1. Linear Embedding - GLoVe

   1.1. GLoVE Parameter Count

   Since the dimension of $\boldsymbol{w}_i$ and $\boldsymbol{w}_i$ are $d$ and the dimension of the biases are 1, there are $2d + 2$ parameters to train for each vocabulary. Thus, we can get that the GLoVe model have $2V(d+1)$ trainable parameters.

   1.2. Expression for gradient $\frac{\partial L}{\partial \boldsymbol{w}_i}$

   According to the given loss function, we can get $\frac{\partial L}{\partial \boldsymbol{w}_i}$ should be

   $$\frac{\partial L}{\partial \boldsymbol{w}_i} = \frac{\partial}{\partial \boldsymbol{w}_i} \sum_{i,j=1}^{V} (\boldsymbol{w}_i^T \tilde{\boldsymbol{w}}_i + b_i + \tilde{b}_j - \log X_{ij})^2$$

   $$= 2 \sum_{i,j=1}^{V} (\boldsymbol{w}_i^T \tilde{\boldsymbol{w}}_i + b_i + \tilde{b}_j - \log X_{ij}) \tilde{\boldsymbol{w}}_i$$

   1.3. Implement the gradient update of GLoVE

   The code for this question is shown below.

```
1  def grad_GLoVE(W, W_tilde, b, b_tilde, log_co_occurence):
2    "Return the gradient of GLoVE objective w.r.t W and b."
3    "INPUT: W − Vxd; W_tilde − Vxd; b − Vx1; b_tilde − Vx1; log_co_occurence: VxV"
4    "OUTPUT: grad_W − Vxd; grad_W_tilde − Vxd, grad_b − Vx1, grad_b_tilde − Vx1"
5    n, _ = log_co_occurence.shape
6
7    if not W_tilde is None and not b_tilde is None:
8    ############################  YOUR CODE HERE  ############################
9      loss = (W @ W_tilde.T + b @ np.ones([1,n]) + np.ones([n,1]) @ b_tilde.T − (
       log_co_occurence))
10     grad_W = 2 ∗ (loss @ W_tilde)
11     grad_W_tilde = 2 ∗ (loss.T @ W)
12     grad_b = 2 ∗ (np.ones([1,n]) @ loss).T
13     grad_b_tilde = 2 ∗ (np.ones([1,n]) @ loss).T
14    ########################################################################
15    else:
16      loss = (W @ W.T + b @ np.ones([1,n]) + np.ones([n,1])@b.T − 0.5∗(
       log_co_occurence + log_co_occurence.T))
17     grad_W = 4 ∗(W.T @ loss).T
18     grad_W_tilde = None
19     grad_b = 4 ∗ (np.ones([1,n]) @ loss).T
20     grad_b_tilde = None
21
22    return grad_W, grad_W_tilde, grad_b, grad_b_tilde
23
```

   1.4. Effects of embedding dimension

2. Network architecture

   2.1. Number of parameters in neural network model

   Let's find the number of the trainable parameters of the word embedding weights first. Since there are $V$ words in the dictionary and the dimension of the word embedding layer is $N \times D$, we can get that there are $V \times D$ trainable parameters.

   For weights between the word embedding layer and the hidden layer, since there are $H$ units in the hidden layer, the dimension of the matrix that connects two layers should be $ND \times H$, which is the number of trainable parameters.

   Thus, for the biases of the hidden layer, there should be $H \times 1$ trainable parameters.

   Similarly, since the output layer consists of $V$ words, there are $V \times H$ trainable parameters for

weights between the hidden layer and the output layer.

Thus, for the biases of the output layer, there should be $V \times 1$ trainable parameters.

Since $V$ is much larger than other variables, we only need to consider the part that depends on $V$.

Thus, we can get that weights between the hidden layer and the output layer, $hid\_to\_output\_weights$, has the largest number of trainable parameters since $H > D$.

2.2. Number of parameters in n-gram mode

For each gram, we can choose any word from $V$ words. Thus, there are $V^N$ number of combinations for the previous $N$ words. For the prediction, since the output layer is a softmax over the $V$ words, there are $V$ words. Thus, there are $V^{N+1}$ entries in the $n$-gram model scale with $N$.

2.3. Comparing neural network and n-gram model scaling

3. Training the Neural Network

3.1. Implement gradient with respect to output layer inputs

The code for this question is shown below.

```
def compute_loss_derivative(self, output_activations, expanded_target_batch,
    target_mask):
    """Compute the derivative of the multiple target position cross-entropy
    loss function \n"

        For example:

            [y_{0}  ....   y_{V-1}] [y_{V}, ..., y_{2*V-1}] [y_{2*V} ... y_{i,3*V-1}] [
    y_{3*V} ... y_{i,4*V-1}]

        Where for colum j + n*V,

            y_{j + n*V} = e^{z_{j + n*V}} / \sum_{m=0}^{V-1} e^{z_{m + n*V}}, for
    n=0,...,N-1

        This function should return a dC / dz matrix of size [batch_size x (
    vocab_size * context_len)],
        where each row i in dC / dz has columns 0 to V-1 containing the gradient
    the 1st output
        context word from i-th training example, then columns vocab_size to 2*
    vocab_size - 1 for the 2nd
        output context word of the i-th training example, etc.

        C is the loss function summed acrossed all examples as well:

            C = -\sum_{i,j,n} mask_{i,n} (t_{i, j + n*V} log y_{i, j + n*V}), for
    j=0,...,V, and n=0,...,N

        where mask_{i,n} = 1 if the i-th training example has n-th context word as
     the target,
        otherwise mask_{i,n} = 0.

        The arguments are as follows:

            output_activations - A [batch_size x (context_len * vocab_size)]
    tensor,
                for the activations of the output layer, i.e. the y_j's.
            expanded_target_batch - A [batch_size (context_len * vocab_size)]
    tensor,
                where expanded_target_batch[i,n*V:(n+1)*V] is the indicator vector
     for
                the n-th context target word position, i.e. the (i, j + n*V) entry
    is 1 if the
                i'th example, the context word at position n is j, and 0 otherwise
    .
            target_mask - A [batch_size x context_len x 1] tensor, where
    target_mask[i,n] = 1
                if for the i'th example the n-th context word is a target position
    , otherwise 0
```

```
34
35          Outputs:
36              loss_derivative − A [batch_size x (context_len * vocab_size)] matrix,
37                  where loss_derivative[i,0:vocab_size] contains the gradient
38                  dC / dz_0 for the i−th training example gradient for 1st output
39                  context word, and loss_derivative[i,vocab_size:2*vocab_size] for
40                  the 2nd output context word of the i−th training example, etc.
41          """
42
43          ############################    YOUR CODE HERE
       ##############################
44          # Loss
45          loss = output_activations − expanded_target_batch
46          expanded_mask = np.repeat(target_mask, self.vocab_size, axis=1).reshape(
       loss.shape)
47          return np.multiply(expanded_mask, loss)
48          #
       ##################################################################################
49
```

## 3.2. Implement gradient with respect to parameters

The code for this question is shown below.

```
1  def back_propagate(self, input_batch, activations, loss_derivative):
2          """Compute the gradient of the loss function with respect to the trainable
       parameters
3          of the model. The arguments are as follows:
4
5              input_batch − the indices of the context words
6              activations − an Activations class representing the output of Model.
       compute_activations
7              loss_derivative − the matrix of derivatives computed by
       compute_loss_derivative
8
9          Part of this function is already completed, but you need to fill in the
       derivative
10         computations for hid_to_output_weights_grad, output_bias_grad,
       embed_to_hid_weights_grad,
11         and hid_bias_grad. See the documentation for the Params class for a
       description of what
12         these matrices represent."""
13
14         # The matrix with values dC / dz_j, where dz_j is the input to the jth
       hidden unit,
15         # i.e. h_j = 1 / (1 + e^{−z_j})
16         hid_deriv = np.dot(loss_derivative, self.params.hid_to_output_weights) \
17                     * activations.hidden_layer * (1. − activations.hidden_layer)
18
19         ############################    YOUR CODE HERE
       ##############################
20         hid_to_output_weights_grad = loss_derivative.T @ activations.hidden_layer
21         output_bias_grad = np.sum(loss_derivative, axis=0)
22         embed_to_hid_weights_grad = hid_deriv.T @ activations.embedding_layer
23         hid_bias_grad = np.sum(hid_deriv, axis=0)
24         #
       ##################################################################################
25
26         # The matrix of derivatives for the embedding layer
27         embed_deriv = np.dot(hid_deriv, self.params.embed_to_hid_weights)
28
29         # Embedding layer
30         word_embedding_weights_grad = np.zeros((self.vocab_size, self.
       embedding_dim))
31         for w in range(self.context_len):
32             word_embedding_weights_grad += np.dot(self.indicator_matrix(
       input_batch[:, w:w+1], mask_zero_index=False).T,
33                                                   embed_deriv[:, w * self.
       embedding_dim:(w + 1) * self.embedding_dim])
```

```
34
35          return Params( word_embedding_weights_grad , embed_to_hid_weights_grad ,
     hid_to_output_weights_grad ,
36                          hid_bias_grad , output_bias_grad )
37
```

### 3.3. Print the gradients The output for print_gradients() is shown below.

```
1  loss_derivative [2 , 5]  0.0
2  loss_derivative [2 , 121]  0.0
3  loss_derivative [5 , 33]  0.0
4  loss_derivative [5 , 31]  0.0
5
6  param_gradient.word_embedding_weights [27 , 2]  0.0
7  param_gradient.word_embedding_weights [43 , 3]  0.011596892511489458
8  param_gradient.word_embedding_weights [22 , 4]  −0.0222670623817297
9  param_gradient.word_embedding_weights [2 , 5]  0.0
10
11 param_gradient.embed_to_hid_weights [10 , 2]  0.3793257091930164
12 param_gradient.embed_to_hid_weights [15 , 3]  0.01604516132110917
13 param_gradient.embed_to_hid_weights [30 , 9]  −0.4312854367997419
14 param_gradient.embed_to_hid_weights [35 , 21]  0.06679896665436337
15
16 param_gradient.hid_bias [10]  0.023428803123345148
17 param_gradient.hid_bias [20]  −0.024370452378874197
18
19 param_gradient.output_bias [0]  0.000970106146902794
20 param_gradient.output_bias [1]  0.16868946274763222
21 param_gradient.output_bias [2]  0.0051664774143909235
22 param_gradient.output_bias [3]  0.15096226471814364
23
```

### 3.4. Run model training

## 4. Arithmetics and Analysis
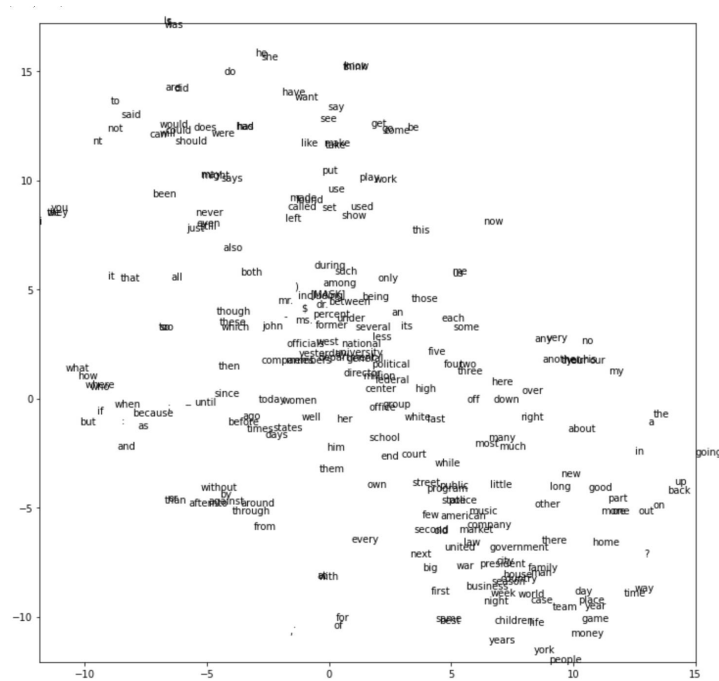
### 4.1. t-SNE



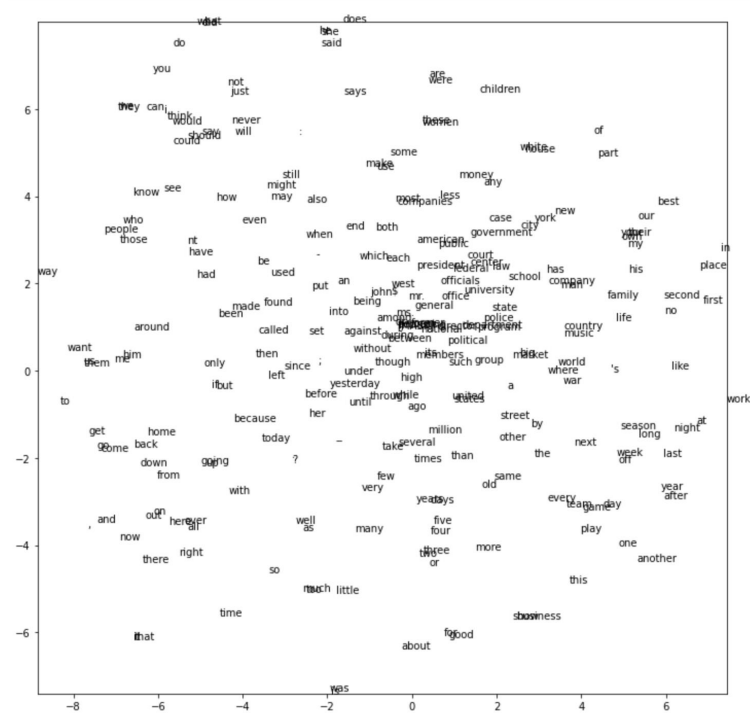Figure 1: The tsne plot representation using the trained weights.
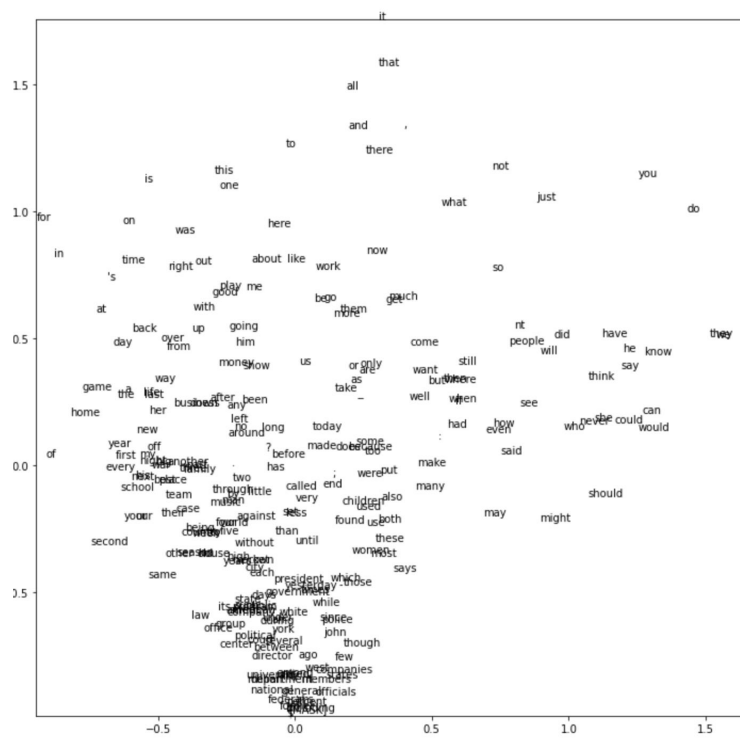
Figure 2: The tsne plot GLoVE representation.



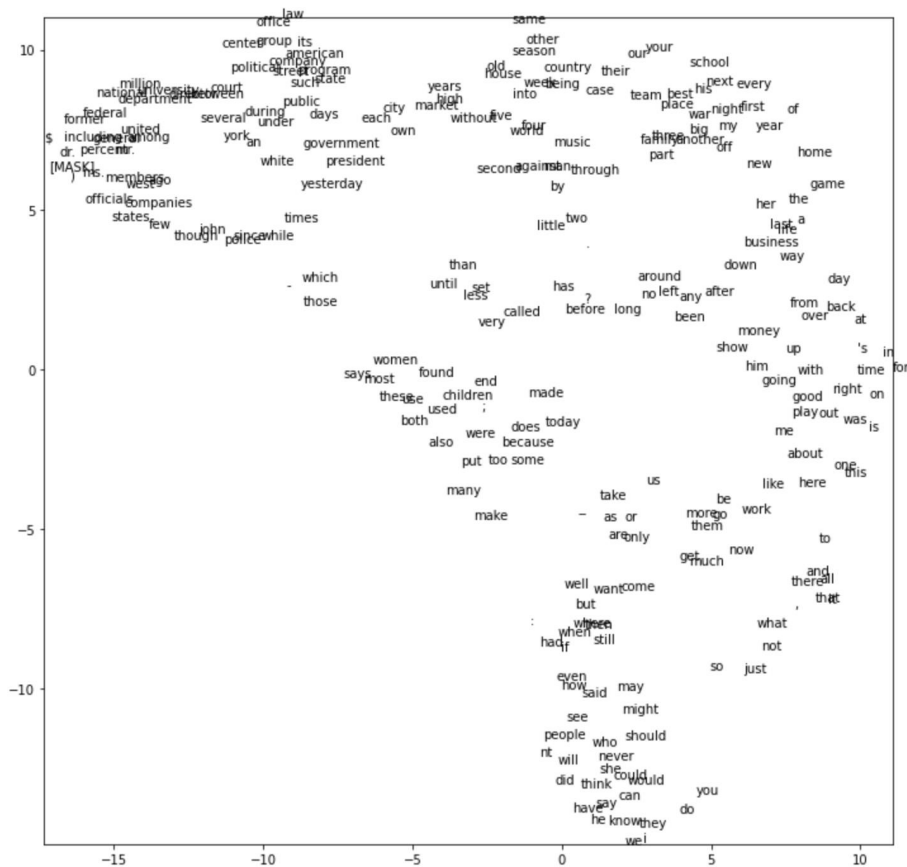Figure 3: The 2d GLoVe representation.

Figure 4: The tsne plot GLoVE representation.

According to Figure 1 (tsne_plot_representation), we can see that words with the similar function in sentences or the structure in terms of the speech gather together. For example, on the top left of the plot, there is a cluster of auxiliary verbs including "would", "could", and "should", and a cluster of question words such as "what", "how", and "who" on the middle left of the plot.

For Figure 1 and 2 (tsne_plot_GLoVE_representation), we can see that words in Figure 1 are distributed like a linear on the main diagonal, while words in Figure 2 diverge circularly around the biggest cluster on the lower middle. In addition, we found that positions of different clusters are different in two plots.

Comparing with Figure 1 and 2, we can see that words in Figure 3 (plot_2d_GLoVe_representation) seem to be more clustered and are distributed like a fan shape. Many nouns gather together on the bottom left cluster, which is the largest cluster and words are fan out upward.

## 4.2. Word Analogy Arithmetic

### 4.2.1. Specific example

The results are shown below.

```
## GloVe embeddings
The top 10 closest words to emb(he) − emb(him) + emb(her) are:
he: 1.4213098857979793
she: 1.48167433432594
said: 2.1025960106397767
then: 2.2720425987761406
does: 2.301964867719902
says: 2.318047293286045
who: 2.328984314854128
where: 2.334702431567161
did: 2.353623598835888
```

```
12  should: 2.4126428205989865
13
14  # Concatenation of W_final_asym, W_tilde_final_asym
15  The top 10 closest words to emb(he) − emb(him) + emb(her) are:
16  he: 2.046826000951795
17  she: 2.3455038844018743
18  i: 3.0624522787351487
19  we: 3.2848647174761094
20  they: 3.390910580609287
21  you: 4.568945007203308
22  john: 4.805241000654006
23  program: 5.084420284826234
24  president: 5.104152796566877
25  never: 5.111163924178705
26
27  # Averaging asymmetric GLoVE vectors
28  The top 10 closest words to emb(he) − emb(him) + emb(her) are:
29  he: 1.0154702232698416
30  she: 1.0744126028585648
31  should: 1.6078139338035942
32  could: 1.6799073061855805
33  i: 1.693953840244398
34  would: 1.70002096216810
35  did: 1.766206937557185
36  can: 1.7744377144463797
37  might: 1.7765376616413824
38  will: 1.7931360498829227
39
40  ## Neural Netework Word Embeddings
41  The top 10 closest words to emb(he) − emb(him) + emb(her) are:
42  he: 2.4284684644619032
43  she: 17.4415802699889
44  have: 25.921497697983263
45  they: 25.981587972296392
46  want: 26.437644546989542
47  we: 27.128094534488834
48  i: 27.215833550319473
49  but: 28.03028938337095
50  about: 28.163403568035555
51  this: 28.531350495330678
52
```

According to the outputs, we can see that the closest word that is not "he", "him", or "her" is "she" for all 4 different arithmetic. The corresponding distances are shown above.

According to four plots, they all show the parallelogram property of the quadruplets approximately. Figure 5, 6, 7, and 8 show how they present in the corresponding plot.
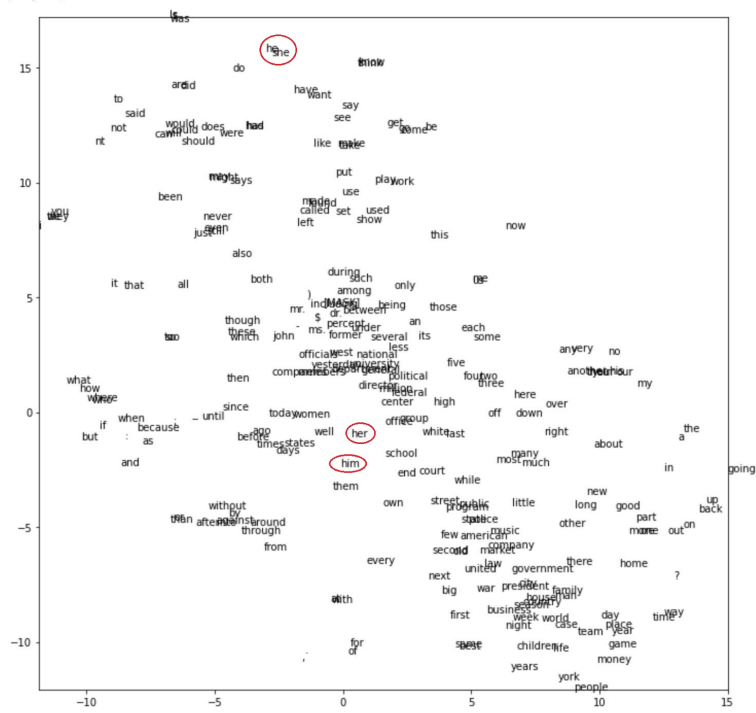
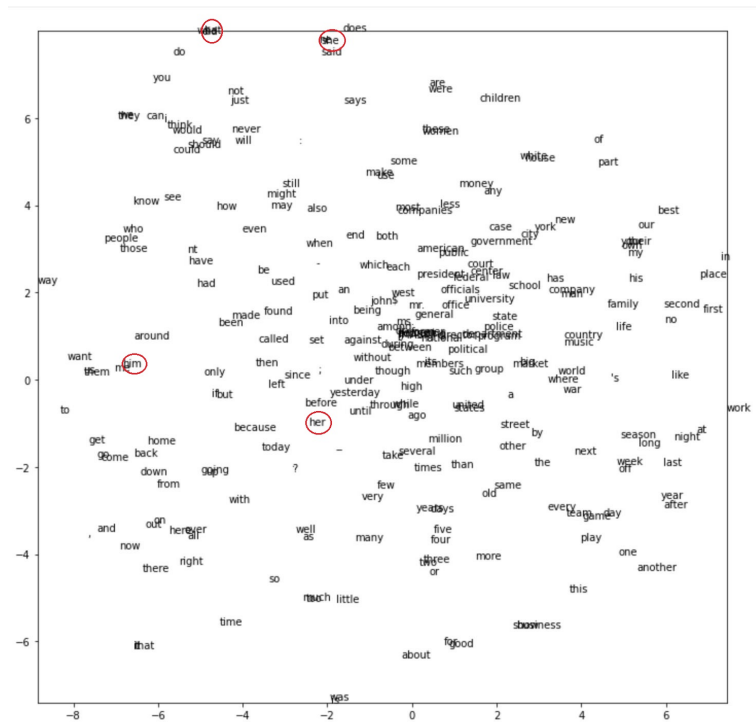Figure 5: The tsne plot representation using the trained weights.



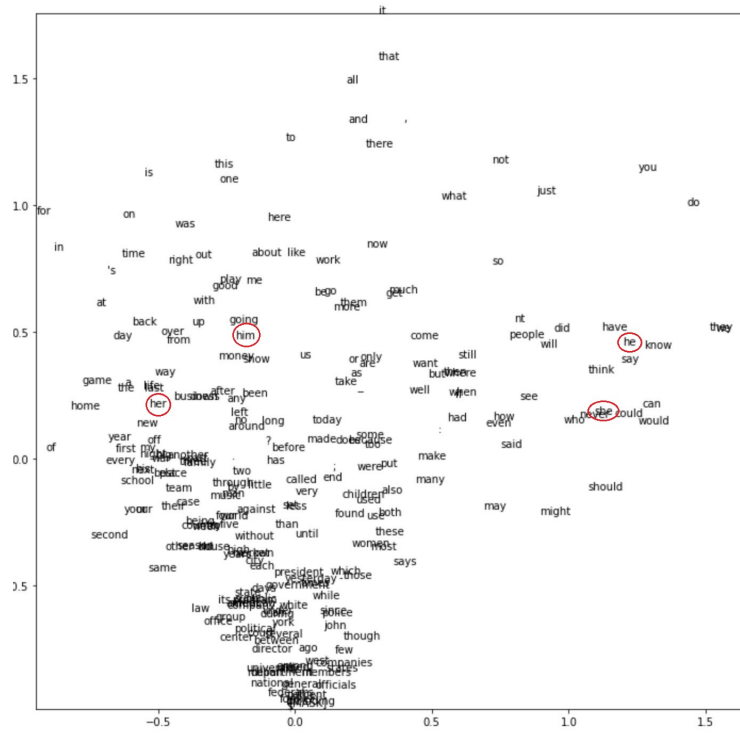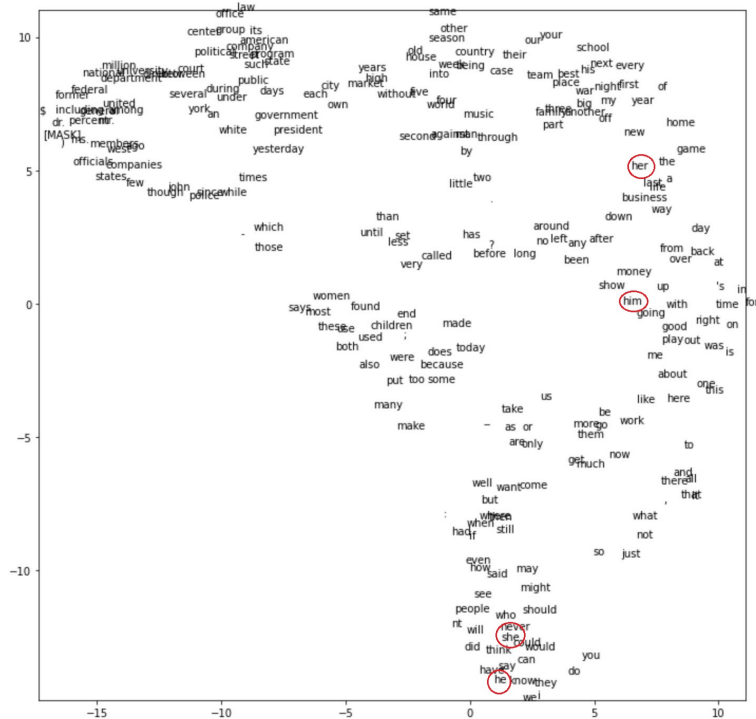Figure 6: The tsne plot GLoVE representation.

Figure 7: The 2d GLoVe representation.



Figure 8: The tsne plot GLoVE representation.