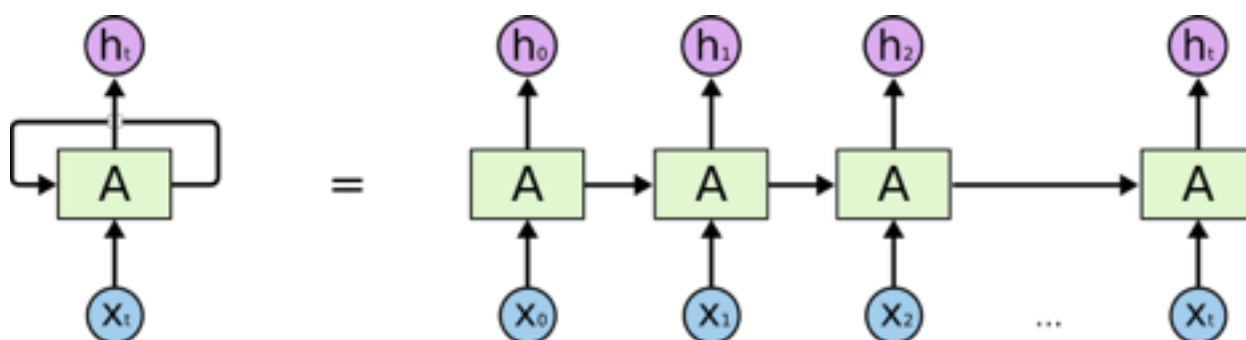


## 1. Model description

### (1) RNN (1%)

RNN(遞歸神經網絡)是時間遞歸神經網絡 (recurrent neural network)，由不同層的神經元間互相連接構成一個矩陣。相較傳統單純的類神經網路(neural network)，RNN透過將上一層的輸出便作下一層的輸入輸進模型，得到擁有前面輸入資料的資訊關聯。

RNN可以描述動態時間行為，因為和傳統類神經前饋網絡 (feedforward neural network) 接受較特定結構的輸入不同，RNN將狀態在自身網絡中循環傳遞，因此可以接受更廣泛的時間序列結構輸入。



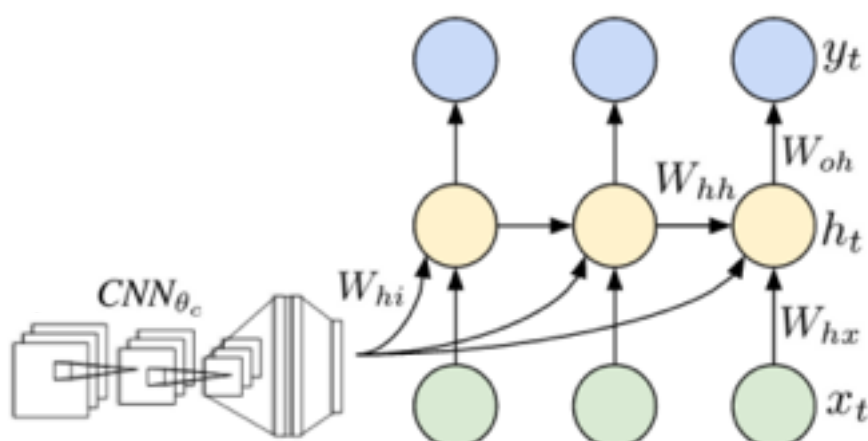
然而最簡單的RNN無法處理因為隨著遞歸，權重指數級爆炸或消失的問題 (Vanishing gradient problem)，難以捕捉長期時間關聯；而LSTM使用了記憶閘、遺忘閘可以很好地解決這個問題。

### (2) RNN+CNN (1%)

卷積神經網路 (Convolutional Neural Network, CNN) 是一種神經網絡，透過用filter一部分覆蓋範圍內的周圍單元，來找特定的特徵以讓模型重複利用。

CNN由一個或多個可供調整大小、橫跨範圍的filters所組成，同時也包括pooling layer來縮減特徵的複雜度。這一結構使得CNN能夠利用輸入數據的二維結構，在相比較不同類型的類神經網路，CNN需要考量的參數更少。

再透過CNN的filter取出關鍵特徵後就可以透過RNN進行進一步的訓練，不僅提升訓練的速度、也更可以降低模型的冗余複雜度。



## 2. How to improve your performance

### (1) Describe the model or technique

我的model不像大家針對每個句子做padding後來當作sequence輸入進RNN，我是使用一個TimeWindow其長度為31，在同一個sentence中將每個frame前後共31個frame(若無相對應的frame則補0)結合成為一個sequence，舉例：Sentence1有3個frame：[A,B,C], Sentence2有2個frame[D,E]，在經過前處理(TimeWindow)後會變成Sequence1：[0,A,B], Sequence2:[A,B,C], Sequence3:[B,C,0], Sequence4:[0,D,E], Sequence5:[D,E,0]。

我也針對這些輸入的Sequence先是逐一進行正規劃(BatchNormalization)、以8個3\*3大小的filters輸入進CNN，並以MaxPooling的方式提取特徵、再輸入進雙向Bidirectional LSTM(為多對一的LSTM結構)，最後經過兩層512大小的full connected layer後得到結果

### (2) Why do you use it

我認為與其直接由padding句子來作為sequence去訓練，並不如直接框一個timeWindow去做預測，因為在這個timeWindow內可以更強烈地依據附近的frame來做判斷phone。

接下來運用CNN以擷取這段TimeWindow特定的特徵。

使用多對一架構的LSTM除了解決vanishing的問題外，也是因為LSTM在經過前後的frame訓練後，將可以更好的判斷當下的frame。

最後為了提升模型複雜度多加了兩層的full connected layer。

## 3. Experimental results and settings

### (1) Compare and analyze the results between RNN and CNN

第一張圖與第二張圖分別為RNN與CNN的訓練過程，可以觀察出CNN明顯在表現上勝過RNN一籌，在耗時上也明顯較短。他們兩個的共通點在於當訓練了短短數次後，val\_loss便開始拉高而不再降低。

我想這除了有著overfitting的因素參雜其中，也跟我所選擇的TimeWindow大小以及模型複雜度有關。

我認為CNN的表現優於RNN的原因主要是因為我選到的filter成功抓取到了frame之間的隱性關聯。

```
1012340/1012340 [=====] - 773s - loss: 0.9506 - acc: 0.6980 - val_loss: 0.8700 - val_acc: 0.7272
Epoch 2/10
1012340/1012340 [=====] - 771s - loss: 0.5046 - acc: 0.8280 - val_loss: 0.8783 - val_acc: 0.7419
Epoch 3/10
1012340/1012340 [=====] - 770s - loss: 0.3533 - acc: 0.8758 - val_loss: 0.9586 - val_acc: 0.7466
Epoch 4/10
1012340/1012340 [=====] - 770s - loss: 0.2747 - acc: 0.9013 - val_loss: 1.0657 - val_acc: 0.7433
Epoch 5/10
1012340/1012340 [=====] - 773s - loss: 0.2324 - acc: 0.9157 - val_loss: 1.0969 - val_acc: 0.7449
Epoch 6/10
1012340/1012340 [=====] - 770s - loss: 0.2016 - acc: 0.9263 - val_loss: 1.1805 - val_acc: 0.7438
Epoch 7/10
1012340/1012340 [=====] - 769s - loss: 0.1812 - acc: 0.9336 - val_loss: 1.2215 - val_acc: 0.7482
```

```
1012340/1012340 [=====] - 531s - loss: 1.3061 - acc: 0.5982 - val_loss: 0.9056 - val_acc: 0.7009
Epoch 2/10
1012340/1012340 [=====] - 525s - loss: 0.6913 - acc: 0.7690 - val_loss: 0.7775 - val_acc: 0.7492
Epoch 3/10
1012340/1012340 [=====] - 524s - loss: 0.4948 - acc: 0.8303 - val_loss: 0.7956 - val_acc: 0.7568
```

## (2) Compare and analyze the results with other models (other models can be variant of basic RNN, like LSTM, or some novel ideas you use)

我曾經使用padding sentence來作為sequence輸入，並搭配建範圍為3的TimeWindow來進行CNN特徵擷取。乍看之下訓練效果不錯(88% validation accuracy及0.4的error)

```
2476/2476 [=====] - 269s - loss: 0.4495 - acc: 0.8598 - val_loss: 0.4423 - val_acc: 0.8617
Epoch 10/20
2476/2476 [=====] - 278s - loss: 0.4261 - acc: 0.8652 - val_loss: 0.4369 - val_acc: 0.8628
Epoch 11/20
2476/2476 [=====] - 268s - loss: 0.4057 - acc: 0.8716 - val_loss: 0.4241 - val_acc: 0.8662
Epoch 12/20
2476/2476 [=====] - 267s - loss: 0.3863 - acc: 0.8773 - val_loss: 0.4252 - val_acc: 0.8656
Epoch 13/20
2476/2476 [=====] - 267s - loss: 0.3693 - acc: 0.8820 - val_loss: 0.3954 - val_acc: 0.8748
Epoch 14/20
2476/2476 [=====] - 263s - loss: 0.3516 - acc: 0.8875 - val_loss: 0.3984 - val_acc: 0.8731
Epoch 15/20
2476/2476 [=====] - 264s - loss: 0.3371 - acc: 0.8922 - val_loss: 0.3842 - val_acc: 0.8777
Epoch 16/20
2476/2476 [=====] - 265s - loss: 0.3245 - acc: 0.8957 - val_loss: 0.3731 - val_acc: 0.8815
Epoch 17/20
2476/2476 [=====] - 266s - loss: 0.3097 - acc: 0.8999 - val_loss: 0.3884 - val_acc: 0.8768
Epoch 18/20
2476/2476 [=====] - 265s - loss: 0.2995 - acc: 0.9031 - val_loss: 0.3833 - val_acc: 0.8779
Epoch 19/20
2476/2476 [=====] - 263s - loss: 0.2864 - acc: 0.9073 - val_loss: 0.3682 - val_acc: 0.8843
Epoch 20/20
2476/2476 [=====] - 263s - loss: 0.2886 - acc: 0.9068 - val_loss: 0.3684 - val_acc: 0.8821
```

但是在經過後處理後在kaggle上最好的表現也在17, 18上下徘徊。

lol\_1509023389.csv  
a day ago by r06725053  
[add submission details](#)

17.94915

以下為我的改善及推論過程：

- I. 在最一開始我發現validation accuracy表現大概有86%但error很高(3點多)，雖然有加Mask layer，但後來發現應該是因為在經過CNN那一層後，由於CNN會帶出bias項目，因此原先設定Mask的遮罩值(0)變失效了，再增加了sample\_weight後，error降到0.6
- II. 後來認為應該是沒有亂數初始化(np.random.seed)以及正規化(BatchNormalization layer)，加上去之後儘管accuracy有了些許上升，但在kaggle上仍然成效不彰。
- III. 由於我的accuracy不錯加上error很低，因此助教建議我使用39維度的label class，但我在使用後在kaggle上的表現反而變糟了，我覺得accuracy高應該只是我的模型成功預測後面一大串padding都為0而造成的假象。

我試著尋找過原因，但因為時限而無法得出確切的結論，若有時間，我接下來會再往增加模型複雜度與結合fbank與mfcc兩種維度的feature產出的結果去做分析跟比對。