

Problem 1

(a) Proof as the following:

$$\begin{aligned}\frac{\partial g}{\partial z} &= -\frac{1}{(1+e^{-z})^2} \cdot (e^{-z}) \cdot -1 \\ &= \frac{(e^{-z})}{(1+e^{-z})^2} \\ &= \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} \\ &= g(z)(1-g(z))\end{aligned}$$

(b) Proof as follows:

$$\begin{aligned}1-g(z) &= \frac{e^{-z}}{1+e^{-z}} \\ &= \frac{1}{e^z+1} \\ &= g(-z)\end{aligned}$$

Problem 2

(a) As g is convex, we have the following relation (equation 1),

$$g(t[\langle w_1, x \rangle + y] + (1-t)[\langle w_2, x \rangle + y]) \leq tg(\langle w_1, x \rangle + y) + (1-t)g(\langle w_2, x \rangle + y) \quad (1)$$

$\forall t \in [0, 1]$ and $\forall w_1, w_2 \in \mathbb{R}^d$. Therefore, we can do the following substitution:

$$\begin{aligned}f(tw_1 + (1-t)w_2) &= g(\langle tw_1 + (1-t)w_2, x \rangle + y) \\ &= g(t[\langle w_1, x \rangle + y] + (1-t)[\langle w_2, x \rangle + y]) \\ &\leq tg(\langle w_1, x \rangle + y) + (1-t)g(\langle w_2, x \rangle + y) \\ &= tf(w_1) + (1-t)f(w_2)\end{aligned}$$

the equality holds when equality of equation 1 holds. Conclude that f is also convex if g is convex.

(b) $\forall t \in [0, 1]$ and $x_1, x_2 \in \mathbb{R}^d$, we have

$$\begin{aligned}g(tx_1 + (1-t)x_2) &= \max_{i \in [r]} f_i(tx_1 + (1-t)x_2) \\ &\leq t \cdot \max_{i \in [r]} f_i(x_1) + (1-t) \cdot \max_{i \in [r]} f_i(x_2) \\ &= tg(x_1) + (1-t)g(x_2)\end{aligned}$$

The equality holds if and only if $\operatorname{argmax}_{i \in [r]} f_i(x_1) = \operatorname{argmax}_{i \in [r]} f_i(x_2)$. Therefore g is also a convex function.

Problem 3

Loss History

Training loss is reported as follows (see figure 1):

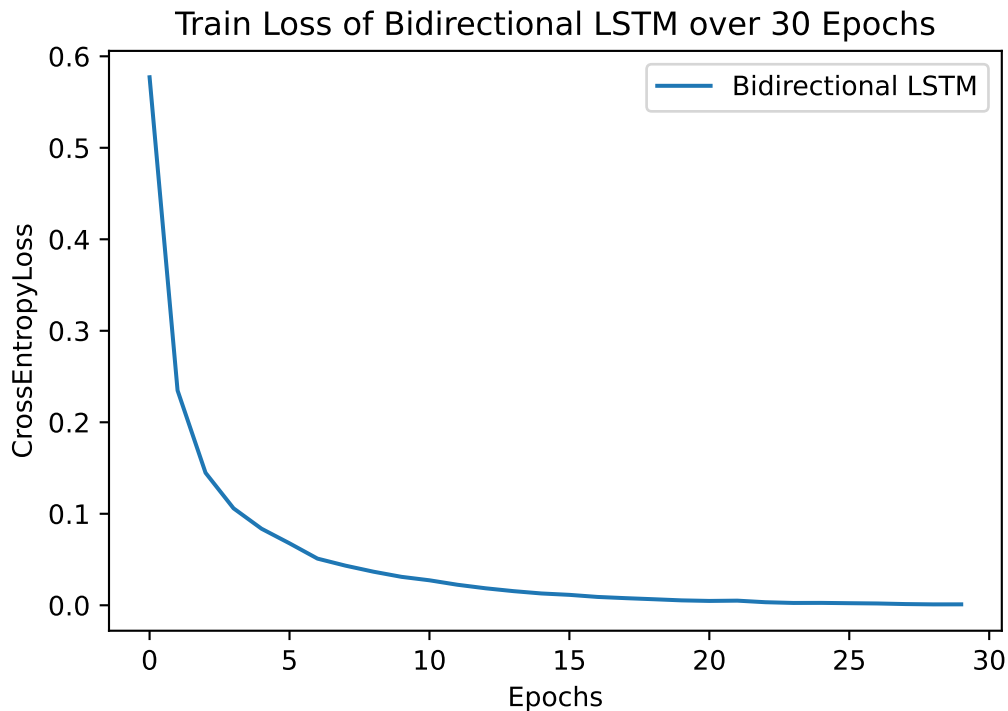


Figure 1: Training loss for 30 epochs

Model structure

```
BidirectionalLSTM(  
    (vocab): Vocab()  
    (word_embeddings): Embedding(13891, 512, padding_idx=1)  
    (lstm): LSTM(512, 256, batch_first=True, bidirectional=True)  
    (lstm2fc): Linear(in_features=512, out_features=256, bias=True)  
    (fc2label): Linear(in_features=256, out_features=2, bias=True)  
)
```

Explanations:

1. **Vocab size** is 13891. I defined 3 special tokens: $\langle pad \rangle$, $\langle bos \rangle$, $\langle eos \rangle$ for padding, begin of sequence, end of sequence.
2. **Embedding dimension** is 512.
3. **LSTM Layer:** $input_size = 512$, $hidden_size = 256$, only one bidirectional layer is adopted.
4. **No dropout layer.**
5. **Fully connected layer:** $input_dim = 512$, $hidden_dim = 256$, $output_dim = 2$ with *relu* activation.

Final Testing Accuracy

Final testing accuracy is 83.83%.