

Problem Set 3

Julian Vanecek

October 2018

1 Flowers

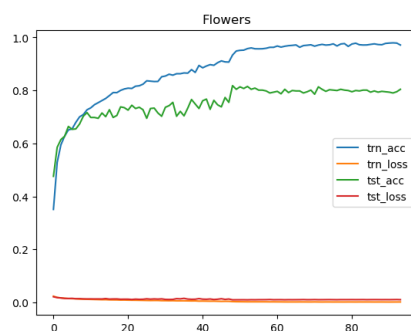


Figure 1: Flowers training

Flowers is using VGG architecture 64, 64, 'M', 128, 128, 'M', 256, 256, 'M', 512, 512, 'M' with batch normalization at every conv. layer no dropout, kernel sizes of 3, padding of 1, and no dropout. The classification layer has 64 hidden nodes, dropout(0.5) and batch normalization. It reshuffles the train data at every epoch, uses an optimizer SGD with momentum 0.9, weight decay, and a learning rate decrease of gamma 0.1 at 50 and 75%. I'm using the class loader and all the class tricks. In the end, this model was a pain to train because I had a bunch of stupid bugs that should have thrown errors but didn't (a.k.a. one time a rounding error, and another my output softmax dimension off by one error, so straight off the bat accuracy+loss was screwed), that encouraged fruitless exploration of model space. I ran into memory out of bounds problems, especially for the larger VGG networks. I believe I solved these with wrapping my `test_model()` in a `"with torch.no_grad()"` clause. So after overcoming these, when the default class implementation worked I felt no reason to improve or innovate beyond my previous adventures. I'm not exactly sure on why the stepwise increase of feature dimensions in VGG would make intuitive sense, except that the step would refine the information for that many features. I guess one conv. layer would take a meta look at what features were extracted before expanding into a higher feature space, thereby improving the quality.

```

[[{"name": "[ 0.8607, -0.1999, -0.1390],",
  [-0.0635, -0.0465, -0.0705],
  [-0.0909, -0.2904, -0.1403]],
 [{"name": "0.2509, 0.0829, 0.0742],",
  [-0.0718, -0.2552, -0.1445],
  [-0.0505, -0.1901, -0.1593]],
 [{"name": "0.1371, -0.0411, 0.0162],",
  [-0.0888, -0.1451, -0.1101],
  [-0.0277, -0.1079, -0.1556]], device="cuda:0",
  [-0.1077, -0.1108, -0.1041],
  [-0.0016, -0.1100, -0.1041]],
 [{"name": "0.0796, -0.1109, 0.1100],",
  [-0.0009, -0.0909, -0.0107],
  [-0.0939, -0.1091, -0.1011]],
 [{"name": "0.0116, -0.0119, 0.0540],",
  [-0.0108, -0.1754, -0.0114],
  [-0.0409, -0.1619, -0.1001]], device="cuda:0",
  [-0.1000, -0.0011, -0.0105],
  [-0.1000, -0.1001, -0.1001]],
 [{"name": "0.0066, -0.1155, 0.0000],",
  [-0.0001, -0.0001, -0.0001],
  [-0.0021, -0.1149, -0.1054]],
 [{"name": "0.0754, -0.1158, -0.0576],",
  [-0.0109, -0.1056, -0.1056],
  [-0.0706, -0.1119, -0.1050]], device="cuda:0",
  [-0.1177, -0.1181, -0.0576],
  [-0.0011, -0.0109, -0.1100]],
  17  0.001788  0.000000  0.110000

```

Figure 2: Some Flowers conv. weights (not zero)

1.1 Failures

To get a rep. of which flowers failed, I was having trouble with displaying the input flowers, because their RGB was scaled to a gaussian distribution (and rescaling was having problems), and showing the original was also difficult because the dataset was being shuffled at every epoch, so here's the best I could do: There is no obvious correlation of a single class being overly predicted or

target:1	output:0
target:1	output:2
target:1	output:2
target:4	output:2
target:3	output:4
target:4	output:0
target:4	output:2
target:4	output:2
target:1	output:4
target:2	output:0
target:4	output:2
target:1	output:0
target:0	output:2
target:4	output:2
target:0	output:2
target:1	output:0
target:0	output:1
target:3	output:0
target:2	output:4
target:0	output:2

Figure 3: 20 failures of flower classification

misclassified.

2 Three Meter

I am not confident in my loss function being correct, or if the produced losses are off by some scalar. From Piazza L1 should work, if you divide it by the number of batches. I don't get good enough loss, but making the network larger or smaller in the number of layers and in the nodes in each layer, in order to expand or compress the data, doesn't improve the loss considerably, so I chose a smaller network. The optimizer is SGD which seems to work better than Adam. No weight decay, but Nesterov momentum. The learning rate decays at 50 and 75% again. You can see a bump in the training loss when the decay hits. The data is centered around a normal distribution. The model is [Relu(33), Relu(24), 16, Relu(16), Relu(24), 33]. I think a slow compression is standard. Looking at which inputs failed in the network, I couldn't find anything abnormal. Many have a few inputs which are 2 std dev away from mean, but that isn't a statistical abnormality with 33 input features.

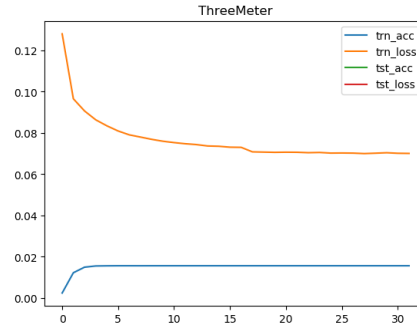


Figure 4: Three Meters training

```
[Tensor[[-0.0702, -0.0105, 0.0534, 0.0640, 0.1191, 0.1215, -0.0154,
-0.0001, -0.0164, 0.1137, 0.1189, -0.0012, 0.1239, 0.0017,
-0.0177, 0.1202, 0.1144, 0.1441, 0.0412, 0.0019, -0.0017,
0.0099, 0.0009, -0.0197, 0.0172, -0.0712, 0.0197, 0.0001,
-0.1197, -0.0209, 0.1121, 0.0918, 0.0001], device='cuda:0'],
-0.0075, 0.0009, 0.1744, 0.1759, 0.1201, 0.0108, -0.1187,
-0.1008, 0.0058, 0.0109, 0.1152, 0.1879, 0.0008, -0.0008,
0.0004, 0.0018, 0.0008, 0.0151, 0.0161, 0.0011, -0.0011,
0.0008, 0.0116, 0.0040, 0.1774, 0.0041], device='cuda:0'],
0.1770, 0.0448, 0.1409, 0.0151, 0.0001, 0.0011, -0.0017,
-0.0134, 0.1412, 0.0109, 0.1119, 0.0042, -0.0134, 0.0108,
0.0100, 0.1170, 0.0044, 0.0070, 0.0001, 0.0070, -0.0137,
0.1770, 0.1100, 0.0044, 0.0172, 0.0001, device='cuda:0'],
0.0450, 0.0709, 0.0164, 0.0478, 0.0009, 0.0045, -0.1072,
-0.0121, 0.0070, 0.1170, 0.0250, 0.0709, 0.0002, -0.1184,
-0.1757, 0.0166, 0.1451, 0.0171, 0.0108, 0.0155, -0.0259,
0.0075], device='cuda:0'], device='cuda:0']
```

Figure 5: Some Three Meter weights (not zero)

3 Adult

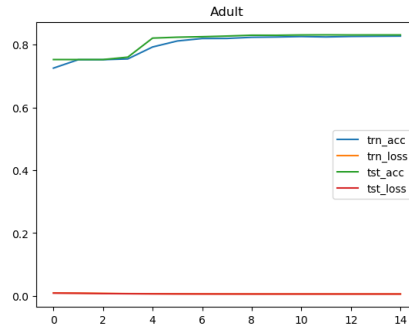


Figure 6: Adult Trainin

I use BCE because it works as a binary softmax if you will, and because it keeps the gradient high with a wrong classification, even though sigmoid gradient is usually zero, making it train faster. Anyways, I use a huge number of nodes, two hidden FC layers of 128 nodes, and a dropout of 0.2. I ended with a test accuracy of 0.83107 and a test loss of 0.00565 after 15 epochs.

3.1 Data

I mapped the features between zero and one. This worked well most of the time because most of the features were one hot encodings, which got mapped to 1, or something like age that exists as an int on a logical normal distribution, well I mapped that to $[0,1]$ as well because, hey, close enough.

3.2 Weights

Parameter	Estimate	Standard Error	z-Statistic	p-Value	95% CI
α_{11}	-0.001	0.001	-0.13	0.89	[-0.002, 0.001]
α_{12}	-0.001	0.002	-0.05	0.96	[-0.003, 0.001]
α_{13}	0.001	0.001	0.10	0.92	[-0.001, 0.002]
α_{14}	0.000	0.001	0.01	0.99	[-0.001, 0.001]
α_{15}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{16}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{17}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{18}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{19}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{21}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{22}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{23}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{24}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{25}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{26}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{27}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{28}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{29}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{31}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{32}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{33}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{34}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{35}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{36}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{37}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{38}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{39}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{41}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{42}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{43}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{44}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{45}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{46}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{47}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{48}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{49}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{51}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{52}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{53}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{54}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{55}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{56}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{57}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{58}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{59}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{61}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{62}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{63}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
α_{64}	0.000	0.001	0.02	0.98	[-0.001, 0.001]
<					

Figure 7: Some Adult weights (not zero)

3.3 Failures

The sum of the 20 failures makes

8.3973,	12.7333,	0.7510,	1.2745,	9.2755,	1.0000,
3.0000,	11.0000,	1.0000,	3.0000,	1.0000,	0.0000,
0.0000,	0.0000,	0.0000,	0.0000,	0.0000,	1.0000,
0.0000,	1.0000,	1.0000,	3.0000,	0.0000,	8.0000,
2.0000,	0.0000,	1.0000,	3.0000,	2.0000,	0.0000,
14.0000,	0.0000,	4.0000,	0.0000,	0.0000,	1.0000,
0.0000,	3.0000,	2.0000,	1.0000,	0.0000,	1.0000,
1.0000,	0.0000,	3.0000,	0.0000,	5.0000,	0.0000,
3.0000,	12.0000,	3.0000,	1.0000,	1.0000,	1.0000,
2.0000,	0.0000,	0.0000,	3.0000,	0.0000,	17.0000,
3.0000,	17.0000,	20.0000,	0.0000,	0.0000,	0.0000,
0.0000					

Figure 8: Adult misclassifications, summation of 20 vectors.

References