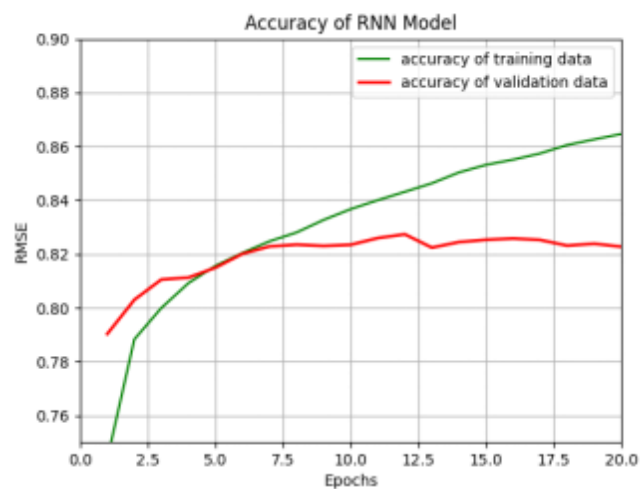
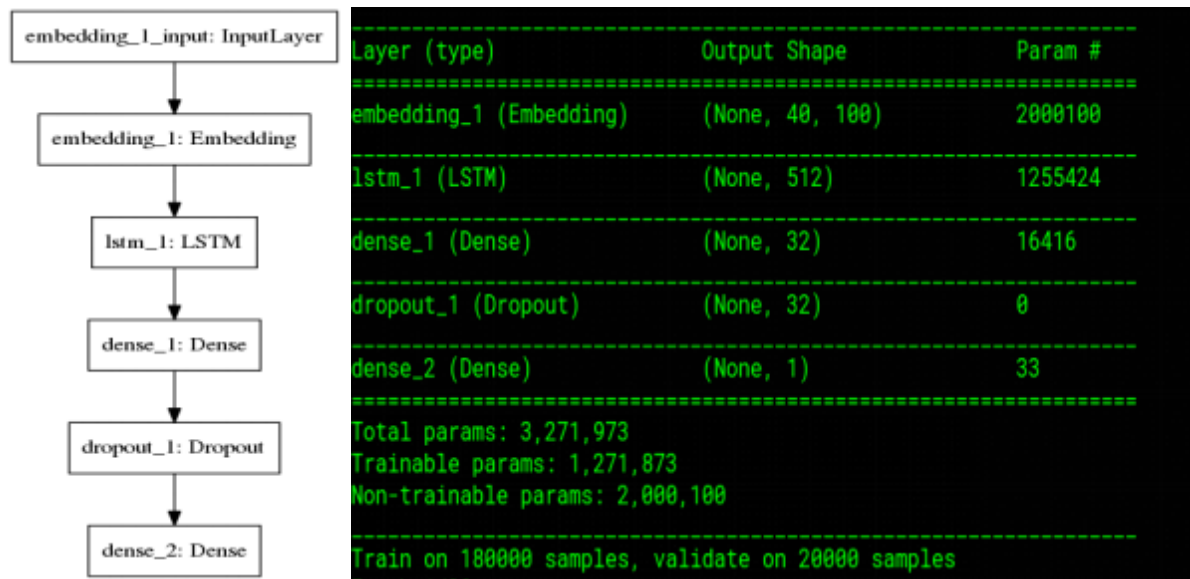


學號：R05922145 系級：資工碩二 姓名：郁錦濤

1. (1%) 請說明你實作的 RNN model，其模型架構、訓練過程和準確率為何？
(Collaborators:)

答：



我的model採用了一層LSTM，然後dense和dropout，train的過程中將embedding也加入一起train，epoch為20，optimizer為adam，adam的參數為默認值lr=1e-3，loss function採用binary_crossentropy。最後在private set上的準確率為0.82122。

2. (1%) 請說明你實作的 BOW model，其模型架構、訓練過程和準確率為何？
(Collaborators:)

答：

3. (1%) 請比較bag of word與RNN兩種不同model對於"today is a good day, but it is hot"與"today is hot, but it is a good day"這兩句的情緒分數，並討論造成差異的原因。
(Collaborators:)

答：

4. (1%) 請比較"有無"包含標點符號兩種不同tokenize的方式，並討論兩者對準確率的影響。
(Collaborators:)

答：

我的model中取tokenize的方式是用keras中的Tokenizer。通過對Tokenizer的參數filter進行修改來實現有無包含標點符號。通過比較，發現含有標點符號的tokenize方式在public set和private set上的performance跟好。仔細分析，我覺得標點符號對於人的情感是具有一定影響的，一個句子以感嘆號結尾和以問號結尾是不同，例如“我很高興。”和“我很高興？”，從字面意義上來看，這兩句話表達的情感是截然不同的，但是如果不包含標點符號的話，就會認為兩句話是一致的，這樣就會產生誤差。

punctions	public	private
yes	0.82355	0.82122
no	0.82031	0.81785

5. (1%) 請描述在你的semi-supervised方法是如何標記label，並比較有無semi-surpervised training對準確率的影響。
(Collaborators:)

答：

我的semi-supervised方法採用了助教的方式，首先用含有label的data進行train，得到用label的data訓練出來的模型，然後用不含label的data通過得到的model進行預測。得到的預測結果，選取合適的threshold值（我取得是0.3和0.7），把機率小於0.2和

大於0.8的data放入training data中再進行fit操作。

按照经验，semi-supervised training会增加training的资料量，對於準確率會有一定程度的提升，但是在我的model中，進行了semi-supervised後準確度卻沒有提升，反而準確率有所降低。究其原因，我認為threshold的取值很重要，對於模稜兩可的data只有選用更高的threshold才能更好的进行分類。另外，數據的相關性也很重要。

