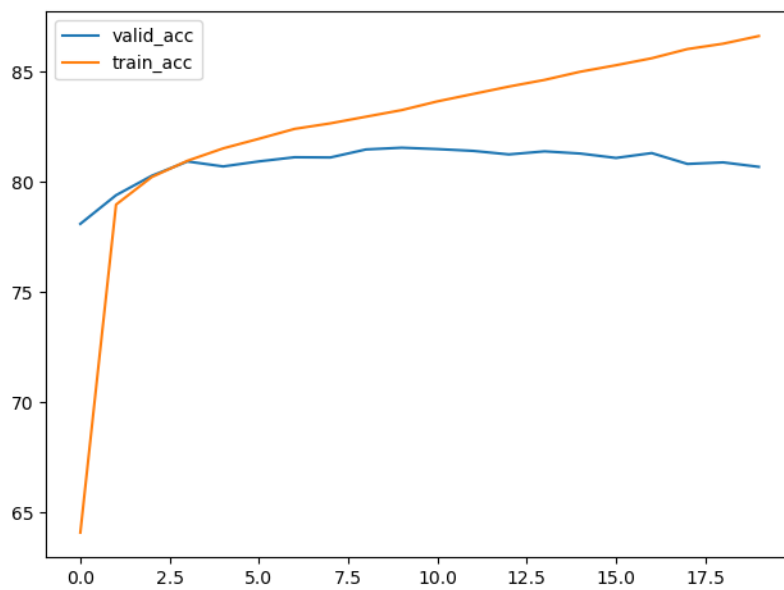
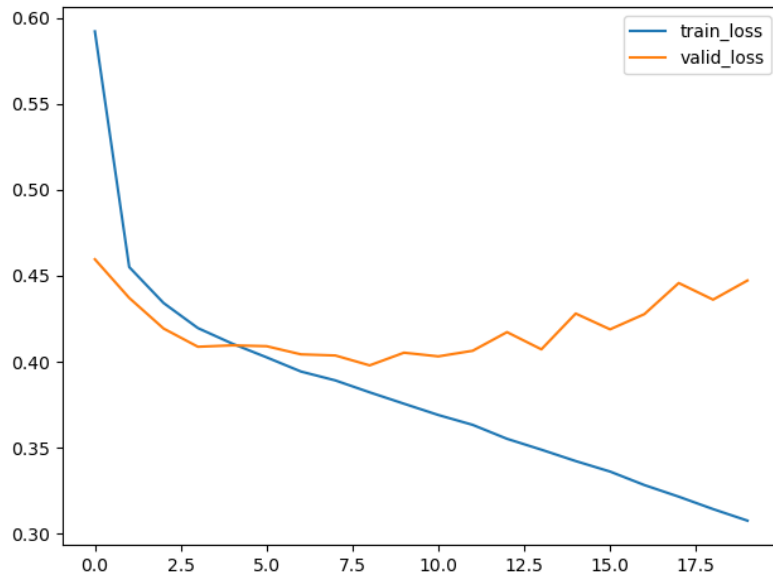


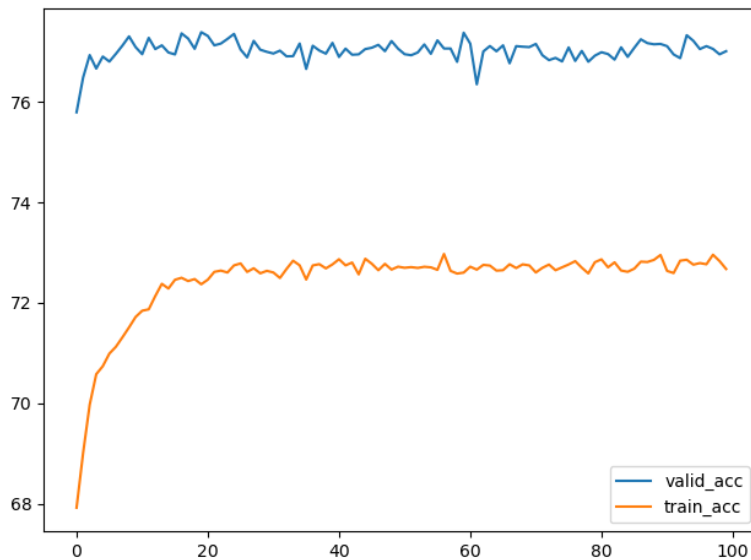
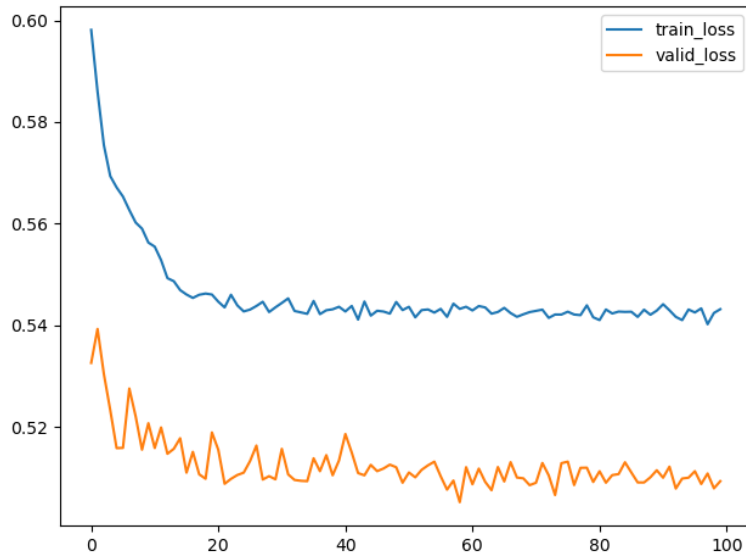
1. (0.5%) 請說明你實作之 RNN 模型架構及使用的 word embedding 方法，回報模型的正確率並繪出訓練曲線\*



模型：1層 Word2vec Embedding (dim 30)  $\rightarrow$  4層 LSTM  $\rightarrow$  Dropout (0.75)  $\rightarrow$  1層線性 (dim 64)  $\rightarrow$  Sigmoid

正確率：Train 0.83266, Validation 0.81559, Kaggle 0.82560

2. (0.5%) 請實作 BOW+DNN 模型，敘述你的模型架構，回報模型的正確率並繪出訓練曲線\*。



模型：4層線性 (dim 128)，每層間間隔dropout 0.5，再接sigmoid

正確率：Train 0.72369, Validation 0.77386, Kaggle 0.77180

這裡發現一個現象，validation的結果比較training還好。

3. (0.5%) 請敘述你如何 improve performance (preprocess, embedding, 架構等)，並解釋為何這些做法可以使模型進步。

i. 預處理

英文有大量的縮寫字，在預處理時先展開這些字，如[can't] → [can, not]，而不是原本的[can, ', t]，因為撇號通常不帶有語意，把它當symbol不太合理。

## ii. 改變word2vec的參數

window size調高到10，dimension調高到300，這樣word可以看見更多的資訊，embedding有機會更準確。

另外，使用unlabeled data一起作word embedding，讓字之間距離的資訊更多。

## iii. 改變rnn的參數與架構

lstm用了4層，每層dimension是64，並加上dropout 0.75和L2 regularization (係數1e-7)，以減少overfitting。效果最明顯的是L2 regularization。

句子長度取45，可以涵蓋data中大部分的句子。

## iv. 使用gru取代lstm

因為gru的參數較少，比較不會overfitting。但在這個dataset似乎沒有明顯進步。

## v. 使用Bidirectional RNN

這樣可以讓training過程中看見之前和之後的word。

4. (0.5%) 請比較 RNN 與 BOW 兩種不同 model 對於 "Today is hot, but I am happy" 與 "I am happy, but today is hot" 這兩句話的分數 (model output)，並討論造成差異的原因。

	RNN	BOW
Today is hot, but I am happy:	0.8716	0.1350
I am happy, but today is hot:	0.1251	0.1350

這兩句話出現的字完全一樣，只有順序不一樣。bow只考慮出現的字，而rnn會參考往前數個字(以lstm為例)，因此rnn可以分辨因語序不同造成的語意差異。

## 5. (3%) Math problem:

<https://drive.google.com/file/d/1fEu87banB4s6Yjku1dA5sMcnwCugEPBF/view?usp=sharing>

ML HW4

1. Let  $X^{\text{raw}} = \begin{pmatrix} 1 & 4 & 3 & 1 & 5 & 7 & 9 & 3 & 11 & 10 \\ 2 & 8 & 12 & 8 & 14 & 4 & 8 & 8 & 5 & 11 \\ 3 & 5 & 9 & 5 & 2 & 1 & 9 & 1 & 6 & 7 \end{pmatrix} = (x_1^{\text{raw}} \dots x_{10}^{\text{raw}})$

Mean  $\mu = \begin{pmatrix} 5.4 \\ 8 \\ 4.8 \end{pmatrix} = \frac{1}{10} \sum_{i=1}^{10} x_i^{\text{raw}}$

$X := X^{\text{raw}} - \mu = \begin{pmatrix} x_1^{\text{raw}} - \mu & \dots & x_{10}^{\text{raw}} - \mu \end{pmatrix} = (x_1 \dots x_{10}) = \begin{pmatrix} -4.4 & \dots \\ -6 & \dots \\ -1.8 & \dots \end{pmatrix}$

Do singular-value decomposition:

$X = UDV^T$

$\Sigma = \frac{1}{10} \sum_{i=1}^{10} x_i x_i^T = \frac{1}{10} X X^T = \frac{1}{10} U D V^T V D U^T = U \left( \frac{1}{10} D D^T \right) U^T$

Where  $U = (u_1 \ u_2 \ u_3) = \begin{pmatrix} -0.62 & 0.68 & -0.40 \\ -0.59 & -0.73 & -0.34 \\ -0.52 & -0.03 & 0.85 \end{pmatrix}$

$D = \begin{pmatrix} 12.37 & 0 & 0 & 0 & 0 \\ 0 & 10.78 & 0 & 0 & 0 \\ 0 & 0 & 7.40 & 0 & 0 \end{pmatrix}_{3 \times 10}$

(a) Principle Axes are:  $u_1 = \begin{pmatrix} -0.62 \\ -0.59 \\ -0.52 \end{pmatrix}$ ,  $u_2 = \begin{pmatrix} 0.68 \\ -0.73 \\ -0.03 \end{pmatrix}$ ,  $u_3 = \begin{pmatrix} -0.40 \\ -0.34 \\ 0.85 \end{pmatrix}$   
(Approximately)

(b)  $U^T X^{\text{raw}} = \begin{pmatrix} \boxed{-3.36} & \boxed{-9.79} & \boxed{-13.62} & -7.94 & -12.37 & -7.19 & -14.96 & -7.08 & -12.86 & \boxed{-16.3} \\ \boxed{-0.71} & \boxed{-3.03} & \boxed{-6.53} & -5.06 & -6.84 & 1.84 & 0.47 & -3.81 & 3.95 & \boxed{-1.11} \\ \boxed{1.48} & \boxed{-0.64} & \boxed{2.42} & 1.16 & -5.02 & -3.30 & 1.37 & -3.05 & -0.97 & \boxed{-1.75} \end{pmatrix}$   
pc of  $x_1^{\text{raw}}$  ... pc of  $x_{10}^{\text{raw}}$

(c) Use  $u_1$  and  $u_2$  since they have largest eigenvalue

$\Rightarrow$  Principle Component (pc)  $= \begin{pmatrix} u_1^T \\ u_2^T \end{pmatrix} X^{\text{raw}}$

Reconstruction Error  $= \left\| \begin{pmatrix} u_1 & u_2 \end{pmatrix} \begin{pmatrix} u_1^T \\ u_2^T \end{pmatrix} X^{\text{raw}} - X^{\text{raw}} \right\|^2 = 60.64$

# ML HW4

2. (a)  $A \in \mathbb{R}^{m \times n}$   
Symmetric:  $(AA^T)^T = AA^T$

$$(A^T A)^T = A^T A$$

Positive Semi-definite:

$$\forall x \quad x^T AA^T x = \|A^T x\|^2 \geq 0$$

$$\forall y \quad y^T A^T A y = \|A y\|^2 \geq 0$$

Sharing same eigenvalues:

For matrix  $A^T A$ , since it's positive semidefinite, its eigenvalues are non-negative. Let positive eigenvalues be  $\lambda_1, \lambda_2, \dots, \lambda_k$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k > 0$ ,  $k \leq n$ , and corresponding eigenvectors be  $v_1, \dots, v_k$

$$\text{So } A^T A v_i = \lambda_i v_i, \quad i \leq k$$

$$\text{Define } \sigma_i = \sqrt{\lambda_i}, \quad u_i = \frac{1}{\sigma_i} A v_i, \quad \text{for } i \leq k$$

$$\Rightarrow A^T u_i = A^T \frac{1}{\sigma_i} A v_i = \frac{1}{\sigma_i} A^T A v_i = \frac{1}{\sigma_i} \lambda_i v_i = \sigma_i v_i$$

$$\Rightarrow AA^T u_i = \sigma_i A v_i = \sigma_i^2 u_i$$

$\Rightarrow \sigma_i^2 = \lambda_i$  is also an eigenvalue of  $AA^T$

$\Rightarrow$  Let  $B = A^T$ , we <sup>can</sup> show the inverse direction (positive eigenvalues of  $AA^T$  are eigenvalues of  $A^T A$ )

$\Rightarrow A^T A$  and  $AA^T$  share same non-zero eigenvalues

2. (b)  $\because \Sigma$  is symmetric and positive definite

$\Rightarrow \Sigma$  can be diagonalize as:

$$\Sigma = U \Lambda U^T \quad \text{where } U \in \mathbb{R}^{m \times m} \text{ is orthonormal}$$

$$\Lambda \in \mathbb{R}^{m \times m} \text{ is diagonal}$$

$$\Lambda = \text{diag}\left(\frac{\sigma_1^2}{n}, \frac{\sigma_2^2}{n}, \dots, \frac{\sigma_m^2}{n}\right) \quad \text{for some } \begin{cases} \sigma_1, \dots, \sigma_m, \sigma_i \geq 0 \\ n, n > m \end{cases}$$

$$\Rightarrow \Sigma = U \left(\frac{1}{n} D D^T\right) U^T \quad \text{where } D \in \mathbb{R}^{n \times n}, D = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_m \\ 0 & 0 & \dots & 0 \end{pmatrix}$$

$$= \frac{1}{n} U \underbrace{D V^T V D^T}_{= I_n} U^T \quad \text{for some orthonormal } V \in \mathbb{R}^{n \times n}$$

$$= \frac{1}{n} (U D V^T) (U D V^T)^T$$

$$= \frac{1}{n} X X^T$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{x}_i \hat{x}_i^T$$

$$= \frac{1}{n} \sum_{i=1}^n (\hat{x}_i - \mu)(\hat{x}_i - \mu)^T \quad \leftarrow \text{Let } \tilde{x}_i = \hat{x}_i - \mu$$

There always exists  $V$  such that  $\frac{1}{n} \sum_{i=1}^n \hat{x}_i = 0$  since:

$$\frac{1}{n} \sum_{i=1}^n \hat{x}_i = 0 \Leftrightarrow X \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = 0 \Leftrightarrow U D V^T \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = 0$$

$$\Leftrightarrow V^T \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \text{Null}(UD)$$

$$\text{Since } \dim(\text{Null}(UD)) = \text{Nullity}(UD) = n - \underbrace{\text{Rank}(UD)}_{\leq m} \geq n - m > 0$$

$$\Rightarrow \exists y = V^T \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\text{So } \{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n\} \text{ satisfies } \frac{1}{n} \sum_{i=1}^n \tilde{x}_i = \frac{1}{n} \sum_{i=1}^n (\hat{x}_i - \mu) = 0$$

$$\frac{1}{n} \sum_{i=1}^n (\tilde{x}_i - \mu)(\tilde{x}_i - \mu)^T = \Sigma$$

2.(c)  $1 \leq k \leq m$

$$\begin{cases} \min: \text{Tr}(\Phi^T \Sigma \Phi) \\ \text{subject to } \Phi^T \Phi = I_k \\ \Phi \in \mathbb{R}^{m \times k} \end{cases}$$

$$\text{Trace}(\Phi^T \Sigma \Phi)$$

$$= \frac{1}{n} \text{Trace}(\Phi^T X X^T \Phi) \leftarrow \text{By 2.(b), } \Sigma \text{ is symmetric and positive semi-definite, so } \Sigma = \frac{1}{n} X X^T \exists X \in \mathbb{R}^{m \times n}$$

$$= \frac{1}{n} \|\Phi^T X\|_F^2$$

$$= \frac{1}{n} \sum_{i=1}^n \|\Phi^T x_i\|^2$$

$$= \frac{1}{n} \sum_{i=1}^n \|\hat{x}_i^{(S)}\|^2$$

By  $\Phi^T \Phi = I_k$ ,  $\Phi^T$  is a projection matrix onto  $k$ -dimensional space of  $\mathbb{R}^m$  (let's say,  $S$ )  
 $\hat{x}_i^{(S)}$  denotes the projection of  $x_i$  onto  $S$

$$\text{Let } \hat{X}^{(S)} = (\hat{x}_1^{(S)} \dots \hat{x}_n^{(S)})$$

$$\hat{X}^{(P)} = (\hat{x}_1^{(P)} \dots \hat{x}_n^{(P)}) = \sum_{i=1}^k \sigma_i u_i v_i^T \quad \text{where } \sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_m \text{ are singular values of } X$$

By Eckart-Young-Minsky thm.

$$\|X - \hat{X}^{(P)}\|_F \geq \|X - \hat{X}^{(S)}\|_F$$

$$\Rightarrow \sum_{i=1}^n \|x_i - \hat{x}_i^{(P)}\|^2 \geq \sum_{i=1}^n \|x_i - \hat{x}_i^{(S)}\|^2$$

$$\Rightarrow \sum_{i=1}^n \|\hat{x}_i^{(P)}\|^2 \leq \sum_{i=1}^n \|\hat{x}_i^{(S)}\|^2 = \text{Trace}(\Phi^T \Sigma \Phi)$$

$$\text{So } \text{Trace}(\Phi^T \Sigma \Phi) \geq \sum_{i=1}^n \|\hat{x}_i^{(P)}\|^2 = \sum_{i=1}^n \left\| \sum_{j=1}^k \sigma_j u_j v_j^T \right\|^2$$

$$= \|\hat{X}^{(P)}\|_F^2 = \left\| \sum_{i=1}^k \sigma_i u_i v_i^T \right\|_F^2 = \sum_{i=1}^k \sigma_i^2 \|u_i v_i^T\|_F^2$$

3. Define: True label  $\hat{y}_i = \begin{pmatrix} -\frac{1}{K-1} \\ \vdots \\ 1 \\ \vdots \\ -\frac{1}{K-1} \end{pmatrix}$  where  $\hat{y}_i^{(p)} = \begin{cases} 1 & \text{if this data is class } p \\ -\frac{1}{K-1} & \text{otherwise} \end{cases}$

Function  $g_t(x) = \begin{pmatrix} g_t^1(x) \\ \vdots \\ g_t^K(x) \end{pmatrix}$ ,  $\alpha_t = \begin{pmatrix} \alpha_t^1 \\ \vdots \\ \alpha_t^K \end{pmatrix}$

The loss function becomes:  $L(g_t) = \sum_{i=1}^n \exp\{-\hat{y}_i^T g_t(x_i)\}$

$$= \sum_{i=1}^n \exp\{-\hat{y}_i^T (g_{t-1}(x_i) + f_t(x_i) \alpha_t)\} = \sum_{i=1}^n w_i \exp\{-f_t(x_i) \hat{y}_i^T \alpha_t\}$$

(Where  $w_i = \exp\{-\hat{y}_i^T g_{t-1}(x_i)\}$ )

$$\Rightarrow \nabla_{\alpha_t} L = \sum_{i=1}^n w_i \exp\{-f_t(x_i) \hat{y}_i^T \alpha_t\} (-f_t(x_i) \hat{y}_i) = 0$$

$$\Rightarrow \alpha_t = \frac{(K-1)^2}{K} \left[ \log\left(\frac{1-\epsilon_t}{\epsilon_t}\right) + \log(K-1) \right]$$

$$\text{where } \epsilon_t = \frac{\sum_{\hat{y}_i \neq f_t(x_i)} w_i}{\sum_{i=1}^n w_i}$$