

Homework 7: Computation Graphs, Backpropagation, and Neural Networks

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. The reason so many people have success is that, thanks to GitHub, they can copy the exact settings that others have used to achieve success. 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 and suggestions 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 is implemented. With this assignment, you'll have an opportunity to linger on these low-level implementation details, which one usually rushes through in more specialized neural network classes. Every major neural network type (RNNs, CNNs, Resnets, etc) can be implemented using the basic building blocks we'll develop in this assignment.

To help things along, we¹ have designed a minimal 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.) In some sense, computation graphs have nothing

¹Philipp Meerkamp, Pierre Garapon, and David Rosenberg

to do with machine learning or neural networks. They are just a way to represent a function that facilitates efficient computation of the function value and it's gradients with respect to inputs. The tutorial takes this perspective, and there is very little 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` are a bit intricate so they can work with any node or `ComputationGraphFunction`, but it's using the exact same `gradient_checker` logic we saw in the first homework assignment.

Once you've understood how linear regression works in our framework, you're ready to start implementing your own algorithms...

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. As we discussed in class, 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. While we don't have this sharing in the multilayer perceptron, we do have it in RNNs and CNNs, the two neural network architectures that have had the most impact. So being able to handle ridge regression is a necessary condition (and possibly sufficient) for being able to represent these important architectures.

We've provided some skeleton code in `ridge_regression.py`, and some test code in `ridge_regression.t.py`, that you should eventually be able to pass.

1. Complete the class `L2NormPenaltyNode` in `nodes.py`.

```
1 class L2NormPenaltyNode(object):
2     """ Node computing  $l2\_reg * ||w||^2$  for scalars  $l2\_reg$  and vector  $w$  """
3     def __init__(self, l2_reg, w, node_name):
4         """
5         Parameters:
6         l2_reg: a scalar value  $\geq 0$  (not a node)
7         w: a node for which  $w.out$  is a numpy vector
8         node_name: node's name (a string)
9         """
10        self.node_name = node_name
11        self.out = None
12        self.d_out = None
```

```

13         self.l2_reg = np.array(l2_reg)
14         self.w = w
15
16     def forward(self):
17         self.out = self.l2_reg * np.dot(self.w.out, self.w.out)
18         self.d_out = np.zeros(self.out.shape)
19         return self.out
20
21     def backward(self):
22         d_w = 2*self.l2_reg*self.d_out*self.w.out
23         self.w.d_out += d_w
24         return self.d_out
25
26     def get_predecessors(self):
27         return [self.w]

```

2. Complete the class SumNode in nodes.py.

```

1 class SumNode(object):
2     """ Node computing a + b, for numpy arrays a and b """
3     def __init__(self, a, b, node_name):
4         """
5         Parameters:
6         a: node for which a.out is a numpy array
7         b: node for which b.out is a numpy array of the same shape as a
8         node_name: node's name (a string)
9         """
10        self.node_name = node_name
11        self.out = None
12        self.d_out = None
13        self.a = a
14        self.b = b
15
16    def forward(self):
17        self.out = self.a.out + self.b.out
18        self.d_out = np.zeros(self.out.shape)
19        return self.out
20
21    def backward(self):
22        d_a = self.d_out
23        d_b = self.d_out
24        self.a.d_out += d_a
25        self.b.d_out += d_b
26        return self.d_out
27
28    def get_predecessors(self):
29        return [self.a, self.b]

```

3. Implement ridge regression with w regularized and b unregularized. Do this by completing the `__init__` method in the `ridge_regression.py`, using the classes created above. When complete, you should be able to pass the tests in `ridge_regression.t.py`. Report the average square error on the **training** set for the parameter settings given in the `main()` function.

```

1 def __init__(self, l2_reg=1, step_size=.005, max_num_epochs = 5000):

```

```

2         self.max_num_epochs = max_num_epochs
3         self.step_size = step_size
4
5         # Build computation graph
6         self.l2_reg = l2_reg
7         self.x = nodes.ValueNode(node_name="x") # to hold a vector input
8         self.y = nodes.ValueNode(node_name="y") # to hold a scalar response
9         self.w = nodes.ValueNode(node_name="w") # to hold the parameter vector
10        self.b = nodes.ValueNode(node_name="b") # to hold the bias parameter (scalar
11    )
12
13    # TODO
14    self.prediction = nodes.VectorScalarAffineNode(x=self.x, w=self.w, b=self.b,
15                                                    node_name="prediction")
16    self.obj_reg = nodes.SquaredL2DistanceNode(a=self.prediction, b=self.y,
17                                              node_name="square loss")
18    self.obj_norm = nodes.L2NormPenaltyNode(l2_reg=self.l2_reg, w=self.w,
19                                           node_name="l2 penalty")
20    self.objective = nodes.SumNode(a=self.obj_reg, b=self.obj_norm,
21                                  node_name="penalized sq loss")
22
23    # TODO
24    self.inputs = [self.x]
25    self.outcomes = [self.y]
26    self.parameters = [self.w, self.b]
27
28    self.graph = graph.ComputationGraphFunction(self.inputs, self.outcomes,
29                                                self.parameters, self.
30        prediction,
31                                                self.objective)

```

Setting 1:

Epoch 1950 : Ave objective= 0.3081544879368522 Ave training loss: 0.199031856598395

setting 2:

Epoch 450 : Ave objective= 0.04005683874474115 Ave training loss: 0.0276461661108703

4. [Optional] Create a new implementation of ridge regression that supports efficient minibatching. You will replace the the ValueNode `x`, which contains a vector, with a ValueNode `X`, which contains a matrix. The convention is that the first dimension indexes examples and the second is features (as we have always done). Many of the nodes will have to be adapted to this use case. Demonstrate its use and speedup.

4 Multilayer Perceptron

In this problem, we'll be implement a multilayer perceptron, with a single hidden layer, for the regression setting. 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 $x \in \mathbf{R}^d$, $W \in \mathbf{R}^{m \times d}$, and $b \in \mathbf{R}^m$, 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 $L(x) = W_1x + b_1$, and the other implementing the nonlinear function $h(L) = \tanh(L)$. In this problem we’ll work out how to implement the backward method for each of these nodes.

4.1.1 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 $f : \mathbf{R}^m \rightarrow \mathbf{R}$, for which $J = f(Wx + b)$. Our goal is to find the partial derivative of J with respect to each element of W , namely $\partial J / \partial W_{ij}$. 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 intermediate variable $y = (y_1, \dots, y_m)^T$, namely $\frac{\partial J}{\partial y_i}$ for $i = 1, \dots, m$. Then by the chain rule, we have

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

1. 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, $\partial_{x_j} (\sum_{i=1}^n x_i^2) = 2x_j \delta_{ij} = 2x_j$.]
For $\frac{\partial J}{\partial W_{ij}} = \sum_{r=1}^m \frac{\partial J}{\partial y_r} \frac{\partial y_r}{\partial W_{ij}}$, only the item with $r = i$ is non-zero. $\implies \frac{\partial J}{\partial W_{ij}} = \frac{\partial J}{\partial y_i} \frac{\partial y_i}{\partial W_{ij}} = \frac{\partial J}{\partial y_i} x_j$
2. Now let’s vectorize this. Let’s write $\frac{\partial J}{\partial y} \in \mathbf{R}^{m \times 1}$ for the column vector whose i th entry is $\frac{\partial J}{\partial y_i}$. Let’s also define the matrix $\frac{\partial J}{\partial W} \in \mathbf{R}^{m \times d}$, 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.]
Based on the rules given in question, item in row i and column j should have expression $\frac{\partial J}{\partial y_i} x_j$.

$$\frac{\partial J}{\partial W_{ij}} = \frac{\partial J}{\partial y_i} x_j \implies \frac{\partial J}{\partial W_{ij}} = \frac{\partial J}{\partial y} \otimes x$$

3. In the usual way, define $\frac{\partial J}{\partial x} \in \mathbf{R}^d$, 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.]

$$\begin{aligned}\frac{\partial J}{\partial x_i} &= \sum_{k=1}^m \frac{\partial J}{\partial y_k} \frac{\partial y_k}{\partial x_i} = \sum_{k=1}^m \frac{\partial J}{\partial y_k} W_{ki} \\ \implies \frac{\partial J}{\partial x} &= \frac{\partial J}{\partial f(x)} \frac{\partial f(x)}{\partial x} = W^T \left(\frac{\partial J}{\partial y} \right)\end{aligned}$$

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

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

4.1.2 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 $\sigma : \mathbf{R} \rightarrow \mathbf{R}$ 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 $x \in \mathbf{R}$, 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 $f : S \mapsto \mathbf{R}$, 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 $\sum_i A_i$.

1. 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.

$$\begin{aligned}\frac{\partial J}{\partial A_i} &= \frac{\partial J}{\partial S_i} \frac{d\sigma(A_i)}{dA_i} \\ \implies \frac{\partial J}{\partial A} &= \frac{\partial J}{\partial S} \frac{\partial S}{\partial A} = \frac{\partial J}{\partial S} \odot \sigma'(A)\end{aligned}$$

4.2 MLP Implementation

1. 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.

```

1  class AffineNode(object):
2      """Node implementing affine transformation (W,x,b)-->Wx+b, where W is a matrix,
3      and x and b are vectors
4      Parameters:
5          W: node for which W.out is a numpy array of shape (m,d)
6          x: node for which x.out is a numpy array of shape (d)
7          b: node for which b.out is a numpy array of shape (m) (i.e. vector of length
            m)
8      """
9      ## TODO
10     def __init__(self, W,x,b,node_name):
11         self.node_name = node_name
12         self.out = None
13         self.d_out = None
14         self.W = W
15         self.x = x
16         self.b = b
17
18     def forward(self):
19         self.out = np.dot(self.W.out, self.x.out)+self.b.out
20         self.d_out = np.zeros(self.out.shape)
21         return self.out
22
23     def backward(self):
24         d_W = np.outer(self.d_out, self.x.out)
25         d_x = np.dot(self.W.out.T, self.d_out)
26         d_b = self.d_out
27         self.W.d_out+=d_W
28         self.x.d_out+=d_x
29         self.b.d_out+=d_b
30         return self.d_out
31     def get_predecessors(self):
32         return [self.W, self.x, self.b]

```

2. 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.

```

34 class TanhNode(object):
35     """Node tanh(a), where tanh is applied elementwise to the array a
36     Parameters:
37         a: node for which a.out is a numpy array
38     """
39     ## TODO
40
41     def __init__(self, a,node_name):
42         self.node_name = node_name
43         self.out = None
44         self.d_out = None
45         self.a = a
46
47     def forward(self):
48         self.out = np.tanh(self.a.out)
49         self.d_out = np.zeros(self.out.shape)
50         return self.out
51

```

```

52     def backward(self):
53         d_a = self.d_out*(1-self.out**2)
54         self.a.d_out+=d_a
55         return self.d_out
56
57
58     def get_predecessors(self):
59         return [self.a]

```

3. Implement an MLP by completing the skeleton code in `mlp_regression.py`, and making use of the nodes above. Your code should pass the tests provided in `mlp_regression.t.py`. Note that to break the symmetry of the problem, we initialize our weights to small random values, rather than all zeros, as we often do for convex optimization problems. Run the MLP for the two settings given in the `main()` function and report the average **training** error. Note that with an MLP, we can take the original scalar as input, in the hopes that it will learn nonlinear features on its own, using the hidden layers. In practice, it is quite challenging to get such a neural network to fit as well as one where we provide features.

```

1  class MLPRegression(BaseEstimator, RegressorMixin):
2      """ MLP regression with computation graph """
3      def __init__(self, num_hidden_units=10, step_size=.005, init_param_scale=0.01,
4                  max_num_epochs = 5000):
5          self.num_hidden_units = num_hidden_units
6          self.init_param_scale = 0.01
7          self.max_num_epochs = max_num_epochs
8          self.step_size = step_size
9
10         # Build computation graph
11         self.x = nodes.ValueNode(node_name="x") # to hold a vector input
12         self.y = nodes.ValueNode(node_name="y") # to hold a scalar response
13         self.w1 = nodes.ValueNode(node_name="w1")
14         self.w2 = nodes.ValueNode(node_name="w2")
15         self.b1 = nodes.ValueNode(node_name="b1")
16         self.b2 = nodes.ValueNode(node_name="b2")
17         self.affine = nodes.AffineNode(W=self.w1, x=self.x, b=self.b1,
18                                       node_name="affine")
19         self.tanh = nodes.TanhNode(a=self.affine, node_name="tanh")
20         self.prediction = nodes.VectorScalarAffineNode(x=self.tanh, w=self.w2, b=
21                                                         self.b2,
22                                                         node_name="prediction")
23         self.objective = nodes.SquaredL2DistanceNode(a=self.prediction, b=self.y,
24                                                       node_name="square loss")
25
26         self.inputs = [self.x]
27         self.outcomes = [self.y]
28         self.parameters = [self.w1, self.b1, self.w2, self.b2]
29
30         self.graph = graph.ComputationGraphFunction(self.inputs, self.outcomes,
31                                                       self.parameters, self.
32                                                       prediction,
33                                                       self.objective)
34
35     def fit(self, X, y):
36         num_instances, num_fts = X.shape
37         y = y.reshape(-1)

```



```

35         num_hidden_units = self.num_hidden_units
36         s = self.init_param_scale
37         ## TODO: Initialize parameters (small random numbers -- not all 0, to break
symmetry )
38         parameter_vals = {"W1": s*np.random.standard_normal((num_hidden_units,
num_fts)),
39                             "b1": s*np.random.standard_normal((num_hidden_units)),
40                             "w2": s*np.random.standard_normal((num_hidden_units)),
41                             "b2": s*np.array(np.random.randn()) }
42
43         self.graph.set_parameters(parameter_vals)
44
45         for epoch in range(self.max_num_epochs):
46             shuffle = np.random.permutation(num_instances)
47             epoch_obj_tot = 0.0
48             for j in shuffle:
49                 obj, grads = self.graph.get_gradients(input_values = {"x": X[j]},
outcome_values = {"y": y[j]})
50
51                 #print(obj)
52                 epoch_obj_tot += obj
53                 # Take step in negative gradient direction
54                 steps = {}
55                 for param_name in grads:
56                     steps[param_name] = -self.step_size * grads[param_name]
57                 self.graph.increment_parameters(steps)
58
59             if epoch % 50 == 0:
60                 train_loss = sum((y - self.predict(X,y)) **2)/num_instances
61                 print("Epoch ",epoch, ": Ave objective=",epoch_obj_tot/num_instances,
"Ave training loss: ",train_loss)

```

Setting 1:

Epoch 4950 : Ave objective= 0.21758927284741764 Ave training loss: 0.20883813853643404

Setting 2:

Epoch 450 : Ave objective= 0.03750042435029375 Ave training loss: 0.024624794808589402

4. [Optional] See if you can get a fit on the training set with an MLP that uses just the scalar input that is about as good as the fit using the featurized inputs. You can do that by tweaking model parameters (e.g. the number of hidden nodes or layers) and/or the parameters of optimization. You **may use** any neural network framework (PyTorch, TensorFlow, etc), which can help by providing more advanced optimization techniques (e.g. Adam), variable initialization methods, and/or various normalization approaches (batch norm, etc).

4.3 [OPTIONAL]

1. [Optional] Implement a Softmax node.
2. [Optional] Implement a negative log-likelihood loss node for multiclass classification.
3. [Optional] Use the classes above to apply an MLP to the simple multiclass classification dataset we had on a previous assignment.