

outline

01



Introduction

02



Data -
preprocessing

03



Model
Architecture

04



Parameters tuning
&
Training experiments

05

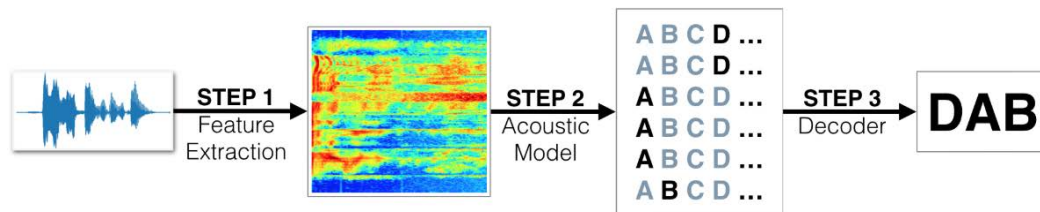
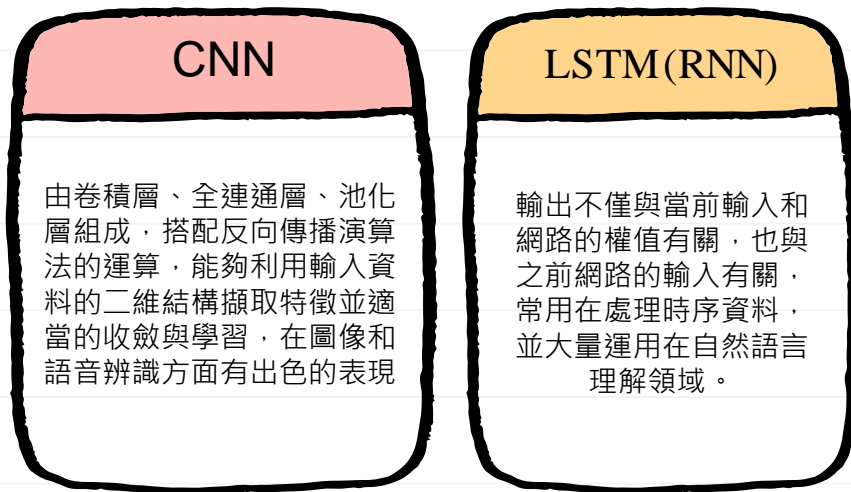


Conclusion

INTRODUCTION

深度學習的應用領域中，最大宗的莫過於圖像辨識 (Image Recognition) 及自然語言處理 (Natural Language Processing) ，因此本研究想藉此專題的實作機會，針對後者的領域進行一簡單的程式實作。

根據課堂所學以及網路上所閱讀之多元教材，利用卷積神經網路 (Convolutional Neural Network · CNN) 及長短期記憶模型 (Long Short-Term Memory · LSTM) 實作一簡易語音辨識模型，希望能對簡單的單詞進行辨識，也透過調參設計及實驗，以期發展一高準確率的辨識模型。



深度網路實現自動語音辨識 (Auto Speech Recognition, ASR) 的主要流程

Model 1 _ CNN

包含三個卷積層、三個最大池化層、
一個扁平層及兩個全連接層。

Training epoch = 15

```
model = deep_cnn(INPUT_SHAPE, NUM_CLASSES)
model.compile(optimizer='Adam', loss='categorical_crossentropy', metrics=['acc'])
```

```
callbacks = [EarlyStopping(monitor='val_acc', patience=4, verbose=1, mode='max')]

history = model.fit_generator(generator=dsGen.generator(BATCH, mode='train'),
                             steps_per_epoch=int(np.ceil(len(dsGen.df_train)/BATCH)),
                             epochs=EPOCHS,
                             verbose=1,
                             callbacks=callbacks,
                             validation_data=dsGen.generator(BATCH, mode='val'),
                             validation_steps=int(np.ceil(len(dsGen.df_val)/BATCH)))
```

```
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Input, Dense, Dropout, Flatten, Conv2D, MaxPooling2D, BatchNormalization

def deep_cnn(features_shape, num_classes, act='relu'):

    x = Input(name='inputs', shape=features_shape, dtype='float32')
    y = x

    # Block 1
    y = Conv2D(32, (3, 3), activation=act, padding='same', strides=1, name='block1_conv', input_shape=features_shape)(y)
    y = MaxPooling2D((3, 3), strides=(2,2), padding='same', name='block1_pool')(y)
    y = BatchNormalization(name='block1_norm')(y)

    # Block 2
    y = Conv2D(32, (3, 3), activation=act, padding='same', strides=1, name='block2_conv')(y)
    y = MaxPooling2D((3, 3), strides=(2,2), padding='same', name='block2_pool')(y)
    y = BatchNormalization(name='block2_norm')(y)

    # Block 3
    y = Conv2D(32, (3, 3), activation=act, padding='same', strides=1, name='block3_conv')(y)
    y = MaxPooling2D((3, 3), strides=(2,2), padding='same', name='block3_pool')(y)
    y = BatchNormalization(name='block3_norm')(y)

    # Flatten
    y = Flatten(name='flatten')(y)

    # Dense Layer
    y = Dense(64, activation=act, name='dense')(y)
    y = BatchNormalization(name='dense_norm')(y)
    y = Dropout(0.2, name='dropout')(y)

    # Predictions
    y = Dense(num_classes, activation='softmax', name='pred')(y)

    # Print network summary
    Model(inputs=x, outputs=y).summary()

    return Model(inputs=x, outputs=y)
```

Python / Keras

Parameters tuning



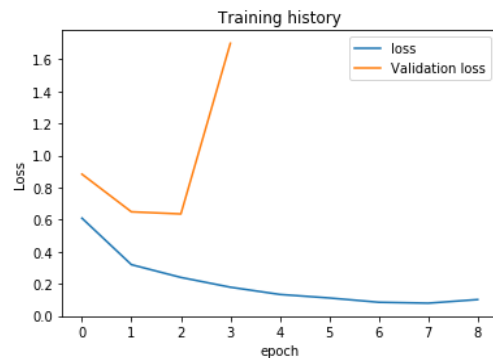
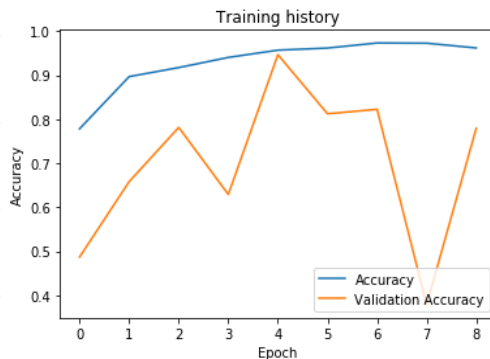
EXPERIMENT

1



Accuracy = 77.8%

- batch size : 32
- Activation : sigmoid
- Optimizer : adam



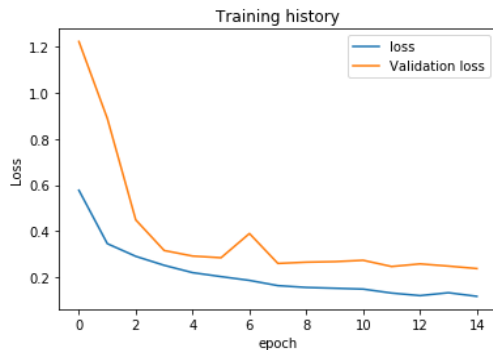
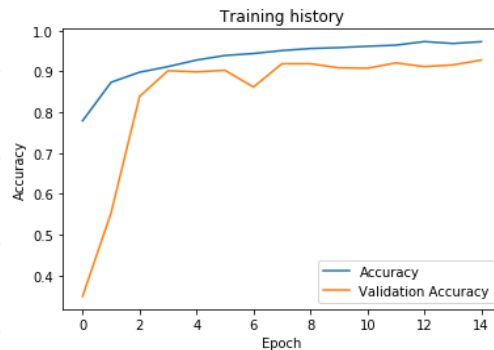
EXPERIMENT

2



Accuracy = 92.0%

- batch size : 32
- Activation : softmax
- Optimizer : adagrad



Parameters tuning



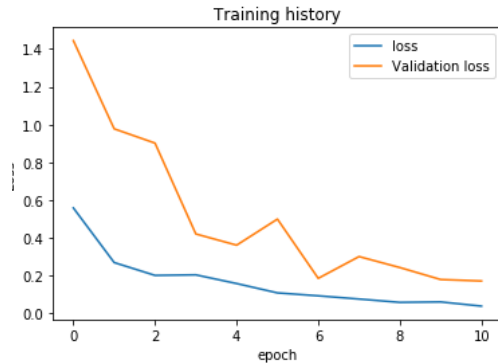
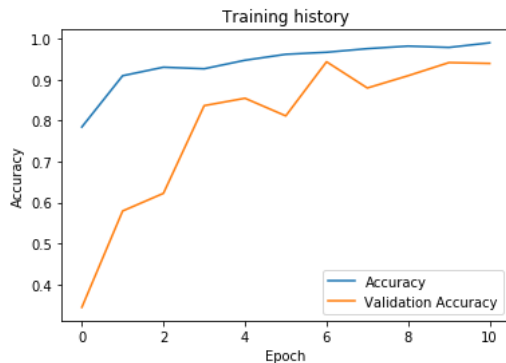
- batch size : 32
- Activation : softmax
- Optimizer : adam

EXPERIMENT

3



Accuracy = 94.4%



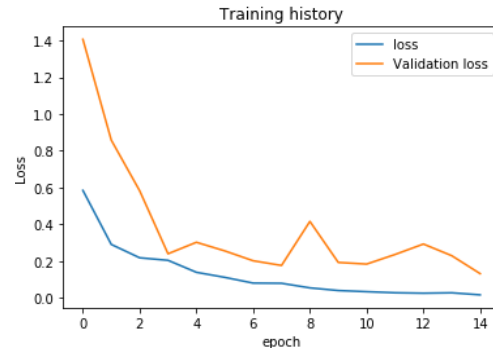
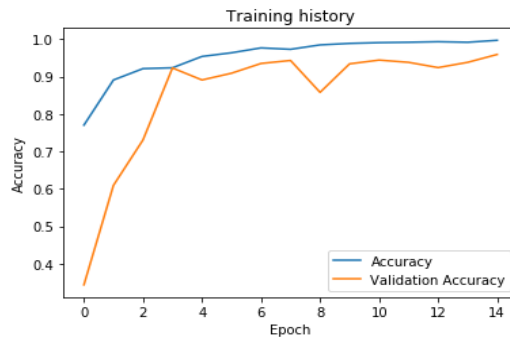
- batch size : 64
- Activation : softmax
- Optimizer : adam

EXPERIMENT

4



Accuracy = 95.6%



Parameters tuning



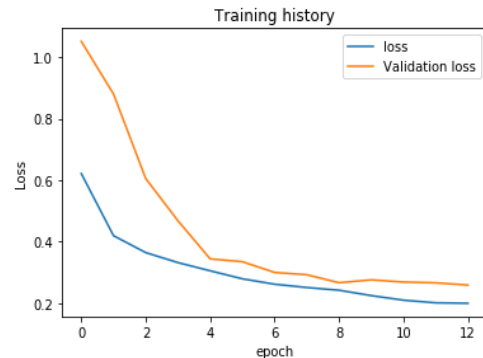
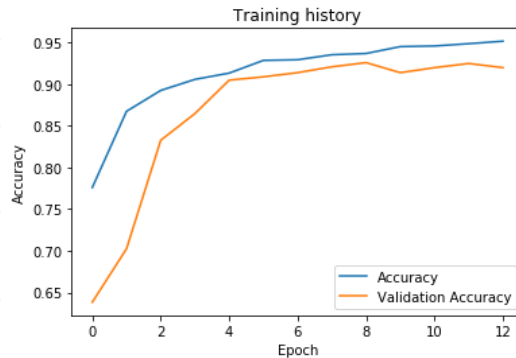
EXPERIMENT

5



Accuracy = 92.7%

- batch size : 32
- Activation : sigmoid
- Optimizer : adagrad



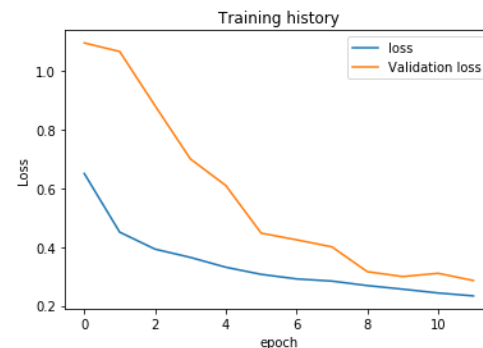
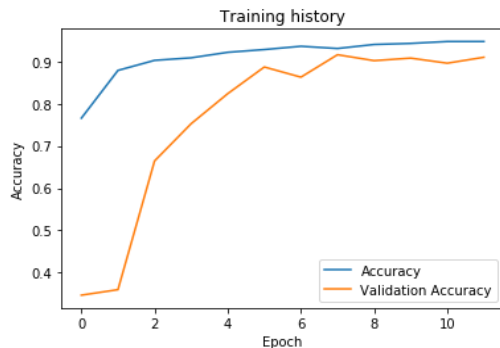
EXPERIMENT

6



Accuracy = 92.5%

- batch size : 64
- Activation : sigmoid
- Optimizer : adagrad



Model 2 _ LSTM

Training epoch = 15

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Activation
from tensorflow.keras.layers import LSTM

model = Sequential()
model.add(LSTM(128, dropout=0.2, recurrent_dropout=0.2))
model.add(Dense(1, activation='sigmoid'))

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
```

Python / Keras

```
callbacks = [EarlyStopping(monitor='val_accuracy', patience=4, verbose=1, mode='max')]
history = model.fit_generator(generator=dsGen.generator(BATCH, mode='train'),
                             steps_per_epoch=int(np.ceil(len(dsGen.df_train)/BATCH)),
                             epochs=EPOCHS,
                             verbose=1,
                             callbacks=callbacks,
                             validation_data=dsGen.generator(BATCH, mode='val'),
                             validation_steps=int(np.ceil(len(dsGen.df_val)/BATCH)))
```

Parameters tuning



- batch size : 32
- Activation : sigmoid
- Optimizer : adam

EXPERIMENT

7



Accuracy = 66.7%



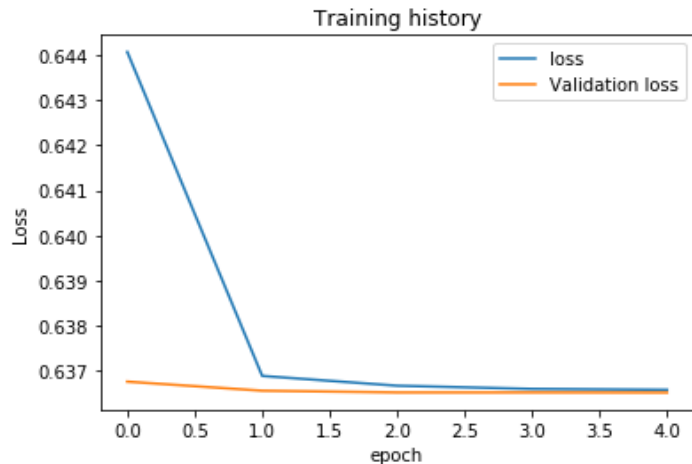
- batch size : 64
- Activation : softmax
- Optimizer : rmsprop

EXPERIMENT

8



Accuracy = 33.3%



SUMMARY

CNN

LSTM

EXP 1

EXP 2

EXP 3

EXP 4

EXP 5

EXP 6

EXP7

EXP 8

Batch size

32

32

32

64

32

64

32

32

Activate
function

Rel u
sigmoid

Rel u
soft max

Rel u
soft max

Rel u
soft max

Rel u
sigmoid

Rel u
sigmoid

sigmoid

soft max

Opt imizer

adam

adagrad

adam

adam

adagrad

adagrad

adam

r mspr op

Test ing
accuracy

77.8%

92.0%

94.4%

95.6%

92.7%

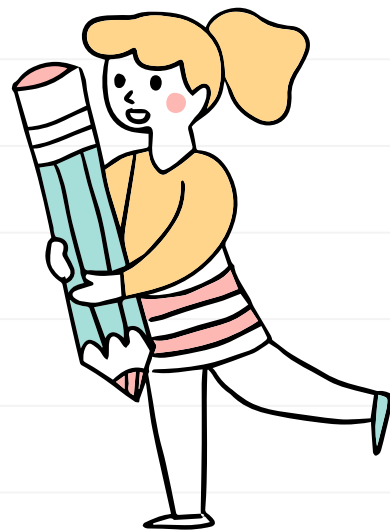
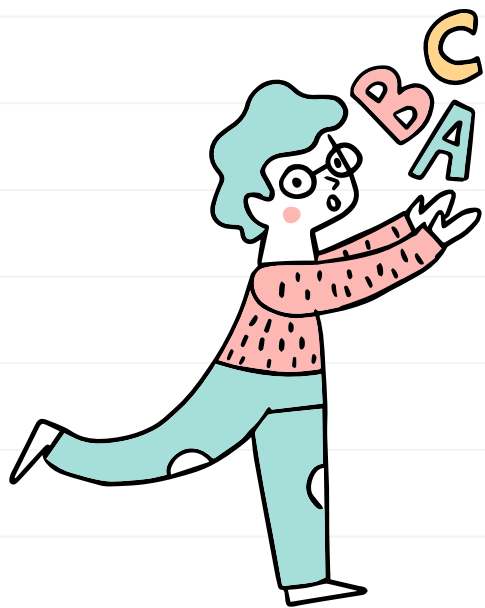
92.5%

66.7%

33.3%

CONCLUSION

本研究利用CNN 及LSTM 兩種深度學習模型進行語音辨識應用的實作，由多次調參的實驗結果可以發現，CNN 對於圖片資料的學習與辨識度較高，大多數皆有9成以上的準確率，經過優化後準確率更可達95.6%。LSTM 雖有記憶特性，但在此例中相對表現較差，未來可考慮結合兩種神經網路，以期進一步提升模型準確率。



THANKS!

