Advanced Methods in NLP

Home Assignment 1: Word Vectors

Lecturer: Jonathan Berant

Due Date: April 4, 2016

In this home assignment we will implement the skip-gram model with negative sampling.

Please copy all the data and supporting code from ~dormuhlg/advanced_nlp/assignment1. To setup the environment, follow the instructions in http://web.stanford.edu/class/cs224n/assignment1/ (we adapted this home assignment from that class).

To submit your solution, create a directory ~dormuhlg/advanced_nlp/assignment1/<id1> _<id2> (where id1 refers to the ID of the first student) and put all relevant files in this directory. The submission directory should include the code and data necessary for running the tests provided out-of-the-box, as well as a written solution, and a text file including an e-mail of one of the students.

1 Basics

(a) Prove that softmax is invariant to constant offset in the input, that is, for any input vector \boldsymbol{x} and any constant c,

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

where x + c means adding the constant c to every dimension of x. Remember that

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

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 x).

(b) 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 softmax.py. You may test by executing 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!

(c) 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:

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

- (d) Implement the sigmoid function in sigmoid.py and test your code by running python sigmoid.py.
- (e) To make debugging easier, we will now implement a gradient checker. Fill in the implementation for gradcheck_naive in gradcheck.py. Test your code using python gradcheck.py.

2 Word2vec

(a) Assume you are given a predicted word vector $\mathbf{v_c}$ corresponding to the center word c for skipgram, and the word prediction is made with the **softmax** function

$$\hat{\boldsymbol{y}}_{\boldsymbol{o}} = p(o|c) = \frac{\exp(\boldsymbol{u}_{\boldsymbol{o}}^{\mathsf{T}} \boldsymbol{v}_{\boldsymbol{c}})}{\sum_{w=1}^{W} \exp(\boldsymbol{u}_{\boldsymbol{w}}^{\mathsf{T}} \boldsymbol{v}_{\boldsymbol{c}})}$$

where o is the expected word, w denotes the w-th word and u_w (w = 1, ..., W) are the "output" word vectors for all words in the vocabulary. Let's define the cross entropy function as:

$$J_{\text{softmax-CE}}(o, v_c, U) = CE(\boldsymbol{y}, \hat{\boldsymbol{y}}) = -\sum_i y_i \cdot \log(\hat{y}_i)$$

where the gold vector \boldsymbol{y} is a one-hot vector, the softmax prediction vector $\hat{\boldsymbol{y}}$ is a probability distribution over the output space, and $\boldsymbol{U} = [u_1, u_2, ..., u_W]$ is the matrix of all the output vectors. Assume cross entropy cost is applied to this prediction, derive the gradients with respect to $\boldsymbol{v_c}$.

- (b) As in the previous part, derive gradients for the "output" word vectors u_w (including u_o).
- (c) Repeat part (a) and (b) assuming we are using the negative sampling loss for the predicted vector v_c . Assume that K negative samples (words) are drawn, and they are $1, \ldots, K$, respectively. For simplicity of notation, assume ($o \notin \{1, \ldots, K\}$). Again, for a given word, o, denote its output vector as u_o . The negative sampling loss function in this case is:

$$\mathbf{J}_{\text{neg-sample}}(o, v_c, U) = -\log(\sigma(\boldsymbol{u_o^\intercal} \boldsymbol{v_c})) - \sum_{k=1}^K \log(\sigma(-\boldsymbol{u_k^\intercal} \boldsymbol{v_c}))$$

where σ is the sigmoid function defined in 1(c).

(d) Derive gradients for all of the word vectors for skip-gram given the previous parts and given a set of context words $[\text{word}_{c-m}, \dots, \text{word}_c, \dots, \text{word}_{c+m}]$ where m is the context size. Denote the "input" and "output" word vectors for $word_k$ as v_k and u_k respectively.

Hint: feel free to use $F(o, \mathbf{v_c})$ (where o is the expected word) as a placeholder for the $J_{softmax-CE}(o, v_c...)$ or $J_{neg-sample}(o, v_c...)$ cost functions in this part – you'll see

that this is a useful abstraction for the coding part. That is, your solution may contain terms of the form $\frac{\partial F(o, \mathbf{v_c})}{\partial \dots}$ Recall that for skip-gram, the cost for a context centered around c is:

$$\sum_{-m \le j \le m, j \ne 0} F(w_{c+j}, v_c)$$

- (e) In this part you will implement the word2vec models and train your own word vectors with stochastic gradient descent (SGD). First, write a helper function to normalize rows of a matrix in word2vec.py. In the same file, fill in the implementation for the softmax and negative sampling cost and gradient functions. Then, fill in the implementation of the cost and gradient functions for the skip-gram model. When you are done, test your implementation by running python word2vec.py.
- (f) Complete the implementation for your SGD optimizer in sgd.py. Test your implementation by running python sgd.py.
- (g) In this part you will implement the k-nearest neighbors algorithm, which will be used for analysis. The algorithm receives a vector, a matrix and an integer k, and returns k indices of the matrix's rows that are closest to the vector. Use the cosine similarity as a distance metric (https://en.wikipedia.org/wiki/Cosine_similarity). Fill the implementation of the algorithm in knn.py.

Note: for this part a naive implementation will be accepted.

(h) Show time! Now we are going to load some real data and train word vectors with everything you just implemented! We are going to use the Stanford Sentiment Treebank (SST) dataset to train word vectors. There is no additional code to write for this part; just run python run.py.

Note: The training process may take a long time depending on the efficiency of your implementation. Plan accordingly!

When the script finishes, a visualization for your word vectors will appear. It will also be saved as word_vectors.png in your project directory. In addition, the script should print the nearest neighbors of a few words (using the knn function you implemented in 2(g)). Include the plot and the nearest neighbors lists in your homework write up, and briefly explain those results.

(i) (optional) Download some pre-trained word vectors such as GloVe (https://nlp.stanford.edu/projects/glove/) and compute the nearest neighbors of the same words with these vectors. This has been trained on much more data. Does it look better? Do not submit pre-trained word vectors as part of the submission.