

Assignment #1

Due date: 4/16 11:59 PM PST

Note: I strongly recommend that you finish these problems before writing the code for the assignment.

These questions require thought, but do not require long answers. Please be as concise as possible. Not all questions will be checked and individual grades will be given.

Setup

Note: Please be sure you have Python 2.7.x installed on your system. The following instructions should work on Mac or Linux. If you have any other system – you are wrong.

Get the code: Download the starter code and the datasets from the course website

Install: To install the required packages locally, run the following:

```
cd assignment1
pip install -r requirements.txt # Install dependencies
```

Start IPython: After you have the Stanford Sentiment data, you should start the IPython notebook server from the `assignment1` directory.

Submitting your work

Once you are done working, put the written part in the same directory as your IPython notebook file, and run the `collectSubmission.sh` script; this will produce a file called `assignment1.zip`.

Submit this file with your written component that should contain your solution.

Tasks

There will be four parts to this assignment, the first three comprise of a written component and a programming component in the IPython notebook. The fourth part is purely programming-based, and we also give you an opportunity to earn extra credits by doing a programming-based optional part. For all of the tasks, you will be using the IPython notebook `wordvec_sentiment.ipynb`.

1 Softmax

Prove that softmax is invariant to constant offsets in the input, that is, for any input vector \mathbf{x} and any constant c ,

$$\text{softmax}(\mathbf{x}) = \text{softmax}(\mathbf{x} + c)$$

where $\mathbf{x} + c$ means adding the constant c to every dimension of \mathbf{x} .

Note: In practice, we make use of this property and choose $c = -\max_i x_i$ when computing softmax probabilities for numerical stability (i.e. subtracting its maximum element from all elements of \mathbf{x}).

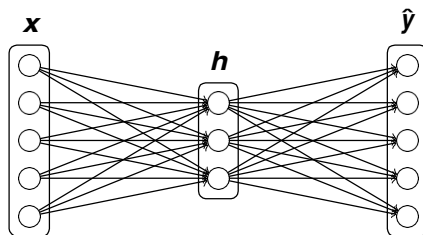
2 Neural Network Basics

- (a) 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.
- (b) 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 $\boldsymbol{\theta}$, when the prediction is made by $\hat{\mathbf{y}} = \text{softmax}(\boldsymbol{\theta})$. Remember the cross entropy function is

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

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

- (c) Derive the gradients with respect to the inputs \mathbf{x} to an one-hidden-layer neural network (that is, find $\frac{\partial J}{\partial \mathbf{x}}$ where J 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 the forward propagation is as follows

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

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

Note that here we're assuming that the input vector (thus the hidden variables and output probabilities) is a row vector to be consistent with the programming assignment. When we apply the sigmoid function to a vector, we are applying it to each of the elements of that vector. \mathbf{W}_i and \mathbf{b}_i ($i = 1, 2$) are the weights and biases, respectively, of the two layers.

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

3 word2vec

- (a) Assume you are given a predicted word vector $\hat{\mathbf{r}}$ (\mathbf{v}_w in the lecture notes in the case of skip-gram), and word prediction is made with the softmax function found in word2vecmodels

$$\Pr(\text{word}_i \mid \hat{\mathbf{r}}, \mathbf{w}) = \frac{\exp(\mathbf{w}_i^\top \hat{\mathbf{r}})}{\sum_{j=1}^{|\mathcal{V}|} \exp(\mathbf{w}_j^\top \hat{\mathbf{r}})}$$

where \mathbf{w}_j ($j = 1, \dots, |\mathcal{V}|$) are the “output” word vectors for all words in the vocabulary (\mathbf{v}_w in the lecture notes).

Assume cross entropy cost is applied to this prediction and word i is the expected word (the i -th element of the one-hot label vector is one), derive the gradients with respect to $\hat{\mathbf{r}}$.

- (b) In the previous problem, derive gradients for the “output” word vectors \mathbf{w}_j 's (including \mathbf{w}_i).
- (c) Repeat part (a) and (b) assuming we are using the negative sampling loss for the predicted vector $\hat{\mathbf{r}}$, and the expected output word is w_i . Assume that K negative samples are drawn, and they are $\mathbf{w}_1, \dots, \mathbf{w}_K$, respectively for simplicity of notation ($i \notin \{1, \dots, K\}$). Recall that the negative sampling loss function in this case is

$$J(\hat{\mathbf{r}}, \mathbf{w}_i, \mathbf{w}_{1 \dots K}) = -\log(\sigma(\mathbf{w}_i^\top \hat{\mathbf{r}})) - \sum_{k=1}^K \log(\sigma(-\mathbf{w}_k^\top \hat{\mathbf{r}}))$$

where $\sigma(\cdot)$ is the sigmoid function.

After you've done this, describe with one sentence why this cost function is much more efficient to compute than the softmax-CE loss (you could provide a speed-up ratio, i.e. the runtime of the softmax-CE loss divided by the runtime of the negative sampling loss).

Note: the cost function here is the negative of what Mikolov et al had in their original paper, because we are doing a minimization instead of maximization in our code.

- (d) Derive gradients for all of the word vectors for skip-gram and CBOW (optional) given the previous parts, given a set of context words $[\text{word}_{i-C}, \dots, \text{word}_{i-1}, \text{word}_i, \text{word}_{i+1}, \dots, \text{word}_{i+C}]$, where C is the context size. You can denote the “input” and “output” word vectors for word_k as \mathbf{v}_{wk} and \mathbf{v}'_{wk} respectively for convenience. (Hint: feel free to use $F(\mathbf{v}_{w_o}|\hat{\mathbf{r}})$ as a placeholder for softmax-CE or negative sampling in this part - you'll see that this is a useful abstraction for the coding part.) Recall that for skip-gram, the cost for a context is

$$J_{\text{skip-gram}}(\text{word}_{i-C \dots i+C}) = \sum_{-C \leq j \leq C, j \neq 0} F(\mathbf{v}'_{w_{i+j}}|\mathbf{v}_{w_i})$$

For (a simpler variant of) CBOW, we sum up the input word vectors in the context

$$\hat{\mathbf{r}} = \sum_{-C \leq j \leq C, j \neq 0} \mathbf{v}_{w_{i+j}}$$

then the CBOW cost is

$$J_{\text{CBOW}}(\text{word}_{i-C \dots i+C}) = F(\mathbf{v}'_{w_i}|\hat{\mathbf{r}})$$

4 Sentiment Analysis

- (a) Explain in less than three sentences why do we want to introduce regularization when doing classification (in fact, most machine learning tasks).
- (b) Plot the classification accuracy on the dev set with respect to the regularization value, using a logarithmic scale on the x-axis. Briefly explain with less than three sentences what you see in the plot.