

# Cs224n - Assignment 2 - version 2021

Artiom Matvei

February 2024

## Note

See a2.pdf file for problem descriptions.

## Understanding word2vec

### Question 1.a

To show that

$$\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_o),$$

note that  $y_w = 0$  for all  $w$  except when  $w = o$  in which case  $y_w = 1$ . This is from the problem statement of question 1 in the assignment 2 pdf file (see 6th paragraph).

### Question 1.b

Recall that

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c),$$

where

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

To compute the partial derivative of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to  $\mathbf{v}_c$  let's first compute

$$\frac{\partial}{\partial \mathbf{v}_c} \log P(o|c)$$

where

$$P(o|c) = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)}.$$

We get

$$\begin{aligned}\frac{\partial}{\partial \mathbf{v}_c} \log P(o|c) &= \frac{\partial}{\partial \mathbf{v}_c} \log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= \frac{\partial}{\partial \mathbf{v}_c} \log \exp(\mathbf{u}_o^T \mathbf{v}_c) - \frac{\partial}{\partial \mathbf{v}_c} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c).\end{aligned}$$

Let's compute the derivatives of the two terms separately. The first term gets simplified because log and exp are inverse operations and because

$$\partial_x u^T x = \partial_x x^T u = u.$$

Thus we get

$$\frac{\partial}{\partial \mathbf{v}_c} \log \exp(\mathbf{u}_o^T \mathbf{v}_c) = \frac{\partial}{\partial \mathbf{v}_c} \mathbf{u}_o^T \mathbf{v}_c = \mathbf{u}_o.$$

Now let's compute the derivative of the second term.

$$\begin{aligned}\frac{\partial}{\partial \mathbf{v}_c} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) &= \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \sum_{w \in \text{Vocab}} \frac{\partial}{\partial \mathbf{v}_c} \exp(\mathbf{u}_w^T \mathbf{v}_c) \\ &= \frac{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) \mathbf{u}_x}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= \sum_{x \in \text{Vocab}} \frac{\exp(\mathbf{u}_x^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{u}_x \\ &= \sum_{x \in \text{Vocab}} P(x|c) \mathbf{u}_x.\end{aligned}$$

Putting the terms of the derivative together we get

$$\frac{\partial}{\partial \mathbf{v}_c} \log P(o|c) = \mathbf{u}_o - \sum_{x \in \text{Vocab}} P(x|c) \mathbf{u}_x \quad (1)$$

Having done this preliminary work we can go back to computing the partial derivative of the cost function by using the result we got in (1). Thus we have

$$\begin{aligned}\frac{\partial}{\partial \mathbf{v}_c} \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\frac{\partial}{\partial \mathbf{v}_c} \log P(o|c) \\ &= -(\mathbf{u}_o - \sum_{x \in \text{Vocab}} P(x|c) \mathbf{u}_x) \\ &= -(\mathbf{U} \mathbf{y} - \mathbf{U} \hat{\mathbf{y}}) \\ &= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y})\end{aligned}$$

Note that the penultimate equality comes from the following derivation where  $V = |\text{Vocab}|$ ,

$$\begin{aligned}
\sum_{x \in \text{Vocab}} P(x|c) \mathbf{u}_x &= P(O = 1|C = c) \mathbf{u}_1 + \dots \\
&\quad + P(O = x|C = c) \mathbf{u}_x + \dots + P(O = V|C = c) \mathbf{u}_V \\
&= \begin{pmatrix} u_1^1 P(O = 1|C = c) \\ \vdots \\ u_1^d P(O = 1|C = c) \end{pmatrix} + \dots + \begin{pmatrix} u_V^1 P(O = V|C = c) \\ \vdots \\ u_V^d P(O = V|C = c) \end{pmatrix} \\
&= \begin{pmatrix} u_1^1 P(O = 1|C = c) + \dots + u_V^1 P(O = V|C = c) \\ \vdots \\ u_1^d P(O = 1|C = c) + \dots + u_V^d P(O = V|C = c) \end{pmatrix} \\
&= \begin{pmatrix} u_1^1 & \dots & u_V^1 \\ \vdots & & \vdots \\ u_1^d & \dots & u_V^d \end{pmatrix} \begin{pmatrix} P(O = 1|C = c) \\ \vdots \\ P(O = V|C = c) \end{pmatrix}
\end{aligned}$$

### Question 1.c

We have

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, U) = -\log P(O = o|C = c) = \log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)}.$$

Thus

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{u}_w} \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, U) &= -\frac{\partial}{\partial \mathbf{u}_w} \log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \\
&= -\frac{\partial}{\partial \mathbf{u}_w} \left[ \log \exp(\mathbf{u}_o^T \mathbf{v}_c) - \log \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) \right] \\
&= \frac{\partial}{\partial \mathbf{u}_w} \log \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) - \frac{\partial}{\partial \mathbf{u}_w} \mathbf{u}_o^T \mathbf{v}_c
\end{aligned}$$

Now we'll compute the two terms separately by looking at the cases where  $w = o$  and where  $w \neq o$ .

If we assume that  $w \neq o$  then the 2nd term is 0 because the dot product is then constant with respect to  $\mathbf{u}_w$ . On the other hand, if we assume that  $w = o$  then the second term yields  $\mathbf{v}_c$

Now consider the 1st term,

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{u}_w} \log \sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) &= \frac{1}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \sum_{y \in \text{Vocab}} \frac{\partial}{\partial \mathbf{u}_w} \exp(\mathbf{u}_y^T \mathbf{v}_c) \\
&= \frac{1}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{u}_w} \exp(\mathbf{u}_w^T \mathbf{v}_c) \\
&= \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c)} \mathbf{v}_c \\
&= P(O = w | C = c) \mathbf{v}_c
\end{aligned}$$

Therefore now putting the two terms together, if  $w = o$  we have

$$\frac{\partial}{\partial \mathbf{u}_w} \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = (P(O = w | C = c) - 1) \mathbf{v}_c = (\hat{y}_w - 1) \mathbf{v}_c,$$

and if  $w \neq o$  we have

$$\frac{\partial}{\partial \mathbf{u}_w} \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = P(O = w | C = c) \mathbf{v}_c = \hat{y}_w \mathbf{v}_c.$$

### Question 1.d

Let  $N = |\text{Vocab}|$ . Since  $\mathbf{U}$  is a matrix of column vectors, i.e.

$$\mathbf{U} = \begin{pmatrix} \vdots & & \vdots \\ \mathbf{u}_1 & \dots & \mathbf{u}_N \\ \vdots & & \vdots \end{pmatrix},$$

the derivatives of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  need to be arranged in a manner to respect the shape convention. That is

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{U}} \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= \\
&= \begin{pmatrix} \vdots & & \vdots \\ \frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\mathbf{u}_1} & \dots & \frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\mathbf{u}_N} \\ \vdots & & \vdots \end{pmatrix} \\
&= \begin{pmatrix} \vdots & & \vdots \\ \hat{y}_1 \mathbf{v}_c & \dots & (\hat{y}_o - 1) \mathbf{v}_c & \dots & \hat{y}_N \mathbf{v}_c \\ \vdots & & \vdots & & \vdots \end{pmatrix} \\
&= \begin{pmatrix} \vdots & & \vdots & & \vdots \\ (\hat{y}_1 - y_1) \mathbf{v}_c & \dots & (\hat{y}_o - y_o) \mathbf{v}_c & \dots & (\hat{y}_N - y_N) \mathbf{v}_c \\ \vdots & & \vdots & & \vdots \end{pmatrix} \\
&= (\hat{\mathbf{y}} - \mathbf{y})^T \mathbf{v}_c
\end{aligned}$$

### Question 1.e

Let

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

First note that

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

Then

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

### Question 1.f

We have

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))$$

**Derivative with respect to central vector**

$$\frac{\partial}{\partial \mathbf{v}_c} \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\frac{\partial}{\partial \mathbf{v}_c} \log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \frac{\partial}{\partial \mathbf{v}_c} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))$$

Let's compute the first term,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{v}_c} \log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) &= \frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sigma(\mathbf{u}_o^T \mathbf{v}_c) \\ &= \frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \sigma(\mathbf{u}_o^T \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o \\ &= (1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o. \end{aligned}$$

Now let's compute the second term, which is very similar to the first one, thus by symmetry we get,

$$\frac{\partial}{\partial \mathbf{v}_c} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) = (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c))(-\mathbf{u}_k)$$

Putting both terms together we get

$$\frac{\partial}{\partial \mathbf{v}_c} \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{u}_o + \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{u}_k$$

**Derivative with respect to  $\mathbf{u}_o$**

$$\frac{\partial}{\partial \mathbf{u}_o} \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\frac{\partial}{\partial \mathbf{u}_o} \log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \frac{\partial}{\partial \mathbf{u}_o} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))$$

Note that the second term above is zero because every summand inside the summation is constant with respect to  $\mathbf{u}_o$  because  $o \notin \{w_1, \dots, w_K\}$  as written in the problem statement. Therefore we get,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{u}_o} \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) &= -\frac{\partial}{\partial \mathbf{u}_o} \log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) \\ &= -\frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{u}_o} \sigma(\mathbf{u}_o^T \mathbf{v}_c) \\ &= -\frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \sigma(\mathbf{u}_o^T \mathbf{v}_c) (1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \frac{\partial}{\partial \mathbf{u}_o} (\mathbf{u}_o^T \mathbf{v}_c) \\ &= (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{v}_c \end{aligned}$$

**Derivative with respect to a negative sample  $\mathbf{u}_k$**

In the problem statement we assumed that the  $K$  negative samples are distinct. Let  $\hat{k} \in \{1, \dots, K\}$  be the variable with respect to which we'll compute the variation. Therefore we have,

$$\frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))$$

Note that the first term is equal to zero because the first term is constant with respect to  $\mathbf{u}_{\hat{k}}$  as  $\hat{k} \neq o$ . Also note that all the terms in the second term are also equal to zero because they are also constant with respect to  $\mathbf{u}_{\hat{k}}$  except for the only one when  $k = \hat{k}$ . Therefore,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) &= -\frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \log(\sigma(-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c)) \\ &= -\frac{1}{\sigma(-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c)} \sigma(-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c) (1 - \sigma(-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c)) \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} (-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c) \\ &= (1 - \sigma(-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c)) \mathbf{v}_c \end{aligned}$$

Note that negative sampling is much more efficient to compute than the naive-softmax loss because the naive-softmax contains a summation spanning over the whole vocabulary while negative sampling contains a summation over  $K$  elements only.

### Question 1.g

We have

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))$$

where  $w_i = w_j$  may be true when  $i \neq j$ . Let  $\hat{k} \in \{1, \dots, K\}$  be the variable with respect to which we'll compute the variation. Note then that the derivative of the first term with respect to  $\mathbf{u}_{\hat{k}}$  is equal to zero. Therefore, letting

$$A = |\{k \in \{1, \dots, K\} | w_k = w_{\hat{k}}\}|,$$

we get

$$\begin{aligned} \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) &= -\frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \left[ \sum_{\{k | w_k = w_{\hat{k}}\}} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) + \sum_{\{k | w_k \neq w_{\hat{k}}\}} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \right] \\ &= - \sum_{\{k | w_k = w_{\hat{k}}\}} \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) - \sum_{\{k | w_k \neq w_{\hat{k}}\}} \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \\ &= - \sum_{\{k | w_k = w_{\hat{k}}\}} \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) - \sum_{\{k | w_k \neq w_{\hat{k}}\}} 0 \\ &= -A \frac{\partial}{\partial \mathbf{u}_{\hat{k}}} \log(\sigma(-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c)) \\ &= A(1 - \sigma(-\mathbf{u}_{\hat{k}}^T \mathbf{v}_c)) \mathbf{v}_c \end{aligned}$$

### Question 1.h

Skip-gram derivative with respect to  $\mathbf{U}$

$$\begin{aligned} \partial_{\mathbf{U}} \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) &= \partial_{\mathbf{U}} \sum_{-m \leq j \leq m, j \neq 0} \mathbf{J}(\mathbf{v}_c, w_t + j, \mathbf{U}) \\ &= \sum_{-m \leq j \leq m, j \neq 0} \partial_{\mathbf{U}} \mathbf{J}(\mathbf{v}_c, w_t + j, \mathbf{U}) \end{aligned}$$

**Skip-gram derivative with respect to  $v_c$**

$$\begin{aligned}\partial_{v_c} \mathbf{J}_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U) &= \partial_{v_c} \sum_{-m \leq j \leq m, j \neq 0} \mathbf{J}(v_c, w_t + j, U) \\ &= \sum_{-m \leq j \leq m, j \neq 0} \partial_{v_c} \mathbf{J}(v_c, w_t + j, U)\end{aligned}$$

**Skip-gram derivative with respect to  $v_w$  for  $w \neq c$**

$$\begin{aligned}\partial_{v_w} \mathbf{J}_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U) &= \partial_{v_w} \sum_{-m \leq j \leq m, j \neq 0} \mathbf{J}(v_c, w_t + j, U) \\ &= \sum_{-m \leq j \leq m, j \neq 0} \partial_{v_w} \mathbf{J}(v_c, w_t + j, U) = 0\end{aligned}$$

where the last line is true because  $\mathbf{J}$  is constant with respect to  $v_w$ .

## Implementing word2vec

After implementing word2vec as per the instruction of the assignment 2 pdf document, I got the word embedding seen in Figure 1. As expected, antonyms group together as they could be used fairly interchangeably in the same sentence (e.g. great vs boring), and synonyms are also close together for the same reason (e.g. amazing and wonderful).



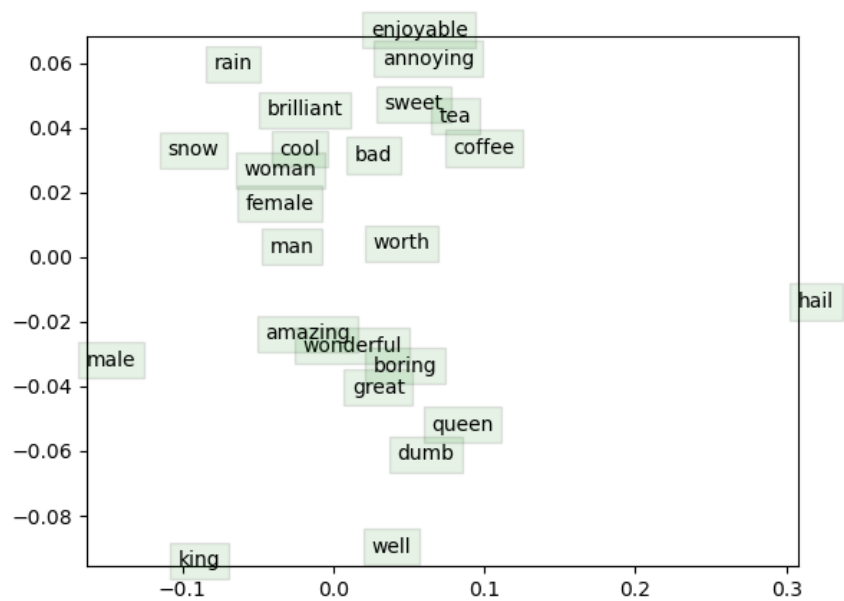


Figure 1: Word embedding after training the word2vec