

Homework4 Report

Professor Pei-Yuan Wu
EE5184 - Machine Learning

姓名：顏宏宇

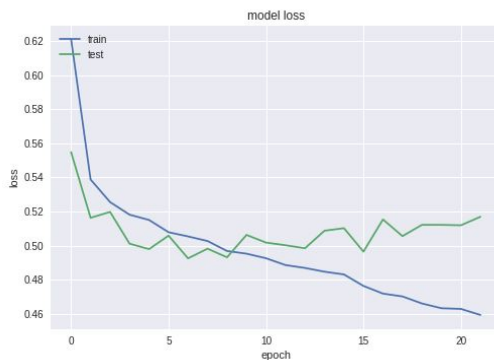
學號：r07942086

Problem 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。(0.5%) 請實作 BOW+DNN 模型,敘述你的模型架構,回報正確率並繪出訓練曲線。

RNN: public score為0.76165

單層LSTM後接Dense layer最後接softmax

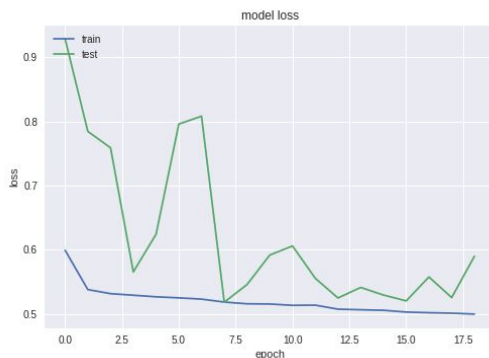
word embedding使用jieba斷字、用word2vec轉換成vector、取前120個word轉成vector當成input



BOW+DNN: public score為0.73827

純Dense layer最後接softmax

word embedding使用jieba斷字、用word2vec轉換成vector、把前120個word vector全部相加當成input, 概念上應該跟BOW一樣

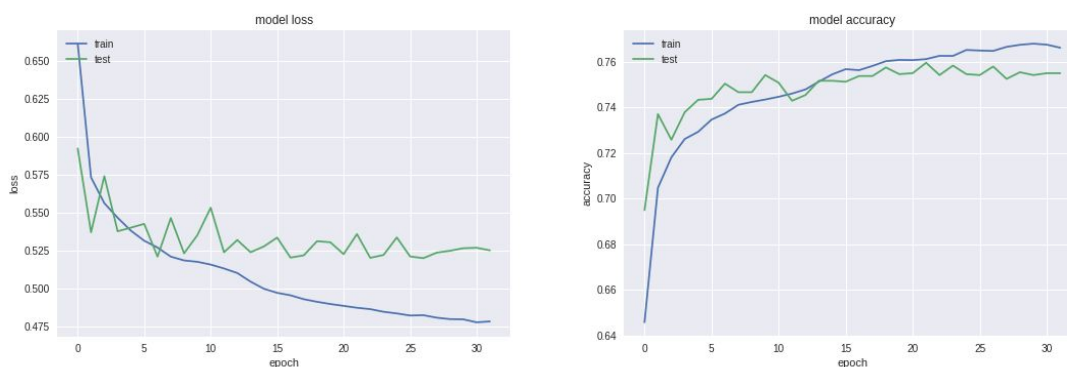


Problem 2. (1%) 請敘述你如何 improve performance(preprocess, embedding, 架構等), 並解釋為何這些做法可以使模型進步。

我在斷字之後有把重複的字詞過濾掉, 例如連續相同的emoji只會抓第一個。
接著是把常常出現的「b01」或「B01」這種沒有特別意義的詞用python re統一換成「Dcard用戶」減少雜訊。
另外我還有試過各種不同的word2vec參數、input的長度等超參數, 想辦法找到最好的。

另外也多加幾層Dense以及Dropout, 再來設定LSTM的Dropout, 效果都會變得更好。

Problem 3. (1%) 請比較不做斷詞 (e.g., 以字為單位) 與有做斷詞,兩種方法實作出來的效果差異,並解釋為何有此差別。



public score為0.75082, 比最好的RNN少了0.01左右, 原因可能是因為中文的字跟詞是不太一樣的單位, 若是只看一個個字沒辦法精確的考慮到差異而被誤導, 例如「白癡」跟「白色」都有白字, 但兩者的意義相差甚遠。

Problem 4. (1%) 請比較 RNN 與 BOW 兩種不同 model 對於”在說別人白痴之前,先想想自己”與”在說別人之前先想想自己,白痴”這兩句話的分數(model output),並討論造成差異的原因。

RNN: 前一句0.45447743、後一句0.5177116

BOW: 前一句0.53776234、後一句0.5377624

BOW只考慮詞出現的頻率, 而沒有考慮到順序關係, 因此兩句話的分數會非常接近。而RNN則會考慮到上下文關係, 所以可以看到RNN前一句分數低於0.5判斷不是惡意留言而後一句則高於0.5判斷為惡意留言。

Problem 5.

Problem 6.