

EEC 289Q Data Analytics for Computer Engineers

Homework 3

Ahmed Mahmoud

May, 7th 2018

Softmax Regression:

The following code shows the implementation of the softmax regression function

```
1 function [f,g] = softmax_regression(theta, X,y)
2     m=size(X,2);
3     n=size(X,1);
4     theta=reshape(theta, n, []);
5     num_classes=size(theta,2)+1;
6     f = 0;
7     g = zeros(size(theta));
8     %%% YOUR CODE HERE %%%
9     theta_x = exp(theta'*X);
10    sum_col = sum(theta_x);
11    for I=1:m
12        theta_x(:,I) = theta_x(:,I)/sum_col(I);
13    end
14    for I=1:m
15        if y(I) < num_classes
16            f = f - log(theta_x(y(I),I));
17        end
18    end
19    %expand y to matrix to allow matrix multiply to obtain the gardient
20    y_mat = full(sparse(y,1:m,1));
21    g = -X * (y_mat(1:num_classes-1,:) - theta_x)';
22    g=g(:); % make gradient a vector for minFunc
```

Using this code, we were able to achieve training accuracy of 87.2% and test accuracy of 87.6% while the optimization took 4.524218 seconds. We used the gradient checker on our implementation and the average absolute error was 0.0330 (of 10 tests).

Supervised Neural Networks:

The following code implements the cost function, forward propagation, and compute the gradients for multiple hidden layers neural network

```
1 function [ cost, grad, pred_prob] = supervised_dnn_cost( theta, ei, ...
    data, labels, pred_only)
2     %SPNETCOSTSLAVE Slave cost function for simple phone net
3     % Does all the work of cost / gradient computation
4     % Returns cost broken into cross-entropy, weight norm, and prox reg
5     % components (ceCost, wCost, pCost)
6
7     %% default values
8     po = false;
9     if exist('pred_only','var')
10         po = pred_only;
11     end
12
13     %% reshape into network
14     stack = params2stack(theta, ei);
15     numHidden = numel(ei.layer_sizes) - 1;
16     hAct = cell(numHidden+1, 1);
17     gradStack = cell(numHidden+1, 1);
18     %% forward prop
19     %%% YOUR CODE HERE %%%
20     for J=1:numHidden
21         if J==1
22             %input
23             z = stack{J}.W*data;
24         else
25             %activation
26             z = stack{J}.W*hAct{J-1};
27         end
28         z = z + stack{J}.b;
29         hAct{J}=sigmoid(z);
30     end
31     z = stack{numHidden+1}.W*hAct{numHidden};
32     z = z + stack{numHidden+1}.b;
33     E = exp(z);
34     pred_prob = E./sum(E,1);
35     hAct{numHidden+1} = pred_prob;
36
37     %% return here if only predictions desired.
38     if po
39         cost = -1; ceCost = -1; wCost = -1; numCorrect = -1;
40         grad = [];
41         return;
42     end
43     %% compute cost
44     %%% YOUR CODE HERE %%%
45     c = log(pred_prob);
46     ind =sub2ind(size(c), labels', 1:size(c,2));
47     values = c(ind);
48     ceCost = -sum(values);
49     %% compute gradients using backpropagation
```

```

50     %%% YOUR CODE HERE %%%
51     d = zeros(size(pred_prob));
52     d(ind)=1;
53     error = (pred_prob-d);
54     for l =numHidden+1:-1:1
55         gradStack{l}.b = sum(error,2);
56         if l ==1
57             gradStack{l}.W = error*data';
58             break;
59         else
60             gradStack{l}.W = error*hAct{l-1}';
61         end
62         error = (stack{l}.W)'*error.*hAct{l-1}.*(1-hAct{l-1});
63     end
64     %% compute weight penalty cost and gradient for non-bias terms
65     %%% YOUR CODE HERE %%%
66     wCost = 0;
67     for l = 1:numHidden+1
68         wCost = wCost + 0.5*ei.lambda*sum(stack{l}.W(:).^2);
69     end
70     cost = ceCost + wCost;
71     for l=numHidden:-1:1
72         gradStack{l}.w = gradStack{l}.W + ei.lambda*stack{l}.W;
73     end
74     %% reshape gradients into vector
75     [grad] = stack2params(gradStack);
76 end

```

Using this code, we were able to get train accuracy of 100% and test accuracy of 0.9712% with weight decay value of 0.0. We can have a better test accuracy by changing the weight decay value to 0.25 which gave test accuracy of 0.973800% while the train accuracy lowered to 0.997717%. We think this is due to over-fitting with weight decay of 0. We used different values of weight decay with two and four hidden layers but not has shown superior performance. Also, changing/reduces the layer size always produced lower accuracy.

AlexNet:

First Layer: The input to AlexNet is a $[227 \times 227 \times 3]$ image (weight W). The first convolutional layer has receptive field (F) of 11, stride (S) of 4 and no zero-padding (P) with 96 kernels (K) (or convolutional layer depth of 96). Thus, the output is

$$\frac{W - F + 2P}{S} + 1 = \frac{227 - 11 + 2 * 0}{4} + 1 = 55$$

The number of neurons in this layer is $55 * 55 * 96 = 290,400$. Each of the 55×55 slice uses a unique $11 * 11 * 3 = 363$ weights and 1 bias. Thus, the total number of parameters in the first layer is $96 * 363 + 96 = 34,944$ parameters.

The operations for this layers are as follows: in order to obtain the one value in the $55 \times 55 \times 96$ output, the filter application will require 11×11 multiply operations and 11 add operations to do the dot product. This will be done 3 times for the three color channels in the input and add up all the result together (3 multiply) in addition to 1 add for the bias. Thus, the number of operations for one entry is $((11 * 11 + 11) * 3) + 1 = 397$ operations. For the whole layer, the total number of operations is $(55 * 55 * 96) * 397 = 115,288,800$.

Note: If we consider multiply and add as one operation (taking one cycle), then the number of operations per one entry is reduced to be $11 * 11 * 3 + 1 = 364$. This will make the total number of computation in this layer to be $(55 * 55 * 96) * 364 = 105,705,600$.

Second Layer: After the max pooling in the first layer, the input to the second layer becomes $[27 \times 27 \times 96]$. Second convolutional layer has receptive field (F) of 5, stride (s) of 1 and zero-padding (P) of 2 with 256 kernels (K). Thus, the output is

$$\frac{W - F + 2P}{S} + 1 = \frac{27 - 5 + 2 * 2}{1} + 1 = 27$$

The number of neurons in this layers is $27 * 27 * 256 = 186,624$. Each of the 27×27 slice uses a unique set of weight of size $5 * 5 * 96 = 2,400$ weights and 1 bias. Thus, the total number of parameters in the second layer is $256 * 2,400 + 256 = 614,656$ parameters.

Following the same computation we have done for the first layer, the total number of computation done in the second layer is $(27 * 27 * 256) * (((5 * 5 + 5) * 96) + 1) = 537,663,744$.

Third Layer: After max pooling, the input to the third layer is $[13 \times 13 \times 256]$. The third convolutional layer has receptive field of (F) 3, stride (S) of 1 and zero-padding (P) of 1 with 384 kernels (K). Thus, the output is

$$\frac{W - F + 2P}{S} + 1 = \frac{13 - 3 + 2 * 1}{1} + 1 = 13$$

The number of neurons in this layer is $13 * 13 * 384 = 64,896$. Each of the 13×13 slice uses a unique set of weight of size $3 * 3 * 256 = 2,304$ weights and 1 bias. Thus, the total number of parameters in the third layer is $2,304 * 384 + 384 = 885,120$ parameters.

The total computation in the third layer is $(13 * 13 * 384) * (((3 * 3 + 3) * 256) + 1) = 199,425,408$.

Fourth Layer: The output of third layer goes straight to fourth layer. The fourth convolutional layer has receptive field of (F) 3, stride (S) of 1 and zero-padding (P) of 1 with 384 kernels (K). The output size is (similar to third layer) 13.

The number of neurons is $13 * 13 * 384 = 64,896$. Each of this 13×13 slice uses a unique set of weights of size $3 * 3 * 384 = 3,456$ weights and 1 bias. Thus, the total number of parameters in the fourth layer is $3,456 * 384 + 384 = 1,327,488$.

The total computation in the forth layer is

$$(13 * 13 * 384) * (((3 * 3 + 3) * 384) + 1) = 299,105,664.$$

Fifth Layer: The output of fourth layers goes straight to the fifth layer. The fifth convolutional layer has receptive field of (F) 3, stride (S) of 1 and zero-padding (P) of 1 with 256 kernels (K). The output size is 13.

The number of neurons is $13 * 13 * 256 = 43,264$. Each of this 13×13 slice uses a unique set of weight of size $3 * 3 * 384 = 3,456$ weights and 1 bias. Thus, the total number of parameters in the fifth layer is $3,456 * 256 + 256 = 884,992$.

The total computation in the fifth layer is $(13*13*256)*(((3 * 3 + 3) * 384) + 1) = 199,403,776$.

Sixth Layer: After max pooling in the fifth layer, the input to the sixth fully connect (FC) layer is $[6 \times 6 \times 256]$. The total number of neurons of the sixth layer is 4096. Since it is fully connected, the number of parameters is $6 * 6 * 256 * 4096 = 37,748,736$. The number of operations for one neuron is $6 * 6 * 256$ multiply followed by $6 * 6 * 256$ accumulate/sum. Thus, the total number of operations are $(6 * 6 * 256) * 2 * 4096 = 75,497,472$.

Seventh Layer: The output of the sixth layer goes to the seventh layer i.e., the input to the seventh layer is $[4096]$. The total number of neurons of the seventh layer is 4096. Since it is fully connected, the number of parameters is $4096 * 4096 = 16,777,216$. The total number of operations are $4096 * 2 * 4096 = 33,554,432$.

Eighth Layer: The output of the seventh layer goes to the eighth layer i.e., the input to the eighth layer is $[4096]$. The total number of neurons of the eighth layer is 1000. Since it is fully connected, the number of parameters is $4096 * 1000 = 4,096,000$. The total number of operations are $4096 * 2 * 1000 = 8,192,000$.

Table 1 shows the number of neurons, parameters, and operations for all layers. We notice that 94% of the weights are in the fully connected layers 89% of the computation (and 98% of the neurons) are withing the convolutional layers.

Layer	Neurons	Parameters	Operations
#1	290,400(44.0481%)	34,944 (0.0560%)	115,288,800(7.8527%)
#2	186,624(28.3071%)	614,656 (0.9855%)	537,663,744(36.6223%)
#3	64,896 (9.8435%)	885,120 (1.4191 %)	199,425,408(13.5836%)
#4	64,896 (9.8435%)	1,327,488 (2.1284 %)	299,105,664(20.3732%)
#5	43,264 (6.5623%)	884,992 (1.41895%)	199,403,776(13.5821%)
#6	4096 (0.6212%)	37,748,736 (60.5246%)	75,497,472 (5.1424%)
#7	4096 (0.6212%)	16,777,216 (26.8998%)	33,554,432 (2.2855%)
#8	1000 (0.1516%)	4,096,000 (6.5673 %)	8,192,000 (0.5579%)
Total	659,272	62,369,152	1,468,131,296

Table 1: AlexNet absolute number of neurons, parameters, and operations for the convolutional and fully connected layers (and percentage).