

Homework 3

Lecturer: Russ Salakhutdinov

Name: Yuan Liu(yuanl4)

1 Problem 1, 4-gram language model

- **Embedding layer:**

- **Input:** three word index: $w_{i-1}, w_{i-2}, w_{i-3}$.
- **Parameters:** A 2-dimensional matrix C of size $V \times D$.
- **Output:** Vector representations for these three words: $C_{w_{i-1},:}, C_{w_{i-2},:}, C_{w_{i-3},:}$, where $C_{j,:}$ means the j^{th} row of the matrix C .

- **Embedding to Hidden:**

- **Input:** $C_{w_{i-1},:}, C_{w_{i-2},:}, C_{w_{i-3},:}$ from the output of the embedding layer. The size of each element of input $1 \times D$.
- **Parameters:** Three embed_to_hidden_weights matrix $W^{(1)}, W^{(2)}, W^{(3)}$, each matrix's size is $D \times H$. A hidden_bias vector b^{hidden} of size $1 \times H$.
- **Output:** $C_{w_{i-1},:}W^{(1)} + C_{w_{i-2},:}W^{(2)} + C_{w_{i-3},:}W^{(3)} + b^{hidden}$

- **Tanh Layer:**

- **Input:** $A = C_{w_{i-1},:}W^{(1)} + C_{w_{i-2},:}W^{(2)} + C_{w_{i-3},:}W^{(3)} + b^{hidden}$ of size $1 \times H$
- **Parameters:** None
- **Output:** $\tanh(A)$

- **Hidden to Output:**

- **Input:** $H = \tanh(A)$ of size $1 \times H$.
- **Parameters:** The hidden_to_output_weight W^{out} of size $H \times V$ and the output_bias b^{out} of size $1 \times V$.
- **Output:** $HW^{out} + b^{out}$

- **Softmax layer:**

- **Input:** $O = HW^{out} + b^{out}$ of size $1 \times V$.
- **Parameters:** None
- **Output:** $S_i = \frac{e^{O_i}}{\sum_{j=1}^V e^{O_j}}$, S is a matrix of size $1 \times V$.

- **Loss:** $loss = -\log S_{w_i}$

Now we can calculate the derivation:

$$\begin{aligned}
\frac{\partial \text{loss}}{\partial \mathbf{S}_{w_i}} &= -\frac{1}{\mathbf{S}_{w_i}} \\
\frac{\partial \mathbf{S}_{w_i}}{\partial \mathbf{O}_j} &= \frac{1(w_i = j)e^{\mathbf{O}_{w_i}}(\sum_{k=1}^V e^{\mathbf{O}_k}) - e^{\mathbf{O}_{w_i}}e^{\mathbf{O}_j}}{(\sum_{k=1}^V e^{\mathbf{O}_k})^2} \\
\frac{\partial \text{loss}}{\partial \mathbf{O}_j} &= \frac{\partial \text{loss}}{\partial \mathbf{S}_{w_i}} \frac{\partial \mathbf{S}_{w_i}}{\partial \mathbf{O}_j} = -(1(w_i = j) - \mathbf{S}_j) \\
\frac{\partial \mathbf{O}_j}{\partial \mathbf{H}_i} &= \mathbf{W}_{i,j}^{\text{out}}, \quad \frac{\partial \mathbf{O}_j}{\partial \mathbf{W}_{i,j}^{\text{out}}} = \mathbf{H}_i, \quad \frac{\partial \mathbf{O}_j}{\partial \mathbf{b}_j^{\text{out}}} = 1 \\
\frac{\partial \text{loss}}{\partial \mathbf{H}_i} &= \sum_{j=1}^V \frac{\partial \text{loss}}{\partial \mathbf{O}_j} \frac{\partial \mathbf{O}_j}{\partial \mathbf{H}_i} \\
\frac{\partial \text{loss}}{\partial \mathbf{W}_{i,j}^{\text{out}}} &= \frac{\partial \text{loss}}{\partial \mathbf{O}_j} \frac{\partial \mathbf{O}_j}{\partial \mathbf{W}_{i,j}^{\text{out}}} = -(1(w_i = j) - \mathbf{S}_j)\mathbf{H}_i \tag{1}
\end{aligned}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{b}_i^{\text{out}}} = \frac{\partial \text{loss}}{\partial \mathbf{O}_i} \frac{\partial \mathbf{O}_i}{\partial \mathbf{b}_i^{\text{out}}} = -(1(w_i = i) - \mathbf{S}_i) \tag{2}$$

$$\frac{\partial \mathbf{H}_j}{\partial \mathbf{A}_i} = (1 - \mathbf{H}_j^2) * 1(i = j)$$

$$\frac{\partial \mathbf{A}_i}{\partial \mathbf{C}_{w_{i-1},j}} = \mathbf{W}_{j,i}^{(1)}, \quad \frac{\partial \mathbf{A}_i}{\partial \mathbf{C}_{w_{i-2},j}} = \mathbf{W}_{j,i}^{(2)}, \quad \frac{\partial \mathbf{A}_i}{\partial \mathbf{C}_{w_{i-3},j}} = \mathbf{W}_{j,i}^{(3)}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{C}_{w_{i-1},j}} = \sum_{i,k} \frac{\partial \text{loss}}{\partial \mathbf{H}_i} \frac{\partial \mathbf{H}_i}{\partial \mathbf{A}_k} \frac{\partial \mathbf{A}_k}{\partial \mathbf{C}_{w_{i-1},j}} = \sum_i \frac{\partial \text{loss}}{\partial \mathbf{H}_i} \frac{\partial \mathbf{H}_i}{\partial \mathbf{A}_i} \frac{\partial \mathbf{A}_i}{\partial \mathbf{C}_{w_{i-1},j}} = \sum_i \frac{\partial \text{loss}}{\partial \mathbf{H}_i} (1 - \mathbf{H}_i^2) \mathbf{W}_{j,i}^{(1)}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{C}_{w_{i-1},j}} = \sum_i \frac{\partial \text{loss}}{\partial \mathbf{H}_i} (1 - \mathbf{H}_i^2) \mathbf{W}_{j,i}^{(1)} \tag{3}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{C}_{w_{i-2},j}} = \sum_i \frac{\partial \text{loss}}{\partial \mathbf{H}_i} (1 - \mathbf{H}_i^2) \mathbf{W}_{j,i}^{(2)} \tag{4}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{C}_{w_{i-3},j}} = \sum_i \frac{\partial \text{loss}}{\partial \mathbf{H}_i} (1 - \mathbf{H}_i^2) \mathbf{W}_{j,i}^{(3)} \tag{5}$$

$$\frac{\partial \mathbf{A}_j}{\partial \mathbf{W}_{i,j}^{(1)}} = \mathbf{C}_{w_{i-1},i}, \quad \frac{\partial \mathbf{A}_j}{\partial \mathbf{W}_{i,j}^{(2)}} = \mathbf{C}_{w_{i-2},i}, \quad \frac{\partial \mathbf{A}_j}{\partial \mathbf{W}_{i,j}^{(3)}} = \mathbf{C}_{w_{i-3},i}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{W}_{i,j}^{(1)}} = \frac{\partial \text{loss}}{\partial \mathbf{H}_j} (1 - \mathbf{H}_j^2) \mathbf{C}_{w_{i-1},i} \tag{6}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{W}_{i,j}^{(2)}} = \frac{\partial \text{loss}}{\partial \mathbf{H}_j} (1 - \mathbf{H}_j^2) \mathbf{C}_{w_{i-2},i} \tag{7}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{W}_{i,j}^{(3)}} = \frac{\partial \text{loss}}{\partial \mathbf{H}_j} (1 - \mathbf{H}_j^2) \mathbf{C}_{w_{i-3},i} \tag{8}$$

$$\frac{\partial \mathbf{A}_i}{\partial \mathbf{b}_j^{\text{hidden}}} = 1(i = j)$$

$$\frac{\partial \text{loss}}{\partial \mathbf{b}_j^{\text{hidden}}} = \frac{\partial \text{loss}}{\partial \mathbf{H}_j} (1 - \mathbf{H}_j^2) \tag{9}$$

The derivation of weights is listed in (1) (9)

2 Problem 2, LSTM & GRU

- (a) LSTM contains 3 gates: *input* gates i_t , *forget* gates f_t and *output* gates o_t .
GRU contains 2 gates: *update* gate z_t and *reset* gate r_t .

- (b) LSTM:

- *forget* gate: to decide what information should be throw away form the cell state. 1 means ‘completely remember’ and 0 means ‘completely forget’.
- *input* gate: to decide which value in the cell state should be updated. 1 means ‘add’ and 0 means ‘ignore’.
- *output* gate: to decide what to output. 1 means ‘output’ and 0 means ‘don’t output’.

GRU:

- *update* gate: to decide which value from \tilde{h}_t should be added to h_t and what value from h_{t-1} should be forget. 1 means the corresponding value from \tilde{h}_t should be remembered and the corresponding value from h_{t-1} should be forgot.
- *reset* gate: to decide what part of h_{t-1} should be computed to get \tilde{h}_t .

- (c) The output of LSTM is the memory unit \mathbf{C} and hidden content h . The output of GRU only contains the hidden content h . The LSTM controls the flow of information according to both \mathbf{C} and h , while the GRU only expose the full hidden content without any control.
- (d) LSTM: $W_f : n \times m$, $U_f : n \times n$, $b_f : n$, $W_i : n \times m$, $U_i : n \times n$, $b_i : n$, $W_o : n \times m$, $U_o : n \times n$, $b_o : n$, $W_c : n \times m$, $U_c : n \times n$, $b_c : n$. Totally, the LSTM contains $4(n^2 + mn + n)$ parameters.
GRU: $W_z : n \times m$, $U_z : n \times n$, $b_z : n$, $W_r : n \times m$, $U_r : n \times n$, $b_r : n$, $W : n \times m$, $U : n \times n$, $b : n$. Totally, the GRU contains $3(n^2 + mn + n)$ parameters.
- (e) The GRU might take less time to train. Because GRU contains fewer parameters than the LSTM. Moreover, GRU only contains 2 gates and its structure is simpler, so it will be more computationally efficient.