
CS 224N: Assignment 1

RYAN McMAHON

MONDAY 30TH JANUARY, 2017

Contents

1	Problem 1: Softmax (10 pts)	2
1.1	(a) Softmax Invariance to Constant (5 pts)	2
1.2	(b) Softmax Coding (5 pts)	2
2	Problem 2: Neural Network Basics (30 pts)	3
2.1	(a) Sigmoid Gradient (3 pts)	3
2.2	(b) Softmax Gradient w/ Cross Entropy Loss (3 pts)	4
2.3	(c) One Hidden Layer Gradient (6 pts)	6
2.4	(d) No. Parameters (2 pts)	6
2.5	(e) Sigmoid Activation Code (4 pts)	7
2.6	(f) Gradient Check Code (4 pts)	7

Problem 1: Softmax (10 pts)

1.1 (a) Softmax Invariance to Constant (5 pts)

Prove that softmax is invariant to constant offsets in the input, that is, for any input vector x and any constant c , $\text{softmax}(x) = \text{softmax}(x + c)$, where $x + c$ means adding the constant c to every dimension of x . Remember that

$$\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}} \quad (1.1)$$

Answer:

We can show that $\text{softmax}(x) = \text{softmax}(x + c)$ by factoring out c and canceling:

$$\begin{aligned} \text{softmax}(x + c)_i &= \frac{e^{x_i + c}}{\sum_j e^{x_j + c}} = \frac{e^{x_i} \times e^c}{e^c \times \sum_j e^{x_j}} \\ &= \frac{e^{x_i} \times \cancel{e^c}}{\cancel{e^c} \times \sum_j e^{x_j}} = \text{softmax}(x)_i \end{aligned}$$

1.2 (b) Softmax Coding (5 pts)

Given an input matrix of N rows and D columns, compute the softmax prediction for each row using the optimization in part (a). Write your implementation in `q1_softmax.py`. You may test by executing `python q1_softmax.py`.

Note: The provided tests are not exhaustive. Later parts of the assignment will reference this code so it is important to have a correct implementation. Your implementation should also be efficient and vectorized whenever possible (i.e., use numpy matrix operations rather than for loops). A non-vectorized implementation will not receive full credit!

Answer:

See code: `~/code/q1_softmax.py`.

Problem 2: Neural Network Basics (30 pts)

2.1 (a) Sigmoid Gradient (3 pts)

Derive the gradients of the sigmoid function and show that it can be rewritten as a function of the function value (i.e., in some expression where only $\sigma(x)$, but not x , is present). Assume that the input x is a scalar for this question. Recall, the sigmoid function is

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (2.1)$$

Answer:

$$\begin{aligned} \sigma(x) &= \frac{1}{1 + e^{-x}} \\ &= \frac{e^x}{1 + e^x} \\ \frac{\partial}{\partial x} \sigma(x) &= \frac{e^x \times (1 + e^x) - (e^x \times e^x)}{(1 + e^x)^2} \\ &= \frac{e^x + \cancel{(e^x \times e^x)} - \cancel{(e^x \times e^x)}}{(1 + e^x)^2} \\ &= \frac{e^x}{(1 + e^x)^2} = \sigma(x) \times (1 - \sigma(x)) \end{aligned}$$

Because $1 - \sigma(x) = \sigma(-x)$ we can show that:

$$\begin{aligned} \frac{\partial}{\partial x} \sigma(x) &= \frac{e^x}{(1 + e^x)^2} \\ &= \sigma(x) \times \sigma(-x) \\ &= \frac{e^x}{1 + e^x} \times \frac{1}{1 + e^{+x}} \\ &= \frac{e^x}{(1 + e^x)^2} \end{aligned}$$

2.2 (b) Softmax Gradient w/ Cross Entropy Loss (3 pts)

Derive the gradient with regard to the inputs of a softmax function when cross entropy loss is used for evaluation, i.e., find the gradients with respect to the softmax input vector θ , when the prediction is made by $\hat{y} = \text{softmax}(\theta)$. Remember the cross entropy function is

$$CE(y, \hat{y}) = - \sum_i y_i \times \log(\hat{y}_i) \quad (2.2)$$

where y is the one-hot label vector; and \hat{y} is the predicted probability vector for all classes. (Hint: you might want to consider the fact many elements of y are zeros, and assume that only the k -th dimension of y is one.)

Answer:

Let S represent the softmax function:

$$\begin{aligned} f_i &= e^{\theta_i} \\ g_i &= \sum_{k=1}^K e^{\theta_k} \\ S_i &= \frac{f_i}{g_i} \\ \frac{\partial S_i}{\partial \theta_j} &= \frac{f'_i g_i - g'_i f_i}{g_i^2} \end{aligned}$$

So if $i = j$:

$$\begin{aligned} f'_i &= f_i; \quad g'_i = e^{\theta_j} \\ \frac{\partial S_i}{\partial \theta_j} &= \frac{e^{\theta_i} \sum_k e^{\theta_k} - e^{\theta_j} e^{\theta_i}}{(\sum_k e^{\theta_k})^2} \\ &= \frac{e^{\theta_i}}{\sum_k e^{\theta_k}} \times \frac{\sum_k e^{\theta_k} - e^{\theta_j}}{\sum_k e^{\theta_k}} \\ &= S_i \times (1 - S_i) \end{aligned}$$

And if $i \neq j$:

$$\begin{aligned} \frac{\partial S_i}{\partial \theta_j} &= \frac{0 - e^{\theta_j} e^{\theta_i}}{(\sum_k e^{\theta_k})^2} \\ &= - \frac{e^{\theta_j}}{\sum_k e^{\theta_k}} \times \frac{e^{\theta_i}}{\sum_k e^{\theta_k}} \\ &= -S_j \times S_i \end{aligned}$$

We can now use these when operating on our loss function (let L represent the cross entropy function):

$$\begin{aligned}\frac{\partial L}{\partial \theta_i} &= - \sum_k y_k \frac{\partial \log S_k}{\partial \theta_i} \\ &= - \sum_k y_k \frac{1}{S_k} \frac{\partial S_k}{\partial \theta_i} \\ &= -y_i(1 - S_i) - \sum_{k \neq i} y_k \frac{1}{S_k} (-S_k \times S_i) \\ &= -y_i(1 - S_i) + \sum_{k \neq i} y_k S_i \\ &= -y_i + y_i S_i + \sum_{k \neq i} y_k S_i \\ &= S_i \left(\sum_k y_k \right) - y_i\end{aligned}$$

And because we know that $\sum_k y_k = 1$:

$$\frac{\partial L}{\partial \theta_i} = S_i - y_i$$

2.3 (c) One Hidden Layer Gradient (6 pts)

Derive the gradients with respect to the inputs \mathbf{x} to a one-hidden-layer neural network (that is, find $\frac{\partial J}{\partial \mathbf{x}}$ where $J = CE(\mathbf{y}, \hat{\mathbf{y}})$ is the cost function for the neural network). The neural network employs sigmoid activation function for the hidden layer, and softmax for the output layer. Assume the one-hot label vector is \mathbf{y} , and cross entropy cost is used. (Feel free to use $\sigma'(x)$ as the shorthand for sigmoid gradient, and feel free to define any variables whenever you see fit.)

Recall that forward propoagation is as follows

$$\mathbf{h} = \text{sigmoid}(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2)$$

Answer:

Let $f_2 = \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1$ and $f_3 = \mathbf{h}\mathbf{W}_2 + \mathbf{b}_2$;

$$\begin{aligned}\frac{\partial J}{\partial f_3} &= \delta_3 = \hat{\mathbf{y}} - \mathbf{y} \\ \frac{\partial J}{\partial \mathbf{h}} &= \delta_2 = \delta_3 \mathbf{W}_2^T \\ \frac{\partial J}{\partial f_2} &= \delta_1 = \delta_2 \circ \sigma'(f_2) \\ \frac{\partial J}{\partial \mathbf{x}} &= \delta_1 \frac{\partial f_2}{\partial \mathbf{x}} \\ &= \delta_1 \mathbf{W}_1^T\end{aligned}$$

2.4 (d) No. Parameters (2 pts)

How many parameters are there in this neural network [from (c) above], assuming the input is D_x -dimensional, the output is D_y -dimensional, and there are H hidden units?

Answer:

$$\begin{aligned}n_{W_1} &= D_x \times H \\ n_{b_1} &= H \\ n_{W_2} &= H \times D_y \\ n_{b_2} &= D_y \\ N &= (D_x \times H) + H + (H \times D_y) + D_y\end{aligned}$$

2.5 (e) Sigmoid Activation Code (4 pts)

Fill in the implementation for the sigmoid activation function and its gradient in `q2_sigmoid.py`. Test your implementation using `python q2_sigmoid.py`. Again, thoroughly test your code as the provided tests may not be exhaustive.

Answer:

See code: `~/code/q2_sigmoid.py`.

2.6 (f) Gradient Check Code (4 pts)

To make debugging easier, we will now implement a gradient checker. Fill in the implementation for `gradcheck_naive` in `q2_gradcheck.py`. Test your code using `python q2_gradcheck.py`.

Answer:

See code: `~/code/q2_gradcheck.py`.