hadamard product**. If A and B are arrays of the same shape, then their Hadamard product $A \odot B$ is an array with the same shape as A and B , and for which $(A \odot B)_i = A_i B_i$. That is, it's just the array formed by multiplying corresponding elements of A and B . Conveniently, in **numpy** if **A** and **B** are arrays of the same shape, then **A*B** is their Hadamard product.

Solution:

A is a vector. $S = \sigma(A)$ is also a vector of the same size, and let it be s

Differentiate S wrt A

$$\frac{\partial S}{\partial A} = \sigma'(A)$$

Consider only 1 element in the resulting vector

$$\frac{\partial S_i}{\partial A_i} = \sigma'(A_i)$$

All other outputs S_j such that $j \neq i$, do not depend on the input A_i and so their differentiation wrt A_i is 0

$$\frac{\partial S_j}{\partial A_i} = \delta_{ij} \sigma'(A_i) \dots (i)$$

Now, consider only 1 element in the vector, and using chain rule

$$\frac{\partial J}{\partial A_i} = \sum_{j=1}^s \frac{\partial J}{\partial S_j} \frac{\partial S_j}{\partial A_i}$$

Now substitute equation (i) in this,

$$\frac{\partial J}{\partial A_i} = \sum_{j=1}^s \frac{\partial J}{\partial S_j} \delta_{ij} \sigma'(A_i)$$

Again, for all $j \neq i$, the differentiation is 0, and so we can get rid of the summation,

$$\frac{\partial J}{\partial A_i} = \frac{\partial J}{\partial S_i} \sigma'(A_i) \dots (\text{ii})$$

All terms in the above equation are scalars. Now, to vectorize this, we use definition of Hadamard product, which is that if A and B are arrays of the same shape, then their Hadamard product is an array of the same shape and each element in the array is a scalar multiplication of the elements of A and B

This is exactly what equation (ii) represents, as $\frac{\partial J}{\partial S_i}$ is a scalar as it represents 1 element in the vector $\frac{\partial J}{\partial S}$ and $\sigma'(A_i)$ is also a scalar which represents 1 element in $\sigma'(A)$, and so, we can say that

$$\frac{\partial J}{\partial A} = \frac{\partial J}{\partial S} \odot \sigma'(A)$$

4.2 MLP Implementation

9. Complete the class `AffineNode` in `nodes.py`. Be sure to propagate the gradient with respect to x as well, since when we stack these layers, x will itself be the output of another node that depends on our optimization parameters. If your code is correct, you should be able to pass `test_AffineNode` in `mlp_regression.t.py`. Please attach a screenshot that shows the test results for this question.

```
class AffineNode(object):

    def __init__(self, x, W, b, node_name):
        self.node_name = node_name
        self.out = None
        self.d_out = None
        self.W = W
        self.x = x
        self.b = b

    def forward(self):
        self.out = np.matmul(self.W.out, self.x.out) + self.b.out
        self.d_out = np.zeros(self.out.shape)
        return self.out

    def backward(self):
        d_W = np.matmul(self.d_out.reshape(-1,1), self.x.out.reshape(-1,1).T)

        d_x = np.matmul(self.W.out.T, self.d_out)
        d_x = d_x.reshape((d_x.shape[0],))

        d_b = self.d_out
        d_b = d_b.reshape((d_b.shape[0],))
        self.x.d_out += d_x
        self.W.d_out += d_W
        self.b.d_out += d_b

    def get_predecessors(self):
        return [self.W, self.x, self.b]
    pass
```

Test results for all 3 tests (as the .t.py file is not to be modified, running all 3 tests at once):

```
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. W is 1.4249409828450496e-08.
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. x is 7.958314233285048e-09.
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. b is 1.6365788156163885e-09.
.DEBUG: (Node tanh) Max rel error for partial deriv w.r.t. a is 4.170376077170414e-09.
.DEBUG: (Parameter W1) Max rel error for partial deriv 4.621822945148994e-06.
DEBUG: (Parameter b1) Max rel error for partial deriv 5.820443779285191e-08.
DEBUG: (Parameter w2) Max rel error for partial deriv 3.1920436296005552e-09.
DEBUG: (Parameter b2) Max rel error for partial deriv 8.27231812855547e-12.
.
-----
Ran 3 tests in 0.006s

OK

Process finished with exit code 0
```

10. Complete the class `TanhNode` in `nodes.py`. As you'll recall, $\frac{d}{dx} \tanh(x) = 1 - \tanh^2 x$. Note that in the forward pass, we'll already have computed \tanh of the input and stored it in `self.out`. So make sure to use `self.out` and not recalculate it in the backward pass. If your code is correct, you should be able to pass `test_TanhNode` in `mlp_regression.t.py`. Please attach a screenshot that shows the test results for this question.

```
class TanhNode(object):
    """Node  $\tanh(a)$ , where  $\tanh$  is applied elementwise to the array  $a$ 
    Parameters:
    a: node for which a.out is a numpy array
    """

    def __init__(self, a, node_name):
        self.node_name = node_name
        self.out = None
        self.d_out = None
        self.a = a

    def forward(self):
        self.out = np.tanh(self.a.out)
        self.d_out = np.zeros(self.out.shape)
        return self.out

    def backward(self):
        d_a = 1 - self.out**2
        self.a.d_out += self.d_out * d_a

    def get_predecessors(self):
        return [self.a]

    pass
```

Test results for all 3 tests (as the .t.py file is not to be modified, running all 3 tests at once):

```
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. W is 1.4249409828450496e-08.
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. x is 7.958314233285048e-09.
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. b is 1.6365788156163885e-09.
.DEBUG: (Node tanh) Max rel error for partial deriv w.r.t. a is 4.170376077170414e-09.
.DEBUG: (Parameter W1) Max rel error for partial deriv 4.621822945148994e-06.
DEBUG: (Parameter b1) Max rel error for partial deriv 5.820443779285191e-08.
DEBUG: (Parameter w2) Max rel error for partial deriv 3.1920436296005552e-09.
DEBUG: (Parameter b2) Max rel error for partial deriv 8.27231812855547e-12.
.
-----
Ran 3 tests in 0.006s

OK

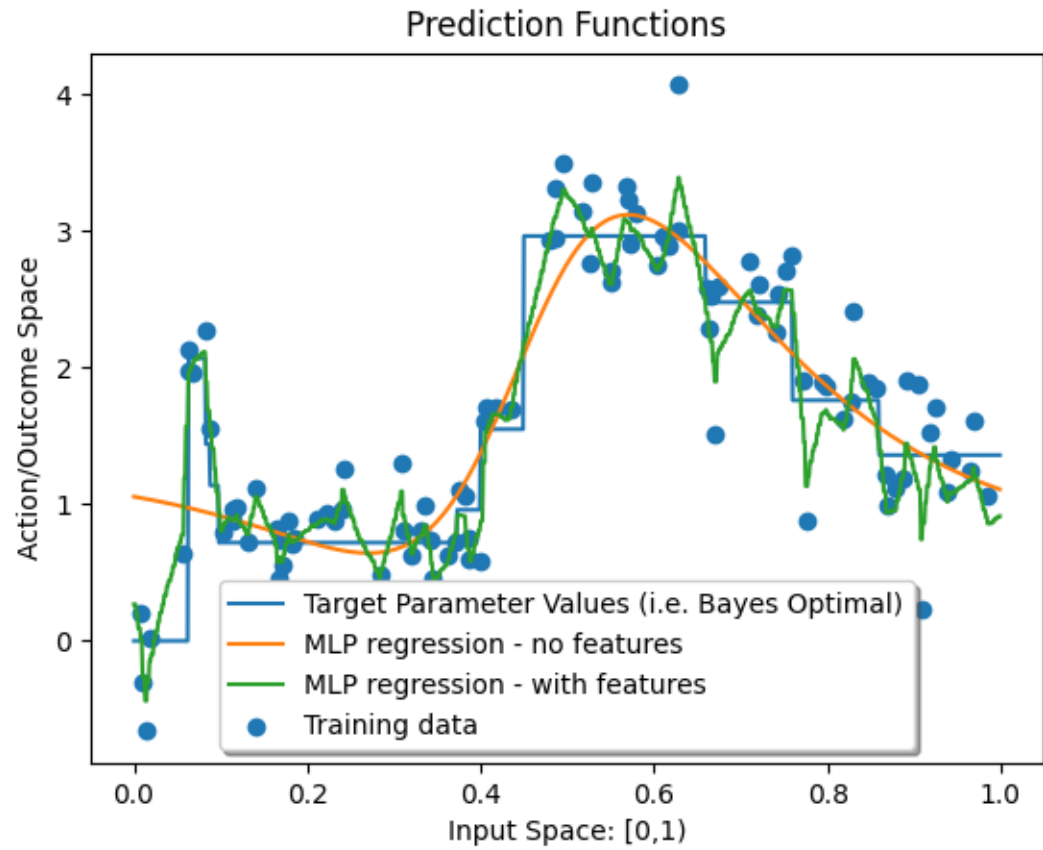
Process finished with exit code 0
```


Test results for all 3 tests (as the .t.py file is not to be modified, running all 3 tests at once):

```
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. W is 1.4249409828450496e-08.
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. x is 7.958314233285048e-09.
DEBUG: (Node affine) Max rel error for partial deriv w.r.t. b is 1.6365788156163885e-09.
.DEBUG: (Node tanh) Max rel error for partial deriv w.r.t. a is 4.170376077170414e-09.
.DEBUG: (Parameter W1) Max rel error for partial deriv 4.621822945148994e-06.
DEBUG: (Parameter b1) Max rel error for partial deriv 5.820443779285191e-08.
DEBUG: (Parameter w2) Max rel error for partial deriv 3.1920436296005552e-09.
DEBUG: (Parameter b2) Max rel error for partial deriv 8.27231812855547e-12.
.
-----
Ran 3 tests in 0.006s

OK

Process finished with exit code 0
```



Final average training square objective and loss for the 2 settings are: (made a minor change in the code to print when epoch+1 is divisible by 50, so that last epochs loss will be printed out)

```
Epoch 5000 : Ave objective= 0.2655365458740911  Ave training loss: 0.26255544601285247
```

```
Epoch 500 : Ave objective= 0.043098622923591144  Ave training loss: 0.04548697071963382
```