## Submission Assignment #2 (word2vec (43 Points))

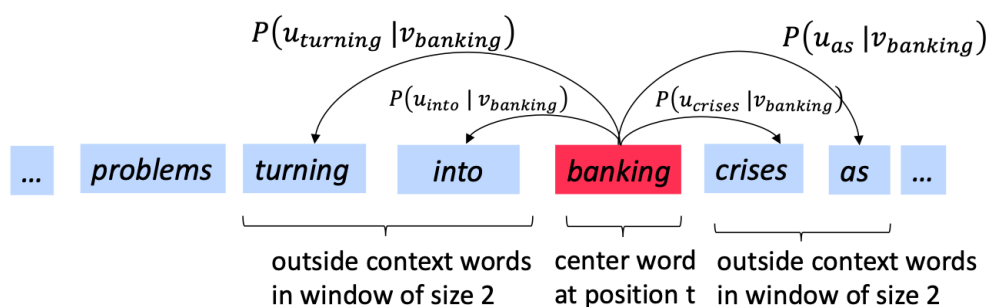*Instructor:* Christopher Manning                                                                 *Name:* Jianpan Gun

---

**Problem 1: Understanding word2vec**                                      (3+5+5+3+4+3=23 points)

Let's have a quick refresher on the `word2vec`algorithm. The key insight behind `word2vec`is that 'a word is known by the company it keeps'. Concretely, suppose we have a 'center' word c and a contextual window surrounding c. We shall refer to words that lie in this contextual window as 'outside words'. For example, in Figure 1 we see that the center word c is 'banking'. Since the context window size is 2, the outside words are 'turning', 'into', 'crises', and 'as'.

The goal of the skip-gram `word2vec`algorithm is to accurately learn the probability distribution $P(O|C)$. Given a specific word o and a specific word c, we want to calculate $P(O = o|C = c)$, which is the probability that word o is an 'outside' word for c, i.e., the probability that o falls within the contextual window of c.



Figure 1: The `word2vec`skip-gram prediction model with window size 2.

In `word2vec`, the conditional probability distribution is given by taking vector dot-products and applying the softmax function:

$$P(O = o|C = c) = \frac{\exp\left(\boldsymbol{u}_o^\top \boldsymbol{v}_c\right)}{\sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right)} \tag{1}$$

Here, $\boldsymbol{u_o}$ is the 'outside' vector representing outside word o, and $\boldsymbol{v_c}$ is the 'center' vector representing center word c. To contain these parameters, we have two matrices, $\boldsymbol{U}$ and $\boldsymbol{V}$. The columns of $\boldsymbol{U}$ are all the 'outside' vectors $\boldsymbol{u_w}$. The columns of $\boldsymbol{V}$ are all of the 'center' vectors $\boldsymbol{v_w}$. Both $\boldsymbol{U}$ and $\boldsymbol{V}$ contain a vector for every $w \in$ Vocabulary.

Recall from lectures that, for a single pair of words c and o, the loss is given by:

$$J_{\text{naive-softmax}}\left(\boldsymbol{v}_c, o, \boldsymbol{U}\right) = -\log P(O = o|C = c) \tag{2}$$

Another way to view this loss is as the cross-entropy between the true distribution $\boldsymbol{y}$ and the predicted distribution $\hat{\boldsymbol{y}}$. Here, both $\boldsymbol{y}$ and $\hat{\boldsymbol{y}}$ are vectors with length equal to the number of words in the vocabulary. Furthermore, the $k^{th}$ entry in these vectors indicates the conditional probability of the $k^{th}$ word being an 'outside word' for the given c. The true empirical distribution $\boldsymbol{y}$ is a one-hot vector with a 1 for the true outside word o, and 0 everywhere else. The predicted distribution $\hat{\boldsymbol{y}}$ is the probability distribution $P(O|C = c)$ given by our model in equation (1).

**(a)** (3 points) **Show that the naive-softmax loss given in Equation (2) is the same as the cross-entropy loss between $y$ and $\hat{y}$; i.e., show that** $-\sum_{w \in V_{ocab}} y_w \log\left(\hat{y}_w\right) = -\log\left(\hat{y}_o\right)$

**A**: Due to the $\boldsymbol{y_i}$ is a one-hot vector, that is, $\mathbb{I}_i$, which means that only have i-th dim is 1. So,

$$\mathcal{L}_{ce} = -\sum_{w \in V_{ocab}} y_w \log(\hat{y}_w) = -\boldsymbol{y}_o \log(\hat{\boldsymbol{y}}_o) = -\log(\hat{y}_o) \tag{3}$$

**(b)** (5 points) **Compute the partial derivative of $J_{\text{naive-softmax}}(\boldsymbol{v}_c, o, \boldsymbol{U})$ with respect to $\boldsymbol{v}_c$. Please write your answer in terms of $\boldsymbol{y}$, $\hat{\boldsymbol{y}}$, and $\boldsymbol{U}$**

$$
\begin{aligned}
\frac{\partial J_{\text{naive-softmax}}}{\partial \boldsymbol{v}_c}(\boldsymbol{v}_c, o, \boldsymbol{U}) &= \frac{\partial}{\partial \boldsymbol{v}_c} - \log \frac{\exp\left(\boldsymbol{u}_o^\top \boldsymbol{v}_c\right)}{\sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right)} = -\boldsymbol{u}_o + \frac{\partial}{\partial \boldsymbol{v}_c} \log \sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right) \\
&= -\boldsymbol{u}_o + \frac{1}{\sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right)} \sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right) \boldsymbol{u}_w \\
&= -\boldsymbol{y}^\top \boldsymbol{U} + \sum_{w \in V_{\text{ocab}}} P(O = w | C = c) \boldsymbol{u}_w \\
&= -\boldsymbol{y}^\top \boldsymbol{U} + \hat{\boldsymbol{y}}^\top \boldsymbol{U} = (\hat{\boldsymbol{y}} - \boldsymbol{y})^\top \boldsymbol{U}
\end{aligned}
\tag{4}
$$

**(c)** (5 points) **Compute the partial derivatives of $J_{\text{naive-softmax}}(\boldsymbol{v}_c, o, \boldsymbol{U})$ with respect to each of the 'outside' word vectors, $\boldsymbol{u}_w$'s. There will be two cases: when $w = o$, the true 'outside' word vector; and $w \neq o$, for all other words. Please write you answer in terms of $\boldsymbol{y}$, $\hat{\boldsymbol{y}}$, and $\boldsymbol{v}_c$.**

$$\frac{\partial J_{\text{naive-softmax}}}{\partial \boldsymbol{u}_w}(\boldsymbol{v}_c, o, \boldsymbol{U}) = \frac{\partial}{\partial \boldsymbol{u}_w} - \log \frac{\exp\left(\boldsymbol{u}_o^\top \boldsymbol{v}_c\right)}{\sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right)} = -\frac{\partial \boldsymbol{u}_o^\top \boldsymbol{v}_c}{\partial \boldsymbol{u}_w} + \frac{\partial}{\partial \boldsymbol{u}_w} \log \sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right) \tag{5}$$

If $w = o$:

$$= -\boldsymbol{v}_c + \frac{\exp\left(\boldsymbol{u}_o^\top \boldsymbol{v}_c\right) \boldsymbol{v}_c}{\sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right)} = (P(O = o | C = c) - 1)\boldsymbol{v}_c \tag{6}$$

If $w \neq o$:

$$= \frac{\exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right) \boldsymbol{v}_c}{\sum_{w \in V_{\text{ocab}}} \exp\left(\boldsymbol{u}_w^\top \boldsymbol{v}_c\right)} = P(O = w | C = c)\boldsymbol{v}_c \tag{7}$$

So,

$$= (\hat{\boldsymbol{y}} - \boldsymbol{y})^\top \boldsymbol{v}_c \tag{8}$$

**(d)** (3 Points) **The sigmoid function is given by Equation 9:**

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \tag{9}$$

Please compute the derivative of $\sigma(x)$ with respect to x, where x is a scalar. **Hint: you may want to write your answer in terms of $\sigma(x)$.**

$$\frac{\partial \sigma(x)}{\partial x} = \frac{\partial}{\partial x} \frac{e^x}{e^x + 1} = \frac{e^x(e^x + 1) - e^x(e^x + 1)'}{(e^x + 1)^2} = \frac{e^{2x}}{(e^x + 1)^2} = \sigma(x)(1 - \sigma(x)) \tag{10}$$

**(e)** (4 points) **Now we shall consider the Negative Sampling loss, which is an alternative to the Naive Softmax loss. Assume that K negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as $w_1, w_2, \ldots, w_K$ and their outside vectors as $u_1, \ldots, u_K$. Note that $o \notin w_1, \ldots, w_K$. For a center word c and an outside word $o$, the negative sampling loss function is given by:**

$$\boldsymbol{J}_{\text{neg-sample}}(\boldsymbol{v}_c, o, \boldsymbol{U}) = -\log\left(\sigma\left(\boldsymbol{u}_o^\top \boldsymbol{v}_c\right)\right) - \sum_{k=1}^K \log\left(\sigma\left(-\boldsymbol{u}_k^\top \boldsymbol{v}_c\right)\right) \tag{11}$$

**for a sample $w_1, w_2, \ldots, w_K$, where $\sigma(\cdot)$ is the sigmoid function.**

Please repeat parts (b) and (c), computing the partial derivatives of $J_{\text{neg-sample}}$ with respect to $v_c$, with respect to $u_o$, and with respect to a negative sample $u_k$. Please write your answers in terms of the vectors $u_o, v_c$, and $u_k$, where $k \in [1, k]$. After you've done this, describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss. Note, you should be able to use your solution to part (d) to help compute the necessary gradients here.

$$\frac{\partial J_{\text{neg-sample}}}{\partial v_c}(v_c, o, U) = \frac{\partial}{\partial v_c}[-\log\left(\sigma\left(u_o^\top v_c\right)\right) - \sum_{k=1}^{K}\log\left(\sigma\left(-u_k^\top v_c\right)\right)] = -\frac{\sigma\left(u_o^\top v_c\right)\left(1 - \sigma\left(u_o^\top v_c\right)\right)u_o}{\sigma\left(u_o^\top v_c\right)}$$
$$+ \sum_{k=1}^{K}\left(1 - \sigma\left(-u_k^\top v_c\right)\right)u_k = -\left(1 - \sigma\left(u_o^\top v_c\right)\right)u_o + \sum_{k=1}^{K}\left(1 - \sigma\left(-u_k^\top v_c\right)\right)u_k \tag{12}$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_o}(v_c, o, U) = \frac{\partial}{\partial u_o}[-\log\left(\sigma\left(u_o^\top v_c\right)\right) - \sum_{k=1}^{K}\log\left(\sigma\left(-u_k^\top v_c\right)\right)] = -\frac{\sigma\left(u_o^\top v_c\right)\left(1 - \sigma\left(u_o^\top v_c\right)\right)v_c}{\sigma\left(u_o^\top v_c\right)}$$
$$= -\left(1 - \sigma\left(u_o^\top v_c\right)\right)v_c \tag{13}$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial u_k}(v_c, o, U) = \frac{\partial}{\partial u_k}[-\log\left(\sigma\left(u_o^\top v_c\right)\right) - \sum_{k=1}^{K}\log\left(\sigma\left(-u_k^\top v_c\right)\right)] = \left(1 - \sigma\left(-u_k^\top v_c\right)\right)v_c \tag{14}$$

**(f)** (3 points) **Suppose the center word is $c = wt$ and the context window is $w_{t-m}, \ldots, w_{t-1}, w_t, w_{t+1}, \ldots, w_{t+m}$, where m is the context window size. Recall that for the skip-gram version of `word2vec`, the total loss for the context window is:**

$$J_{\text{skip-gram}}(v_c, w_{t-m}, \ldots w_{t+m}, U) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(v_c, w_{t+j}, U) \tag{15}$$

Here, $J(v_c, w_{t+j}, U)$ represents an arbitrary loss term for the center word $c = wt$ and outside word $w_{t+j}$. $J(v_c, w_{t+j}, U)$ could be $J_{\text{neg-sample}}(v_c, w_{t+j}, U)$ or $J_{\text{naive-softmax}}(v_c, w_{t+j}, U)$, depending on your implementation.

Write down three partial derivatives:

(I) $\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \ldots w_{t+m}, U)/\partial U$

(II) $\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \ldots w_{t+m}, U)/\partial v_c$

(III) $\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \ldots w_{t+m}, U)/\partial v_w$, when $w \neq c$

Write your answers in terms of $\partial J(v_c, w_{t+j}, U)/\partial U$ and $\partial J(v_c, w_{t+j}, U)/\partial v_c$. This is very simple. Each solution should be one line. Once you're done: Given that you computed the derivatives of $J(v_c, w_{t+j}, U)$ with respect to all the model parameters $U$ and $V$ in parts (a) to (c), you have now computed the derivatives of the full loss function $J_{skip-gram}$ with respect to all parameters. You're ready to implement `word2vec`!

(I)
For native-softmax:
$$\frac{\partial J_{\text{skip-gram}}}{\partial U}(v_c, o, U) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} (\hat{y} - y)^\top v_j \tag{16}$$

For neg-sample:
$$\frac{\partial J_{\text{skip-gram}}}{\partial U}(v_c, o, U) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} (1 - \sigma\left(u_o^\top v_j\right))v_j \tag{17}$$

(II)
For native-softmax:
$$\frac{\partial J_{\text{skip-gram}}}{\partial v_c}(v_c, o, U) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} (\hat{y} - y)^\top U \tag{18}$$

For neg-sample:

$$\frac{\partial J_{\text{skip-gram}}}{\partial \boldsymbol{v}_c}\left(\boldsymbol{v}_c, o, \boldsymbol{U}\right) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \left(-(1 - \sigma\left(\boldsymbol{u}_o^\top \boldsymbol{v}_j\right))\boldsymbol{u}_o + \sum_{k=1}^{K}\left(1 - \sigma\left(-\boldsymbol{u}_k^\top \boldsymbol{v}_j\right)\right)\boldsymbol{u}_k\right) \tag{19}$$

(III)

$$\frac{\partial J_{\text{skip-gram}}}{\partial \boldsymbol{v}_w}\left(\boldsymbol{v}_c, o, \boldsymbol{U}\right) = 0 \tag{20}$$

| **Problem Implementing word2vec (20 points)** | (12+4+4=20 points) |

**(a)** (12 points) **We will start by implementing methods in** `word2vec.py`. **First, implement the** `sigmoid` **method, which takes in a vector and applies the sigmoid function to it. Then implement the softmax loss and gradient in the** `naiveSoftmaxLossAndGradient` **method, and negative sampling loss and gradient in the** `negSamplingLossAndGradient` **method. Finally, fill in the implementation for the skip-gram model in the** `skipgram` **method. When you are done, test your implementation by running python** `word2vec.py`

Listing 1: word2vec.py function.

```python
def sigmoid(x):
    # \sigma(x)=\frac{1}{1+e^{-x}}=\frac{e^{x}}{e^{x}+1}
    s = 1 / (1 + np.exp(-x))
    return s

def naiveSoftmaxLossAndGradient(centerWordVec, outsideWordIdx, outsideVectors, ↩
    dataset):
    y_hat = softmax(np.dot(centerWordVec, outsideVectors.T))

    loss = -np.log(y_hat[outsideWordIdx]) # ref to equation(2)
    y = np.zeros_like(y_hat)
    y[outsideWordIdx] = 1
    gradCenterVec = np.dot((y_hat - y).T, outsideVectors) # ref to equation(4)
    gradOutsideVecs = np.dot(
        (y_hat - y).reshape(-1, 1), centerWordVec.reshape(1, -1)
    )  # (T, 1) * (1, H) ref to equation(8)

    return loss, gradCenterVec, gradOutsideVecs

def negSamplingLossAndGradient(
    centerWordVec, outsideWordIdx, outsideVectors, dataset, K=10
):
    negSampleWordIndices = getNegativeSamples(outsideWordIdx, dataset, K)
    indices = [outsideWordIdx] + negSampleWordIndices

    negSampleWordIndices = np.array(negSampleWordIndices)

    pos = sigmoid(np.dot(outsideVectors[outsideWordIdx], centerWordVec))
    neg = sigmoid(-np.dot(outsideVectors[negSampleWordIndices], centerWordVec))
    loss = -np.log(pos) - np.log(neg).sum() # ref to equation(11)
    gradCenterVec = -np.dot((1 - pos), outsideVectors[outsideWordIdx]) + np.dot(
        1 - neg, outsideVectors[negSampleWordIndices]
    ) # ref to equation(12)
    gradOutsideVecs = np.zeros_like(outsideVectors)
    gradOutsideVecs[outsideWordIdx] = - np.dot(1 - pos, centerWordVec) # ref to ↩
        equation(13)
    gradOutside = np.dot((1 - neg).reshape(-1, 1), centerWordVec.reshape(1, -1))
    for idx, neg_idx in enumerate(negSampleWordIndices):
        gradOutsideVecs[neg_idx] += gradOutside[idx] # ref to equation(14)
```

```
38
39      return loss, gradCenterVec, gradOutsideVecs
40
41  def skipgram(
42      currentCenterWord,
43      windowSize,
44      outsideWords,
45      word2Ind,
46      centerWordVectors,
47      outsideVectors,
48      dataset,
49      word2vecLossAndGradient=naiveSoftmaxLossAndGradient,
50  ):
51      loss = 0.0
52      gradCenterVecs = np.zeros(centerWordVectors.shape)
53      gradOutsideVectors = np.zeros(outsideVectors.shape)
54
55      centerWordIdx = word2Ind[currentCenterWord]
56      centerWordVector = centerWordVectors[centerWordIdx]
57      for outsideWord in outsideWords:
58          outsideWordIdx = word2Ind[outsideWord]
59          tmp_loss, gradCenterVec, gradOutsideVector = word2vecLossAndGradient(
60              centerWordVector, outsideWordIdx, outsideVectors, dataset
61          )
62          loss += tmp_loss
63          gradCenterVecs[centerWordIdx, :] += gradCenterVec # ref to equation(18,19)
64          gradOutsideVectors += gradOutsideVector
65
66      return loss, gradCenterVecs, gradOutsideVectors
```

**(b)** (4 points) **Complete the implementation for your** `SGD` **optimizer in the** `sgd` **method of** `sgd.py`. **Test your implementation by running** `python sgd.py`.

Listing 2: sgd.py function.

```
1
2   def sgd(f, x0, step, iterations, postprocessing=None, useSaved=False, PRINT_EVERY↩
        =10):
3       # Anneal learning rate every several iterations
4       ANNEAL_EVERY = 20000
5
6       if useSaved:
7           start_iter, oldx, state = load_saved_params()
8           if start_iter > 0:
9               x0 = oldx
10              step *= 0.5 ** (start_iter / ANNEAL_EVERY)
11
12          if state:
13              random.setstate(state)
14      else:
15          start_iter = 0
16
17      x = x0
18
19      if not postprocessing:
20          postprocessing = lambda x: x
21
22      exploss = None
23
24      for iter in range(start_iter + 1, iterations + 1):
25          # You might want to print the progress every few iterations.
26
```

```
27            # c, gin, gout = word2vecModel(
28            #     centerWord,
29            #     windowSize1,
30            #     context,
31            #     word2Ind,
32            #     centerWordVectors,
33            #     outsideVectors,
34            #     dataset,
35            #     word2vecLossAndGradient,
36            # )
37            # loss += c / batchsize
38            # grad[: int(N / 2), :] += gin / batchsize
39            # grad[int(N / 2) :, :] += gout / batchsi
40            loss, grad = f(x)
41            x -= step * grad
42
43            x = postprocessing(x)
44            if iter % PRINT_EVERY == 0:
45                if not exploss:
46                    exploss = loss
47                else:
48                    exploss = 0.95 * exploss + 0.05 * loss
49                print("iter %d: %f" % (iter, exploss))
50
51            if iter % SAVE_PARAMS_EVERY == 0 and useSaved:
52                save_params(iter, x)
53
54            if iter % ANNEAL_EVERY == 0:
55                step *= 0.5
56        return x
```

**(c)** (4 points) **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, and later apply them to a simple sentiment analysis task. You will need to fetch the datasets first. To do this, run sh get** `datasets.sh`. **There is no additional code to write for this part; just run python** `run.py`.
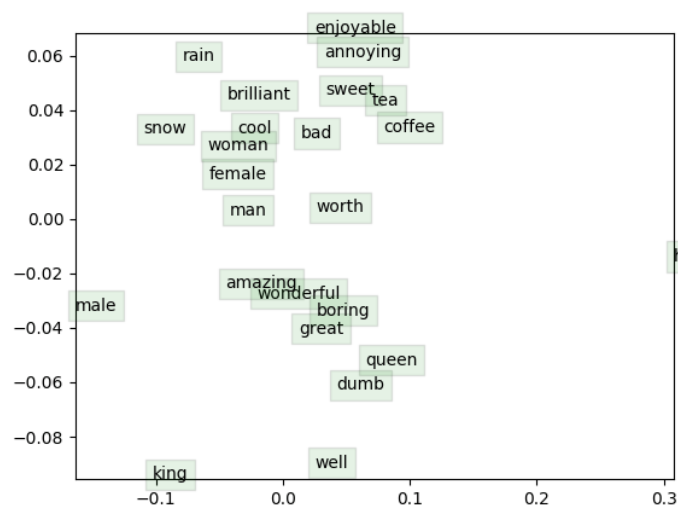


Figure 2: The word vector distribution figure.