

CS 224n Assignment #2: word2vec (44 Points)

Carolyn Kao (#06366163) chkao83@stanford.edu

Jan 26, 2021

1 Written: Understanding word2vec (26 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 .

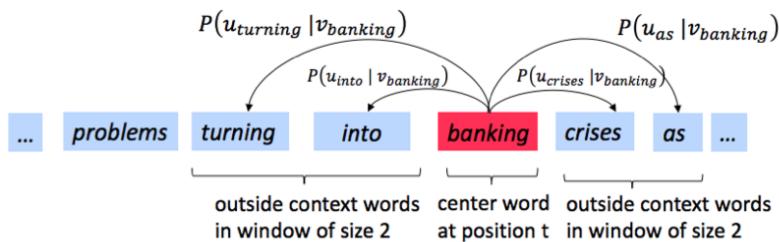


Figure 1: The word2vec skip-gram prediction model with window size 2

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

$$P(O = o | C = c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (1)$$

Here, \mathbf{u}_o is the 'outside' vector representing outside word o , and \mathbf{v}_c is the 'center' vector representing center word c . To contain these parameters, we have two matrices, \mathbf{U} and \mathbf{V} . The columns of \mathbf{U} are all the 'outside' vectors \mathbf{u}_w . The columns of \mathbf{V} are all of the 'center' vectors \mathbf{v}_w . Both \mathbf{U} and \mathbf{V} contain a vector for every $w \in \text{Vocabulary}$.¹

U_K : K^{th} column of \mathbf{U} ; the outside vector for the word indexed by K

V_K : K^{th} column of \mathbf{V} ; the center vector for the word indexed by K

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

$$J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c). \quad (2)$$

We can view this loss as the cross-entropy² between the true distribution \mathbf{y} and the predicted distribution $\hat{\mathbf{y}}$. Here, both \mathbf{y} and $\hat{\mathbf{y}}$ 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 \mathbf{y} is a one-hot vector with a 1 for the true outside word o , and 0 everywhere else. The predicted distribution $\hat{\mathbf{y}}$ 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 \mathbf{y} and $\hat{\mathbf{y}}$; i.e., show that

$$-\sum_{w \in Vocab} y_w \log(\hat{y}_w) = -\log(\hat{y}_o). \quad (3)$$

Your answer should be one line.

$\because \mathbf{y}$ is a one-hot vector in which $\begin{cases} y_w = 1 & \text{when } w=o \\ y_w = 0 & \text{when } w \neq o \end{cases}$

$$\begin{aligned} \therefore -\sum_{w \in Vocab} y_w \log(\hat{y}_w) &= -[y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + \dots + y_o \log \hat{y}_o + \dots + y_w \log \hat{y}_w] \\ &= -y_o \log \hat{y}_o \\ &= -\log \hat{y}_o \quad \# \end{aligned}$$

- (b) (5 points) Compute the partial derivative of $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{v}_c . Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{U} . Note that in this course, we expect your final answers to follow the shape convention.³ This means that the partial derivative of any function $f(x)$ with respect to x should have the same shape as x . For this subpart, please present your answer in vectorized form. In particular, you may not refer to specific elements of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{U} in your final answer (such as y_1 , y_2 , ...).

Method I : Presented in Matrix-Vector Multiplication form

From (a), we see that the naive-softmax loss function can be expressed as cross-entropy (CE) loss btw \mathbf{y} and $\hat{\mathbf{y}}$.

Hence denote $J = \text{CE}(\mathbf{y}, \hat{\mathbf{y}})$

Then, denote the input vector as $\theta = \mathbf{U}^T \mathbf{v}_c$ and the predicted function $\hat{\mathbf{y}}$ as function of θ .

$$\frac{\partial J}{\partial \theta} = (\hat{\mathbf{y}} - \mathbf{y})^T$$

(ref: <http://courses.cs.ut.ee/MTAT.03.277/2015-fall/uploads/Main/word2vec.pdf>)

$$\begin{aligned} \text{Using chain rule, } \frac{\partial J}{\partial \mathbf{v}_c} &= \frac{\partial J}{\partial \theta} \frac{\partial \theta}{\partial \mathbf{v}_c} \\ &= (\hat{\mathbf{y}} - \mathbf{y})^T \frac{\partial \mathbf{U}^T \mathbf{v}_c}{\partial \mathbf{v}_c} \\ &= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}) \# \end{aligned}$$

(Akshay's OH on Nook)

Method II : Presented in Vector-Vector Addition Form

$$\begin{aligned} J(\mathbf{v}_c, o, \mathbf{U}) &= -\log(\hat{y}_o) \\ &= -\log \left(\frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \right) \\ &= -[\log(\exp(\mathbf{u}_o^T \mathbf{v}_c)) - \log(\sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c))] \\ &= -\mathbf{u}_o^T \mathbf{v}_c + \log(\sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)) \end{aligned}$$

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \frac{1}{\sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)$$

$$= -u_0 + \frac{1}{\sum_{w'} \exp(u_{w'}^T v_c)} \sum_w \frac{\partial}{\partial v_c} \exp(u_w^T v_c)$$

$$= -u_0 + \frac{1}{\sum_{w'} \exp(u_w^T v_c)} \sum_w \exp(u_w^T v_c) u_w$$

$$= -u_0 + \sum_w \frac{\exp(u_w^T v_c)}{\sum_{w'} \exp(u_w^T v_c)} u_w$$

$$= \underbrace{\sum_w \hat{y}_w u_w}_{\text{It is the observed representation of the outside word}} - \underbrace{u_0}_{\#}$$

It is
the
expected
context
word

It is the
observed
representation
of the outside
word

(c) (5 points) Compute the partial derivatives of $J_{\text{naive-softmax}}(v_c, o, U)$ with respect to each of the ‘outside’ word vectors, 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 your answer in terms of y , \hat{y} , and v_c . In this subpart, you may use specific elements within these terms as well, such as (y_1, y_2, \dots) .

(Akshay's DM on Nook)

$$J(v_c, o, U) = -\log(\hat{y}_o)$$

$$= -\log \left(\frac{\exp(u_o^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \right)$$

$$= -[\log(\exp(u_o^T v_c)) - \log(\sum_w \exp(u_w^T v_c))]$$

$$= -u_o^T v_c + \underbrace{\log(\sum_w \exp(u_w^T v_c))}_{\text{take derivative of this}}$$

↓

$$\frac{\partial}{\partial u_w} \log(\sum_w \exp(u_w^T v_c))$$

$$= \frac{1}{\sum_w \exp(u_w^T v_c)} \frac{\partial}{\partial u_w} \sum_w \exp(u_w^T v_c)$$

$$= \frac{1}{\sum_w \exp(u_w^T v_c)} \frac{\partial}{\partial u_w} \exp(u_w^T v_c)$$

$$= \underbrace{\left[\frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} \right]}_{P(W=w | C=c)} v_c = \hat{y}_w v_c$$

$$P(W=w | C=c)$$

$$\text{If } w=o, \quad \frac{\partial J}{\partial v_w} = [P(W=w | C=c) - 1] v_c$$

$$\text{If } w \neq o, \quad \frac{\partial J}{\partial v_w} = P(W=w | C=c) v_c$$

$$\text{Hence, } \frac{\partial J_{\text{naive-softmax}}(v_c, o, U)}{\partial v_w} = (\hat{y} - y)^T v_c = \hat{y}_w v_c - y_w v_c$$

- (d) (1 point) Compute the partial derivative of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{U} . Please write your answer in terms of $\frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial u_1}, \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial u_2}, \dots, \frac{\partial J(\mathbf{v}_c, o, \mathbf{U})}{\partial u_{|\text{Vocab}|}}$. The solution should be one or two lines long.

(Rui Wang's OH on Nook)

$$\mathbf{U} = \left[\begin{array}{c} \overbrace{\quad \quad \quad}^N \end{array} \right] \quad \left. \right\} \text{dim} = M \quad \text{by shape convention,} \\ \frac{\partial J}{\partial \mathbf{U}} \quad \text{is } M \times N$$

$$\frac{\partial J}{\partial \mathbf{U}} = \left[\frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_{|\text{Vocab}|}} \right] \quad \#$$

$\left(\begin{array}{l} \text{This is a row matrix formed by} \\ \text{column vec's, each w/ dim M} \end{array} \right)$

(e) (3 Points) The sigmoid function is given by Equation 4:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (4)$$

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 $\sigma(x)$.

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

- (f) (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, \dots, w_K and their outside vectors as $\mathbf{u}_1, \dots, \mathbf{u}_K$. For this question, assume that the K negative samples are distinct. In other words, $i \neq j$ implies $w_i \neq w_j$ for $i, j \in \{1, \dots, K\}$. Note that $o \notin \{w_1, \dots, w_K\}$. For a center word c and an outside word o , the negative sampling loss function is given by:

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \quad (5)$$

for a sample w_1, \dots, w_K , where $\sigma(\cdot)$ is the sigmoid function.⁴

Please repeat parts (b) and (c), computing the partial derivatives of $J_{\text{neg-sample}}$ with respect to \mathbf{v}_c , with respect to \mathbf{u}_o , and with respect to a negative sample \mathbf{u}_k . Please write your answers in terms of the vectors \mathbf{u}_o , \mathbf{v}_c , and \mathbf{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 (e) to help compute the necessary gradients here.

$$J = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

with respect to \mathbf{v}_c

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{v}_c} &= \frac{\partial}{\partial \mathbf{v}_c} \left[-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\ &= \underbrace{\frac{\partial}{\partial \mathbf{v}_c} [-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))]}_{=} - \underbrace{\frac{\partial}{\partial \mathbf{v}_c} \left[\sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right]}_{=} \\ &= -\mathbf{u}_o(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \left[-\sum_{k=1}^K \mathbf{u}_k(1 - \sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\ &= (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{u}_o - \sum_{k=1}^K (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1) \mathbf{u}_k \end{aligned}$$

with respect to \mathbf{u}_o

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_o} &= \frac{\partial}{\partial \mathbf{u}_o} \left[-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \right] \\ &= \frac{\partial}{\partial \mathbf{u}_o} [-\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c))] - \mathbf{v}_c \\ &= -\left\{ \frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} [\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \mathbf{v}_c] \right\} \\ &= (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{v}_c \end{aligned}$$

with respect to u_k

\because All words are distinct

$\therefore u_1, u_2 \dots$ each corresponds to a specific word —

$$\text{Do on } - \sum_{i \in \{1, \dots, K \} \text{ where } w_i = w_k} \log(\sigma(-u_i^T v_c))$$

$$\Rightarrow \frac{\partial J}{\partial u_k} = - \frac{\partial}{\partial u_k} \left[\sum_{k=1}^K \log(\sigma(-u_k^T v_c)) \right]$$

$$= - \frac{1}{\sigma(-u_k^T v_c)} \cdot \frac{\partial}{\partial u_k} (\sigma(-u_k^T v_c))$$

$$= \frac{\sigma(-u_k^T v_c)(1 - \sigma(-u_k^T v_c))v_c}{\sigma(-u_k^T v_c)}$$

$$= (1 - \sigma(-u_k^T v_c))v_c \quad \#$$

Negative Sampling loss only computes a fixed size of the outside vectors V instead on the whole thing as in naive softmax loss; Hence, its computationally more efficient.

(g) (2 point) Now we will repeat the previous exercise, but without the assumption that the K sampled words are distinct. 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, \dots, w_K and their outside vectors as $\mathbf{u}_1, \dots, \mathbf{u}_K$. In this question, you may not assume that the words are distinct. In other words, $w_i = w_j$ may be true when $i \neq j$ is true. Note that $o \notin \{w_1, \dots, w_K\}$. For a center word c and an outside word o , the negative sampling loss function is given by:

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)) \quad (6)$$

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

Compute the partial derivative of $J_{\text{neg-sample}}$ with respect to a negative sample \mathbf{u}_k . Please write your answers in terms of the vectors \mathbf{v}_c and \mathbf{u}_k , where $k \in [1, K]$. Hint: break up the sum in the loss function into two sums: a sum over all sampled words equal to u_k and a sum over all sampled words not equal to u_k .

$$J = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

\downarrow
 k is no longer distinct

Now $\frac{\partial J}{\partial \mathbf{u}_k}$ changes

$$i = 1, 2, \dots, K$$

\uparrow \uparrow \uparrow

$$\left\{ \text{condition: } w_i = w_K \right\}$$

$$\text{Take derivative on } - \sum_{\substack{i=1 \\ |w_i=w_K}}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

Rui Wang's DH
on Nook

$$\text{Equivalently, } - \sum_{\substack{i \in \{1, \dots, K\} : \\ |w_i=w_K}} \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

$$\begin{aligned} \text{Hence, } \frac{\partial J}{\partial \mathbf{u}_k} &= -\frac{\partial}{\partial \mathbf{u}_k} \left[\sum_{\substack{i=1 \\ |w_i=w_K}}^K \log(\sigma(-\mathbf{u}_i^\top \mathbf{v}_c)) \right] \\ &= - \left[\sum_{\substack{i=1 \\ |w_i=w_K}}^K \frac{-\mathbf{v}_c \sigma(-\mathbf{u}_i^\top \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_i^\top \mathbf{v}_c))}{\sigma(-\mathbf{u}_i^\top \mathbf{v}_c)} \right] \end{aligned}$$

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

$$J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) \quad (7)$$

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

Write down three partial derivatives:

- (i) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U}$
- (ii) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c$
- (iii) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w$ when $w \neq c$

Write your answers in terms of $\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{U}$ and $\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{v}_c$. This is very simple – each solution should be one line.

$$\left. \begin{aligned} \partial J_{\text{skip-gram}} / \partial \mathbf{v} &= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{J(\mathbf{v}_c, w_{t+j}, \mathbf{v})}{\partial \mathbf{v}} \\ \partial J_{\text{skip-gram}} / \partial \mathbf{v}_c &= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{J(\mathbf{v}_c, w_{t+j}, \mathbf{v})}{\partial \mathbf{v}_c} \\ \partial J_{\text{skip-gram}} / \partial w_c &= 0 \end{aligned} \right\} \quad w \neq c$$

2 Coding: Implementing word2vec

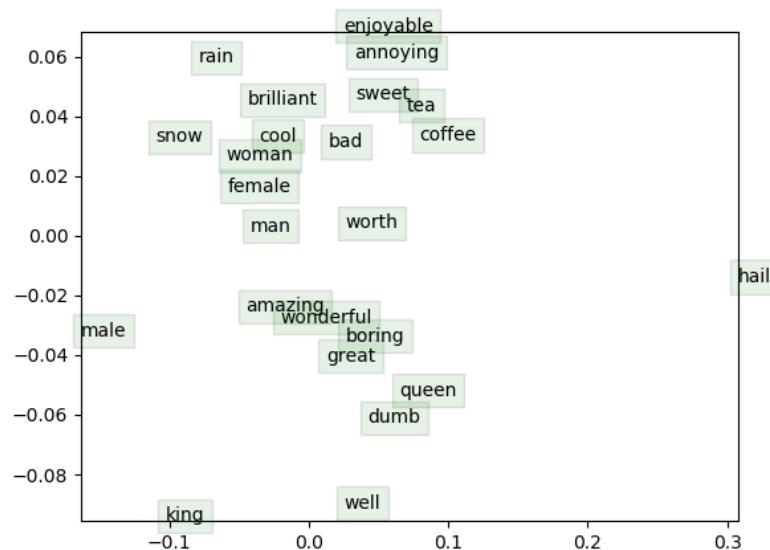
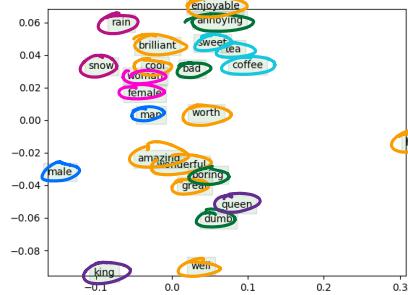
(c) (2 points) Show time! Now we are going to load some real data and train word vectors with everything you just implemented! We are going to use the Stanford Sentiment Treebank (SST) dataset to train word vectors, and later apply them to a simple sentiment analysis task. You will need to fetch the 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`.

Note: The training process may take a long time depending on the efficiency of your implementation and the compute power of your machine (an efficient implementation takes one to two hours). Plan accordingly!

After 40,000 iterations, the script will finish and a visualization for your word vectors will appear. It will also be saved as `word-vectors.png` in your project directory. **Include the plot in your homework write up.** Briefly explain in at most three sentences what you see in the plot.

```
Terminal Shell Edit View Window Help
a2 -- bash -- 112x46
iter 39580: 9.289658
iter 39590: 9.298853
iter 39600: 9.269859
iter 39610: 9.288015
iter 39620: 9.279245
iter 39630: 9.248681
iter 39640: 9.243792
iter 39650: 9.262415
iter 39660: 9.286713
iter 39670: 9.303074
iter 39680: 9.347061
iter 39690: 9.351811
iter 39700: 9.346526
iter 39710: 9.348751
iter 39720: 9.376277
iter 39730: 9.377825
iter 39740: 9.403598
iter 39750: 9.449631
iter 39760: 9.411052
iter 39770: 9.478179
iter 39780: 9.515252
iter 39790: 9.602017
iter 39800: 9.575306
iter 39810: 9.535311
iter 39820: 9.545193
iter 39830: 9.634330
iter 39840: 9.667367
iter 39850: 9.728679
iter 39860: 9.773566
iter 39870: 9.726360
iter 39880: 9.816947
iter 39890: 9.796976
iter 39900: 9.756038
iter 39910: 9.709420
iter 39920: 9.708453
iter 39930: 9.681717
iter 39940: 9.611494
iter 39950: 9.629689
iter 39960: 9.717683
iter 39970: 9.776979
iter 39980: 9.813174
iter 39990: 9.854022
iter 40000: 9.812206
sanity check: cost at convergence should be around or below 10
training took 3913 seconds
(base) DNa80d7af:a2 chih-hsunkao$
```

Observation :



I roughly used different colors to create clusters that I deem proper — And from this, there're indeed good clusters such as [amazing, wonderful], [woman, female], [sweet, tea, coffee], etc. However, an innegligible portion of the word embedding doesn't seem to be that ideal, for example, 'boring' clusters with positive words. Also, 'hail' is out of frame due to word2vec's inability to handle out-of-vocabulary words, as word2vec algo is based on the contextual window structure.