

# CS 577 — Deep Learning — Final exam practice problems

## 1 Automatic differentiation for scalars

### 1.1 $x^4$

Draw the computational graph for the following function. Then compute `x.grad` using backpropagation.

```
1 x = ag.Scalar(2.0, label="z1\nleaf(x)")
2 z1 = x
3 z2 = z1*z1
4 z3 = z2*z2
5 z3.backward()
6 print(x.grad)
```

### 1.2 $\log(\exp(xy))$

Draw the computational graph for the following function. Then compute `x.grad` and `y.grad` using backpropagation.

```
7 x = ag.Scalar(2.0, label="z1\nleaf(x)")
8 y = ag.Scalar(3.0, label="z2\nleaf(y)")
9
10 z1 = x
11 z2 = y
12 z3 = z1*z2
13 z4 = ag.exp(z3)
14 z5 = ag.log(z4)
15 z5.backward()
16 print(x.grad, y.grad)
```

### 1.3 Recurrent neural networks

Draw the computational graph for the following function. Then compute `wr.grad`, `wi.grad`, and `wo.grad` using backpropagation.

```
17 x1 = ag.Scalar(2.0, label="z1\nleaf(x1)")
18 h0 = ag.Scalar(3.0, label="z2\nleaf(h0)")
19 wr = ag.Scalar(4.0, label="z3\nleaf(wr)")
20 wi = ag.Scalar(5.0, label="z4\nleaf(wi)")
21 wo = ag.Scalar(6.0, label="z5\nleaf(wo)")
22
23 z1 = x1
24 z2 = h0
25 z3 = wr
26 z4 = wi
27 z5 = z3*z2 # wr*h0
28 z6 = z4*z1 # wi*x1
29 z7 = z5+z6
30 z8 = ag.relu(z7) # relu(wr*h0 + wi*x1)
31 z9 = wo
32 z10 = z8*z9
33 z10.backward()
34
35 print(wr.grad, wi.grad, wo.grad)
```

## 2 Convolutional networks

### 2.1 im2col

What is the shape of the `im2col_mat_sparse` matrix constructed below? Give your answer in terms of `N`, `Cin`, `Hin`, `Win`, `K` and `S`.

```
36 def im2col_matrix_sparse(Xin, K, S=1):
37     N, Cin, Hin, Win = Xin.shape
38     CHW = Cin * Hin * Win
39     Hout = (Hin - K) // S + 1
40     Wout = (Win - K) // S + 1
41     P = Hout * Wout # Total number of patches per image
42     patch_size = Cin * K * K # Size of each flattened patch
43
44     data = [1 for _ in range(P*patch_size)]
45     row_indices = []
46     col_indices = list(range(P*patch_size))
47
48     patch_idx = 0
49     for hout in range(Hout):
50         for wout in range(Wout):
51             for cin in range(Cin):
52                 for hker in range(K):
53                     for wker in range(K):
54                         input_index = cin * Hin * Win + hout * S * Win + wout * S + hker *
55                         Win + wker
56                         row_indices.append(input_index)
57                         patch_idx += 1
58
59     im2col_mat_sparse = csr_matrix((data, (row_indices, col_indices)), shape=(CHW, P *
60     patch_size))
61     return im2col_mat_sparse
```

## 2.2 Convolution layer time and space complexities

- What is the time complexity of computing `Xout_flat` in the 2D convolution layer defined below? Give your answer in terms of `N`, `Cin`, `Cout`, `Hin`, `Win`, `K` and `S`.
- How much memory is required to store `Xin_im2col`? Give your answer in terms of bytes. (Assume entries of the array are in `float64`, which requires 8 bytes)

```
60 class Conv2d(Module):
61     def __init__(self, in_channels, out_channels, kernel_size, stride=1):
62         super().__init__()
63         self.in_channels = in_channels # Cin
64         self.out_channels = out_channels # Cout
65         self.kernel_size = kernel_size # K
66         self.stride = stride # S
67
68         kaiming_he_init_constant = np.sqrt(2 / (in_channels * kernel_size**2))
69
70         self.weight = ag.Tensor(np.random.randn(in_channels*kernel_size**2, out_channels
71 ) * kaiming_he_init_constant)
72
73         self.bias = ag.Tensor(np.zeros(out_channels))
74
75         self._parameters['weight'] = self.weight
76         self._parameters['bias'] = self.bias
77
78         self.im2col_mat = None # im2col_mat will be cached
79
80     def forward(self, Xin): # Xin.shape == (N, Cin, Hin, Win)
81         N, Cin, Hin, Win = Xin.shape
82         assert(Cin == self.in_channels)
83
84         K = self.kernel_size
85         S = self.stride
86         Hout = (Hin - K) // S + 1
87         Wout = (Win - K) // S + 1
88         P = Hout * Wout # Total number of patches per image
89         patch_size = Cin * K * K # Size of each flattened patch
90         Cout = self.out_channels
91
92         Xin_flat = Xin.reshape(-1, Cin * Hin * Win)
93
94         if self.im2col_mat is None: # Cache the im2col matrix
95             self.im2col_mat = im2col_matrix_sparse(Xin, K, S)
96
97         Xin_im2col = ag.spcmatmul(Xin_flat, self.im2col_mat)
98         Xin_patches_flat = Xin_im2col.reshape(N, P, patch_size)
99
100         Xout_flat = (Xin_patches_flat @ self.weight) + self.bias
101         Xout_flat = ag.moveaxis(Xout_flat, 1, 2)
102         Xout = Xout_flat.reshape(N, Cout, Hout, Wout)
103
104         return Xout # Xout.shape == (N, Cout, Hout, Wout)
```

### 3 Attention

- How much memory is required to store the following intermediate activations in SingleHeadAttention portion of the transformer block?: Queries, Keys, KQ, expKQ, softmaxKQ, P.
- How about the MLP? Namely how much memory is required to store the following?: hidden

```
104 class SingleHeadAttention:
105     def __init__(self, n_features):
106         self.Wq = ag.Tensor(np.random.randn(n_features, n_features), label="Wq")
107         self.Wk = ag.Tensor(np.random.randn(n_features, n_features), label="Wk")
108         self.Wv = ag.Tensor(np.random.randn(n_features, n_features), label="Wv")
109     def __call__(self, Xin):
110         # Xin is a (n_samples, n_context, n_features) tensor
111         # Xout is *also* a (n_samples, n_context, n_features) tensor
112         Queries = Xin @ self.Wq
113         Keys = Xin @ self.Wk
114         KQ = (Keys @ ag.moveaxis(Queries, 1,2))
115         expKQ = ag.exp(KQ)
116         softmaxKQ = expKQ / ag.sum(expKQ, axis=1, keepdims=True)
117         P = ag.moveaxis(Xin,1,2) @ softmaxKQ
118         Xout = ag.moveaxis(P, 1,2) @ self.Wv
119         return Xout
120
121 class MLP:
122     def __init__(self, n_features, n_hidden):
123         self.Wh = ag.Tensor(np.random.randn(n_features, n_hidden), label="Whidden")
124         self.bh = ag.Tensor(np.random.randn(n_hidden), label="bhidden")
125         self.wo = ag.Tensor(np.random.randn(n_hidden, n_features), label="Wout")
126         self.bo = ag.Tensor(np.random.randn(n_features), label="bout")
127
128     def __call__(self, Xin):
129         hidden = ag.relu((Xin @ self.Wh) + self.bh)
130         return hidden @ self.wo + self.bo
131
132 class TransformerBlock:
133     def __init__(self, n_features, n_hidden):
134         self.att = SingleHeadAttention(n_features)
135         self.mlp = MLP(n_features, n_hidden)
136     def __call__(self, Xin):
137         return self.mlp(self.att(Xin))
```

## 4 Tensor rematerialization (aka checkpointing)

### 4.1 No rematerialization

How much memory (in MB) does computing the forward function below require for the `.value` fields for the intermediate activations (so excluding the input and the output `ag.sum(x)`)?

```
138 class Tensor: # Tensor with grads
139     def __init__(self,
140                   value,
141                   requires_grad=False,
142                   rematerializer = None, # None means don't rematerialize, i.e., keep
143                   op="",
144                   _backward= lambda : None,
145                   inputs=[],
146                   label=""):
147
148         if type(value) in [float ,int]:
149             value = np.array(value)
150
151         self.requires_grad = requires_grad
152         self.rematerializer = rematerializer
153
154         self.value = 1.0*value
155         self.grad = None
156
157         if self.requires_grad:
158             self.grad = np.zeros_like(self.value)
159
160 # [...]
161
162 num_layers = 10
163 num_samples = 4096
164 dim_hidden = 1000
165
166 weights = [ag.Tensor(0.02*np.random.randn(dim_hidden, dim_hidden),
167                      requires_grad = True) for _ in range(num_layers)]
168 X = ag.Tensor(np.random.randn(num_samples, dim_hidden))
169
170 def forward(x, weights):
171     for w in weights:
172         x = ag.matmul(x, w)
173     return ag.sum(x)
```

## 4.2 With rematerialization

One way to implement rematerialization for the model considered in the previous part is to pass a “rematerializer” to the intermediate tensor to reconstruct the `.value` field:

```
174 def forward_with_rematerializer(x, weights, checkpoints):
175
176     farthest_checkpoint = 0
177     x_at_farthest_checkpoint = x
178
179     for i, w in enumerate(weights):
180         if i in checkpoints:
181             x = ag.matmul(x, w)
182             farthest_checkpoint = i
183             x_at_farthest_checkpoint = x
184         else:
185             def _rematerializer():
186                 xval = x_at_farthest_checkpoint.value
187                 for w in weights[farthest_checkpoint:(i+1)]:
188                     xval = np.matmul(xval, w.value)
189                 return xval
190             x = ag.matmul(x, w)
191             x.rematerializer = _rematerializer
192
193     return ag.sum(x)
```

Once the `.value` field is no longer needed during the forward computation, it is discarded:

```
194 # inside ag.Tensor class definition:
195     def discard_value_if_has_rematerializer(self):
196         if self.rematerializer is not None:
197             self.value = None
198         return None
```

How much memory (in MB) is used for allocating the `.value` fields if we call `forward_with_rematerializer` with the following inputs?

```
199 num_layers = 10
200 num_samples = 4096
201 dim_hidden = 1000
202
203 weights = [ag.Tensor(0.02*np.random.randn(dim_hidden, dim_hidden),
204                 requires_grad = True) for _ in range(num_layers)]
205 x = ag.Tensor(np.random.randn(num_samples, dim_hidden))
206 checkpoints = [5]
```