

HW3 Handwritten Assignment Solution

Lecturer: Pei-Yuan Wu

TAs: Yuan-Chia Chang(Problem 1, 2), Chun-Lin Huang(Problem 3, 4, 5)

November 2023, First Edition

Problem 1 (LSTM Cell)(0.5%)

In this exercise, we will simulate the forward pass of a simple LSTM cell. Figure.1 shows a single LSTM cell, where z is the cell input, z_i, z_f, z_o are the control inputs of the gates, c is the cell memory, and f, g, h are activation functions. Given an input x , the cell input and the control inputs can be calculated by the following equations :

- $z = w \cdot x + b$
- $z^i = w_i \cdot x + b_i$
- $z^f = w_f \cdot x + b_f$
- $z^o = w_o \cdot x + b_o$

where w, w_i, w_f, w_o are weights and b, b_i, b_f, b_o are biases. The final output can be calculated by

$$y = f(z^o) h(c')$$

where the value stored in cell memory is updated by

$$c' = f(z^i)g(z) + cf(z^f)$$

Note that $f(z) = \frac{1}{1+e^{-z}}$, $g(z) = z$, $h(z) = z$

Given an input sequence x^t ($t = 1, 2, 3, 4$), please derive the output sequence y_t . The input sequence, the weights, and the activation functions are provided below.

$$\begin{array}{ll} w = [0, 0, 1, 0] & , b = 0 \\ w_i = [50, 50, 0, 0] & , b_i = -5 \\ w_f = [-50, -50, 0, 0], & , b_f = 120 \\ w_o = [0, 0, 200, 0] & , b_o = -30 \\ x^1 = [0, 0, 1, 3] & , x^2 = [0, 1, -1, 2] \\ x^3 = [2, 1, 3, 4] & , x^4 = [0, 1, 0, 0] \end{array}$$

The initial value in cell memory is 0. **Please note that your calculation process is required to receive full credit.**

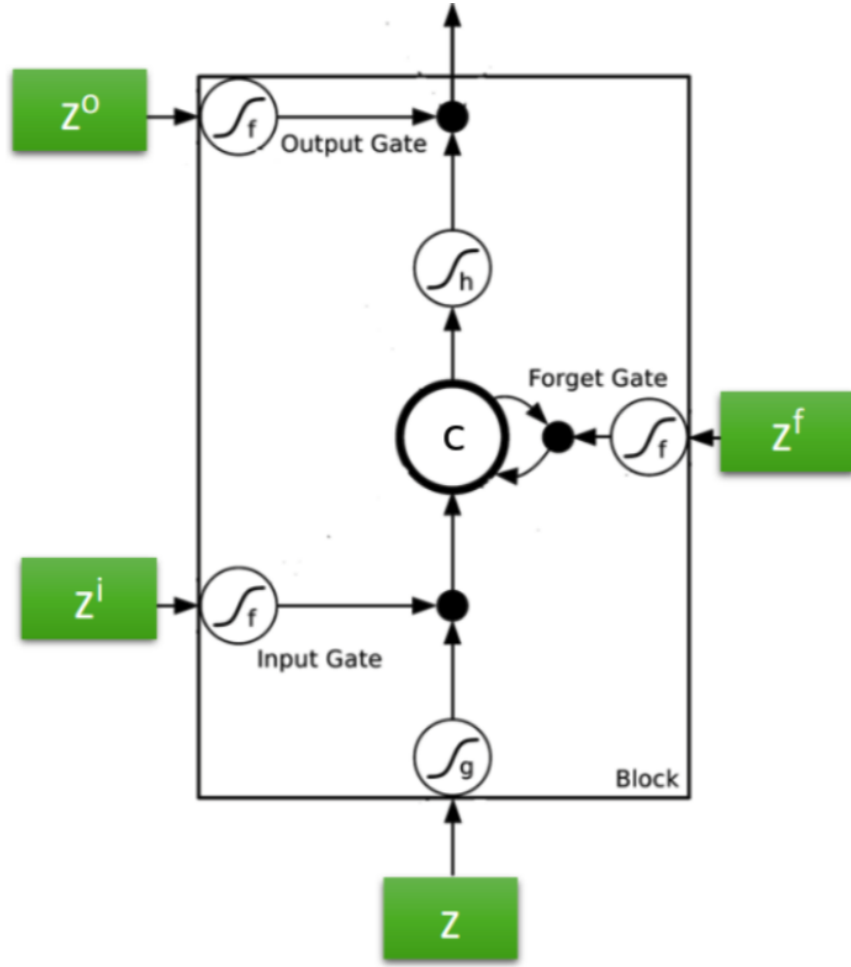


Figure 1: Problem 1 LSTM model

Solution: when $t = 1$, $z_1 = 1$, $z_1^i = -5$, $z_1^f = 120$, $z_1^o = 170$, $c_1 = 0.0067$, $y_1 = 0.0067$.
 $t = 2$, $z_2 = -1$, $z_2^i = 45$, $z_2^f = 70$, $z_2^o = -230$, $c_2 = -1$, $y_2 = 0$.
 $t = 3$, $z_3 = 3$, $z_3^i = 145$, $z_3^f = -30$, $z_3^o = 570$, $c_3 = 3$, $y_3 = 3$.
 $t = 4$, $z_4 = 0$, $z_4^i = 45$, $z_4^f = 70$, $z_4^o = -30$, $c_4 = 0$, $y_4 = 0$.
Hence, $y_1 = 0.0067$, $y_2 = 0$, $y_3 = 0$, $y_4 = 0$.

Problem 2 (Laplacian Eigenmaps)(1.5%)

Consider an undirected connected graph G , which is shown below. We want to utilize Laplacian Eigenmaps method to reduce these 10 points to 3-dimensional space. Here, undirected graph means that edges in the graph do not have a direction, and connected graph means that there is a path from any node to any other node in the graph.

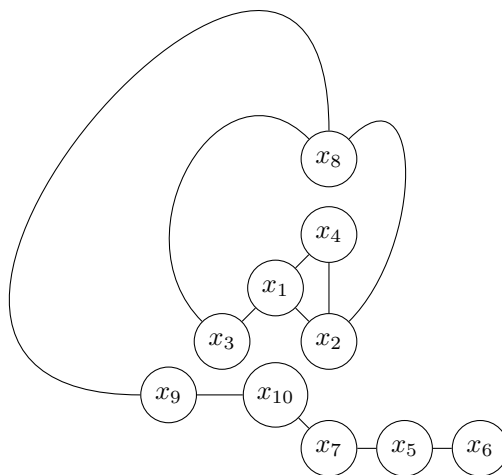


Figure 2: Problem 2 undirected connected graph G

1. Write down the adjacency matrix \mathbf{W}

Solution: \mathbf{W} is shown below.

$$\begin{bmatrix} 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

Figure 3: \mathbf{W}

2. Write down the diagonal matrix $\mathbf{D} = \text{diag}(d_1, \dots, d_{10})$, where $d_i = \sum_{j=1}^{10} \frac{\mathbf{W}_{ij} + \mathbf{W}_{ji}}{2}$ and the Laplacian $\mathbf{L} = \mathbf{D} - \mathbf{W}$.

Solution: \mathbf{D} and \mathbf{L} is shown below.

$$\begin{bmatrix} 3. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 3. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 2. & 0. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 2. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & 2. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 1. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 2. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 0. & 3. & 0. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 2. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 0. & 2. \end{bmatrix}$$

Figure 4: \mathbf{D}

$$\begin{bmatrix} 3. & -1. & -1. & -1. & 0. & 0. & 0. & 0. & 0. & 0. \\ -1. & 3. & 0. & -1. & 0. & 0. & 0. & -1. & 0. & 0. \\ -1. & 0. & 2. & 0. & 0. & 0. & 0. & -1. & 0. & 0. \\ -1. & -1. & 0. & 2. & 0. & 0. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & 2. & -1. & -1. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & -1. & 1. & 0. & 0. & 0. & 0. \\ 0. & 0. & 0. & 0. & -1. & 0. & 2. & 0. & 0. & -1. \\ 0. & -1. & -1. & 0. & 0. & 0. & 0. & 3. & -1. & 0. \\ 0. & 0. & 0. & 0. & 0. & 0. & 0. & -1. & 2. & -1. \end{bmatrix}$$

Figure 5: \mathbf{L}

3. By HW2 Problem 3, Neighbor Embedding Slide p.7-p.10 and programming tools(MATLAB, Python...), solve the optimization problem
 minimize $\text{Trace}(\mathbf{\Psi}^T \mathbf{L} \mathbf{\Psi})$
 subject to $\mathbf{\Psi}^T \mathbf{D} \mathbf{\Psi} = \mathbf{I}_3$
 variables $\mathbf{\Psi} \in \mathbb{R}^{10 \times 3}$
 Also, please plot the reduced points $\mathbf{z}_1, \dots, \mathbf{z}_{10}$ in 3-D scatter plot.

Solution: the scatter plot is shown below.

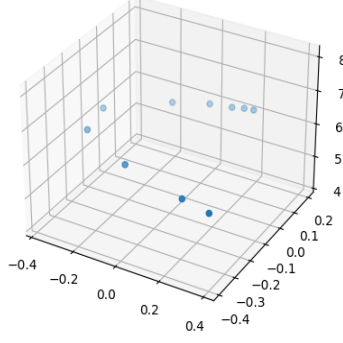


Figure 6: Scatter plot

4. You may find that the minimal eigenvalue of \mathbf{L} is 0, and the corresponding eigenvector is

$$\begin{bmatrix} c \\ c \\ c \\ \vdots \\ c \end{bmatrix} \quad (1)$$

where c is a constant. Since all the points fall into a plane, the span of these points is \mathbb{R}^2 . In order to construct $\mathbf{z}_1, \dots, \mathbf{z}_{10}$ such that $\text{span}\{\mathbf{z}_1, \dots, \mathbf{z}_{10}\} = \mathbb{R}^3$, we need choose the second, third, fourth smallest eigenvalue and the corresponding eigenvectors. Please plot the reduced points by the updated $\mathbf{z}_1, \dots, \mathbf{z}_{10}$ in 3-D scatter plot and verify that whether $\text{Trace}(\Psi^T \mathbf{L} \Psi) = 1.098$ and $\Psi^T \mathbf{D} \Psi = \mathbf{I}_3$.

Solution: Ψ and the scatter plot of adjusted reduced points is shown below. It is true that $\text{Trace}(\Psi^T \mathbf{L} \Psi) = 1.098$ and $\Psi^T \mathbf{D} \Psi = \mathbf{I}_3$.

```
[[ 0.09069099  0.14189576  0.18831754]
 [ 0.09972695  0.09357859  0.17464619]
 [-0.30654733 -0.00046639  0.15764162]
 [ 0.29471481  0.18210673  0.19459486]
 [-0.09353918  0.25247422 -0.3627844 ]
 [-0.28954605  0.39050752 -0.38899742]
 [ 0.22910952 -0.06404401 -0.28767812]
 [-0.28875405 -0.14249883  0.10571999]
 [-0.0730296  -0.36950084 -0.03650001]
 [ 0.24156895 -0.33528676 -0.17380083]]
```

Figure 7: Psi

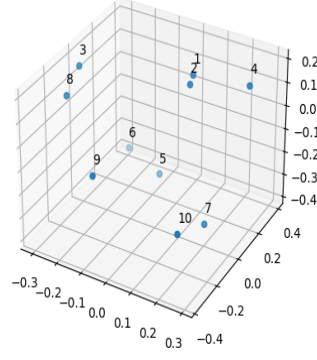


Figure 8: Scatter plot

5. Show that for no matter the graph is, there is an eigenvector of \mathbf{L}

$$\begin{bmatrix} c \\ c \\ \vdots \\ c \end{bmatrix} \quad (2)$$

where c is a constant, and the corresponding eigenvalue is 0.

Solution: Let $\mathbf{v} =$

$$\begin{bmatrix} c \\ c \\ \vdots \\ c \end{bmatrix} \quad (3)$$

$\mathbf{L}\mathbf{v}_i = \sum_{k=1}^N \mathbf{L}_{ik}c = c \sum_{k=1}^N \mathbf{L}_{ik} = c(\sum_{k=1}^N \mathbf{D}_{ik} - \sum_{k=1}^N \mathbf{D}_{ik}) = c(\text{deg of node } i - \text{deg of node } i) = c \times 0 = 0$.
Hence, $\mathbf{L}\mathbf{v} = \mathbf{0} = 0\mathbf{v}$, which indicates that \mathbf{v} is an eigenvector of \mathbf{L} that corresponds to eigenvalue 0.

6. By Neighbor Embedding Slide p.9, please show that

$$\forall \mathbf{f} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_N \end{bmatrix} \in \mathbb{R}^N, \mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{1 \leq i, j \leq N} w_{ij} (f_i - f_j)^2.$$

Solution: Under the knowledge on Neighbor Embedding Slide p.9 that

$$\text{Trace}(\mathbf{\Psi}^T \mathbf{L} \mathbf{\Psi}) = \frac{1}{2} \sum_{i, j=1, \dots, N} w_{ij} \|z_i - z_j\|_2^2$$

. Let $\mathbf{f} = \mathbf{\Psi} \in \mathbb{R}^{N \times 1}$, $\text{Trace}(\mathbf{f}^T \mathbf{L} \mathbf{f}) = \mathbf{f}^T \mathbf{L} \mathbf{f}$, $\sum_{i, j=1, \dots, N} w_{ij} \|z_i - z_j\|_2^2 = \sum_{i, j=1, \dots, N} w_{ij} (f_i - f_j)^2$.
Hence we get $\mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{1 \leq i, j \leq N} w_{ij} (f_i - f_j)^2$.

7. Show that if \mathbf{f} is an eigenvector of \mathbf{L} which corresponds to eigenvalue 0, then $\mathbf{f}^T \mathbf{L} \mathbf{f} = 0$.

Solution: $\mathbf{L} \mathbf{f} = 0 \mathbf{f} = \mathbf{0}$, hence $\mathbf{f}^T \mathbf{L} \mathbf{f} = 0$.

8. Show that if the graph is connected, the second smallest eigenvalue of \mathbf{L} will be nonzero.

Solution: By 6, we know that \mathbf{L} is positive semidefinite matrix, which indicates that all the eigenvalues of \mathbf{L} is non-negative. From the above result and 5. we deduced that 0 is the smallest eigenvalue, which corresponds to eigenvector $\mathbf{v}_1 =$

$$\begin{bmatrix} c \\ c \\ \vdots \\ c \end{bmatrix} \quad (4)$$

Also, because \mathbf{L} is a symmetric matrix, it is diagonalizable, which means that \mathbf{L} has N independent eigenvectors. Suppose \mathbf{v}_2 is an eigenvector corresponding to the second smallest eigenvalue λ_2 . From the above facts, \mathbf{v}_2 can not be the multiply of \mathbf{v}_1 , which means \mathbf{v}_2 can not be the form

$$\begin{bmatrix} d \\ d \\ \vdots \\ d \end{bmatrix} \quad (5)$$

On the other hand, because G is a connected graph, we can find a path from node i to node j . We define the nodes on path be $n_1 = i, n_2, \dots, n_{k-1}, n_k = j$. $\sum_{1 \leq i, j \leq N} w_{ij}(v_{2i} - v_{2j})^2 \geq \sum_{1 \leq r \leq k-1} (v_{2n_r} - v_{2n_{r+1}})^2 > 0$ due to the fact that v_2 is not the multiply of all one vector. Also by 6., we get $\mathbf{v}_2^T \mathbf{L} \mathbf{v}_2 = \mathbf{v}_2^T \lambda_2 \mathbf{v}_2 = \lambda_2 \|\mathbf{v}_2\|_2 = \sum_{1 \leq i, j \leq N} w_{ij}(v_{2i} - v_{2j})^2 > 0$, where λ_2 is the second smallest eigenvalue. Hence by the above inequality, we show that the second smallest eigenvalue of \mathbf{L} will be nonzero.

Problem 3 (Multiclass AdaBoost)(1.5%)

Let \mathcal{X} be the input space, \mathcal{F} be a collection of multiclass classifiers that map from \mathcal{X} to $[1, K]$, where K denotes the number of classes. Let $\{(x_i, \hat{y}_i)\}_{i=1}^m$ be the training data set, where $x_i \in \mathcal{X}$ and $\hat{y}_i \in [1, K]$. Given $T \in \mathbb{N}$, suppose we want to find functions

$$g_{T+1}^k(x) = \sum_{t=1}^T \alpha_t f_t^k(x), \quad k \in [1, K]$$

where $f_t \in \mathcal{F}$ and $\alpha_t \in \mathbb{R}$ for all $t \in [1, T]$. Here for $f \in \mathcal{F}$, we denote $f^k(x) = \mathbf{1}\{f(x) = k\}$, where $\mathbf{1}(\cdot)$ is an indicator function, as the k 'th element in the one-hot representation of $f(x) \in [1, K]$. The aggregated classifier $h : \mathcal{X} \rightarrow [1, K]$ is defined as

$$x \mapsto \operatorname{argmax}_{1 \leq k \leq K} g_{T+1}^k(x)$$

Please apply gradient boosting to show how the functions f_t and coefficients α_t are computed with an aim to minimize the following loss function

$$L((g_{T+1}^1, \dots, g_{T+1}^K)) = \sum_{i=1}^m \exp \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{T+1}^k(x_i) - g_{T+1}^{\hat{y}_i}(x_i) \right)$$

Solution: Given $\mathbf{g}_{t-1} = \{g_{t-1}^k\}_{k=1}^K$, we update $\mathbf{g}_t = \{g_{t-1}^k + \alpha_t f_t^k\}_{k=1}^K = \mathbf{g}_{t-1} + \alpha_t \mathbf{f}_t$ as follows:

$$\begin{aligned}
\mathbf{f}_t &\in \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} \frac{\partial}{\partial \alpha} L(\mathbf{g}_{t-1} + \alpha \mathbf{f}) \Big|_{\alpha=0} \\
&= \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} \frac{\partial}{\partial \alpha} \sum_{i=1}^n \exp \left(\left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right) + \alpha \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i) \right) \right) \Big|_{\alpha=0} \\
&= \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} \sum_{i=1}^n \exp \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right) \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i) \right) \\
&= \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} Z_t \mathbb{E}_{i \sim D_t} \left[\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f^k(x_i) - f^{\hat{y}_i}(x_i) \right] \\
&= \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} Z_t \mathbb{E}_{i \sim D_t} \left[\frac{1}{K-1} \cdot 1\{f(x_i) \neq \hat{y}_i\} - 1\{f(x_i) = \hat{y}_i\} \right] \\
&= \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} Z_t \left(\frac{K}{K-1} \mathbb{P}_{i \sim D_t} [f(x_i) \neq \hat{y}_i] - 1 \right) = \operatorname{argmin}_{\mathbf{f} \in \mathcal{F}} \mathbb{P}_{i \sim D_t} [f(x_i) \neq \hat{y}_i] \\
\alpha_t &\in \operatorname{argmin}_{\alpha \in \mathbb{R}} L(\mathbf{g}_{t-1} + \alpha \mathbf{f}_t) \\
&= \operatorname{argmin}_{\alpha \in \mathbb{R}} \sum_{i=1}^n \exp \left(\left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right) + \alpha \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f_t^k(x_i) - f_t^{\hat{y}_i}(x_i) \right) \right) \\
&= \operatorname{argmin}_{\alpha \in \mathbb{R}} Z_t \mathbb{E}_{i \sim D_t} \left[e^{\alpha \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} f_t^k(x_i) - f_t^{\hat{y}_i}(x_i) \right)} \right] \\
&= \operatorname{argmin}_{\alpha \in \mathbb{R}} Z_t \mathbb{E}_{i \sim D_t} \left[e^{\frac{\alpha}{K-1}} \cdot 1\{f_t(x_i) \neq \hat{y}_i\} + e^{-\alpha} \cdot 1\{f_t(x_i) = \hat{y}_i\} \right] \\
&= \operatorname{argmin}_{\alpha \in \mathbb{R}} Z_t (\epsilon_t e^{\frac{\alpha}{K-1}} + e^{-\alpha} (1 - \epsilon_t)) = \left\{ \frac{K-1}{K} \log \frac{(K-1)(1 - \epsilon_t)}{\epsilon_t} \right\}
\end{aligned}$$

where

$$Z_t = \sum_{i=1}^n \exp \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right)$$

and that D_t is a probability distribution for $t = 1, \dots, n$ given by

$$D_t(i) = \frac{1}{Z_t} \exp \left(\frac{1}{K-1} \sum_{k \neq \hat{y}_i} g_{t-1}^k(x_i) - g_{t-1}^{\hat{y}_i}(x_i) \right)$$

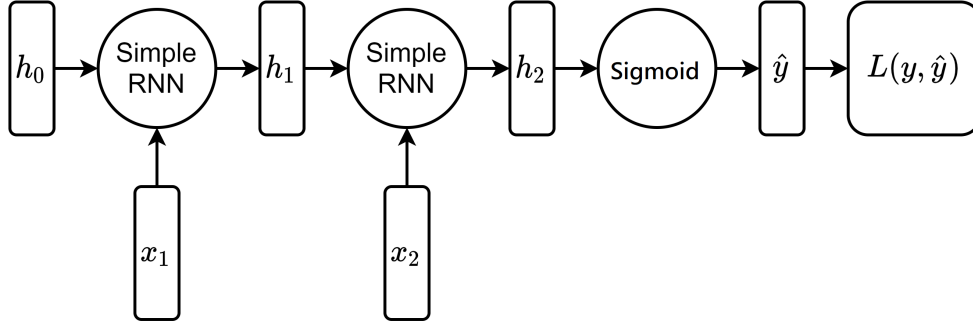
and that $\epsilon_t = \mathbb{P}_{i \sim D_t} [f_t(x_i) \neq \hat{y}_i]$ is the error of f_t on training sample weighted by the distribution D_t .

Problem 4 (Backpropagation through time via Simple RNN)(1%)

Backpropagation through time is a critical concept to know as we train a recurrent network. Here, we set a toy case of prediction problem. The Simple RNN module has two kinds of weights, w_x and w_h , such that $h_t = \tanh(w_x x_t + w_h h_{t-1})$, where t represents the index of steps. The module has the weight w_o such that $\hat{y} = \sigma(w_o h_2)$, where $\sigma(w_o h_2) = \frac{1}{1 + \exp(-w_o h_2)}$. The initial state h_0 is set to be 0. The sequential input only contains $\{x_1, x_2\}$; the label is y ; the loss function is MSE. Please derive $\frac{\partial L(y, \hat{y})}{\partial w_o}$, $\frac{\partial L(y, \hat{y})}{\partial w_h}$, $\frac{\partial L(y, \hat{y})}{\partial w_x}$ in terms of $x_1, x_2, h_0, h_1, h_2, w_x, w_o$, and w_h .

Solution:

$$\frac{\partial L(y, \hat{y})}{\partial w_o} = \frac{\partial \hat{y}^2}{\partial w_o} - 2y \frac{\partial \hat{y}}{\partial w_o} = 2(\hat{y} - y) \frac{\exp(-w_o h_2) h_2}{(1 + \exp(-w_o h_2))^2}$$



For the rest of the part, we first calculate two recursive relation.

$$\frac{\partial h_t}{\partial w_x} = \frac{\partial \tanh(w_x x_t + w_h h_{t-1})}{\partial (w_x x_t + w_h h_{t-1})} \frac{\partial (w_x x_t + w_h h_{t-1})}{\partial w_x} = \text{sech}^2(w_x x_t + w_h h_{t-1}) x_t + \text{sech}^2(w_x x_t + w_h h_{t-1}) w_h \frac{\partial h_{t-1}}{\partial w_x}$$

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial \tanh(w_x x_t + w_h h_{t-1})}{\partial (w_x x_t + w_h h_{t-1})} \frac{\partial (w_x x_t + w_h h_{t-1})}{\partial w_h} = \text{sech}^2(w_x x_t + w_h h_{t-1}) h_{t-1} + \text{sech}^2(w_x x_t + w_h h_{t-1}) \frac{\partial h_{t-1}}{\partial w_h}$$

Actually if you can identify this two relation I will give you full credit. The followings are the rest of the part.

$$\frac{\partial L(y, \hat{y})}{\partial w_h} = (\hat{y} - y) \frac{\exp(-w_o h_2) w_o}{(1 + \exp(-w_o h_2))^2} \frac{\partial h_2}{\partial w_h}$$

$$\frac{\partial L(y, \hat{y})}{\partial w_h} = (\hat{y} - y) \frac{\exp(-w_o h_2) w_o}{(1 + \exp(-w_o h_2))^2} \frac{\partial h_2}{\partial w_x}$$

Problem 5 (Loss function of Decision tree) (1.5%)

It is known that decision tree is still a powerful classification model now a days. There are two different loss functions when it comes to entropy counting, which are Shannon information gain and Gini index. Following are their definition:

$$\text{Gini index} = \frac{N_{left}}{N} \left(1 - \sum_{i=1}^c (p_{left}^i)^2 \right) + \frac{N_{right}}{N} \left(1 - \sum_{i=1}^c (p_{right}^i)^2 \right)$$

$$\text{Shannon information gain} = \frac{N_{left}}{N} \left(- \sum_{i=1}^c p_{left}^i \log_2 p_{left}^i \right) + \frac{N_{right}}{N} \left(- \sum_{i=1}^c p_{right}^i \log_2 p_{right}^i \right)$$

p_i^j := the proportion of class j in the node i .

N_i := the number of cases in the node i .

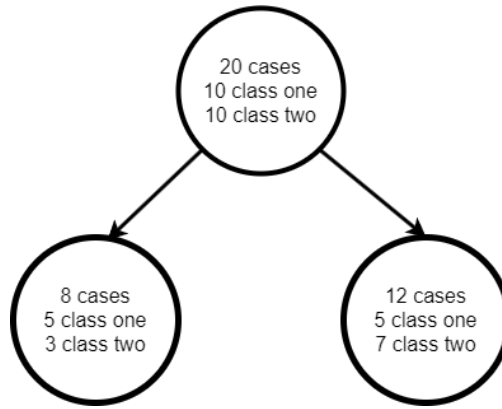
Now we give a toy example. In this case $N_{left} = 8, p_{left}^1 = \frac{5}{8}, p_{left}^2 = \frac{3}{8}, N_{right} = 12, p_{right}^1 = \frac{5}{12}, p_{right}^2 = \frac{7}{12}$

In the following questions we consider classification of two cases. Please calculate the entropy of the following questions using two above-mentioned loss functions.

- (a) (i) A 50/50 split with the first part containing 80% of positive examples and the second part containing 75% of positive examples.
- (ii) A 80/20 split with the first part containing 0% of positive examples and the second part containing 90% of positive examples.
- (iii) A 90/10 split with the first part containing 1% of positive examples and the second part containing 100% of positive examples.

Solution:

- (i) Gini index = 0.35, Shannon information gain = 0.77



(ii) Gini index = 0.036, Shannon information gain = 0.094

(iii) Gini index = 0.018, Shannon information gain = 0.073

- (b) However, now suppose that our case is to detect the covid-19. Thus, we want our entropy function can have a higher loss on (iii) than (ii). Please decide a function that fulfill this criteria and write down the loss.

Solution: The most easy way is adding a weight in front of different classes. You, however, can design your own loss function to satisfy this criteria. It is better if you can make sure that it keeps the same order with the origin when we fix the split or fix the proportion in two nodes.

Version Description

1. First Edition: Finish Problem Solution 1 to 5