

TRƯỜNG ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG



BÁO CÁO
PROJECT III

ĐỀ TÀI: Hệ thống nhận dạng giọng nói (Speech-to-text)

Giảng viên hướng dẫn:

TS. Nguyễn Tuấn Dũng

Sinh viên thực hiện:

Nguyễn Đức Quyết – 20177021

Hà Nội, ngày 04 tháng 12 năm 2020

Mục lục

I.	Nhận xét của giáo viên.....	3
II.	Lời mở đầu	4
III.	Nội dung chi tiết	4
1.	Lý thuyết	4
2.	Giải pháp	7
3.	Souce code.....	8
3.1	Dataset.....	8
3.2	Tiền xử lý dữ liệu	8
3.3	Mô hình	11
3.4	Đánh giá	13
3.5	Nhận xét	14
3.6	Phương hướng phát triển	14

This image shows a full page of primary-ruled paper. It features multiple horizontal rows, each defined by two parallel dotted lines. The rows are evenly spaced across the entire page, providing a guide for handwriting practice. There are no margins, text, or other markings present.

II. Lời mở đầu

Với sự hướng dẫn tận tình của thầy **Nguyễn Tuấn Dũng**, đã hỗ trợ em rất nhiều để hoàn thành báo cáo môn học. Tuy nhiên trong quá trình tìm hiểu, học tập bản thân em còn nhiều thiếu sót. Em rất mong nhận được sự đóng góp của thầy. Em xin chân thành cảm ơn.

Nhận diện giọng nói đang xâm nhập vào cuộc sống hiện đại. Nó được cài đặt trong những chiếc điện thoại, điều khiển trò chơi hay những chiếc đồng hồ thông minh. Chỉ với khoảng \$50, bạn có thể có Amazon Echo Dot - một chiếc hộp thần kỳ cho phép bạn đặt pizza, nhận thông tin dự báo thời tiết hoặc thậm chí mua những vật dụng - chỉ bằng cách đưa ra mệnh lệnh.

Speech-to-text (STT) hay còn được xem là speech recognition – nhận dạng giọng nói là 1 lĩnh vực con trong NLP – xử lý ngôn ngữ tự nhiên, trong đó STT phát triển lý thuyết và công nghệ nhận dạng và chuyển dữ liệu ngôn ngữ nhận dạng âm thanh thành ký tự.

III. Nội dung chi tiết

1. Lý thuyết

Tín hiệu âm thanh là gì?

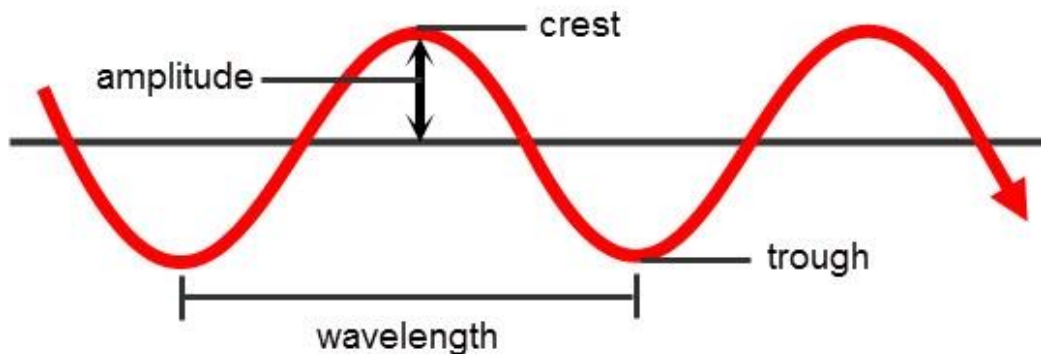
Điều này khá trực quan - bất kỳ đối tượng nào rung đều tạo ra sóng âm thanh. Bạn đã bao giờ nghĩ về cách chúng ta có thể nghe thấy giọng nói của ai đó? Đó là do sóng âm thanh. Hãy nhanh chóng hiểu quá trình đằng sau nó.

Khi một vật thể rung lên, các phân tử không khí dao động đến và đi từ vị trí nghỉ ngơi của chúng và truyền năng lượng của nó đến các phân tử lân cận. Điều này dẫn đến việc truyền năng lượng từ phân tử này sang phân tử khác, từ đó tạo ra sóng âm.

Âm thanh được truyền dưới dạng sóng, có thể lưu trữ âm thanh được ghi lại bằng độ cao của sóng âm sau các khoảng thời gian bằng nhau, phương thức này được gọi là lấy mẫu – **sampling**.

Thông số của tín hiệu âm thanh

- **Biên độ:** Biên độ đề cập đến sự dịch chuyển tối đa của các phân tử không khí từ vị trí còn lại
- **Đỉnh và máng:** Đỉnh là điểm cao nhất trong sóng trong khi máng là điểm thấp nhất
- **Bước sóng:** Khoảng cách giữa 2 đỉnh hoặc máng liên tiếp được gọi là bước sóng
- **Chu kỳ:** Mỗi tín hiệu âm thanh đi qua trong các hình thức của chu kỳ. Một chuyển động lên hoàn toàn và chuyển động xuống của tín hiệu tạo thành một chu kỳ
- **Tần số:** Tần số đề cập đến tốc độ thay đổi tín hiệu trong một khoảng thời gian



Một số thách thức

Cùng một nội dung nhưng khác nhau về âm điệu, cao độ, âm lượng,...

Cùng một từ nhưng có thể phát âm khác nhau tùy theo vùng, lãnh thổ, cùng nói một ngôn ngữ.

Chuyển âm thanh thành số

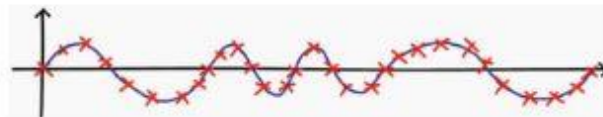
Bước đầu tiên trong nhận diện giọng nói là chúng ta cần truyền sóng âm vào máy tính.

Sóng âm có **một chiều** dữ liệu. Ở mỗi thời điểm, chúng có một giá trị cao độ. Để chuyển sóng âm thành số, chúng ta chỉ cần ghi lại cao độ(sampling).

Step 1: Analog audio signal - Continuous representation of signal



Step 2: Sampling - Samples are selected at regular time intervals

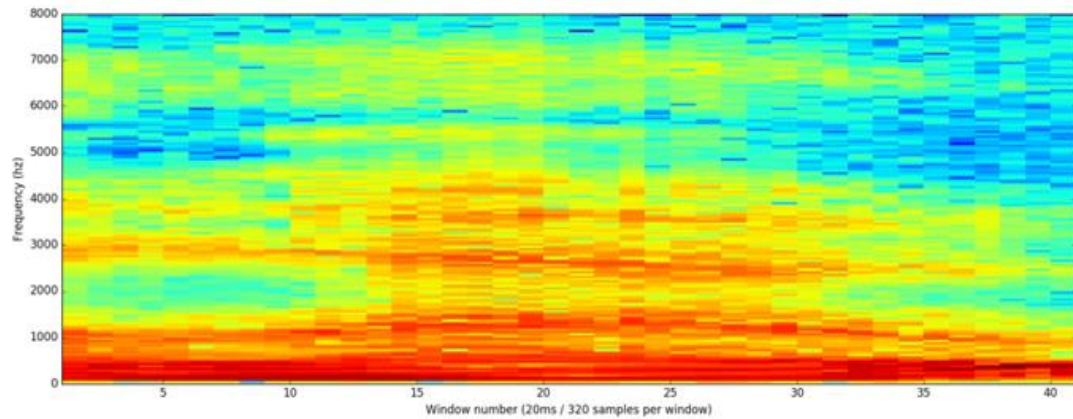


Step 3: Digital audio signal - The way it is stored in memory



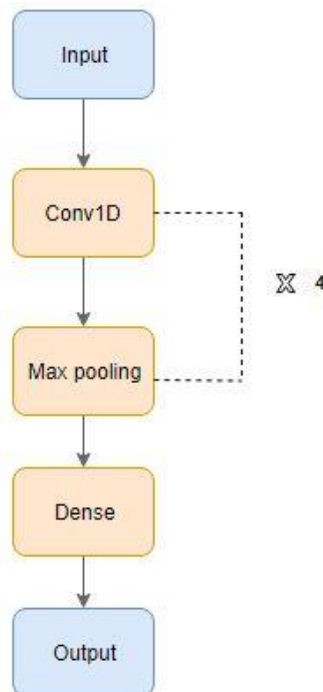
Do sóng âm rất phức tạp, kể cả những đoạn ghi âm ngắn cũng là một hỗn độn vì chứa âm thấp, âm trung và cả âm cao. Để máy tính có thể dễ dàng hiểu và học được các đặc trưng của sóng âm, ta sử dụng biến đổi Fourier transform trong toán

học thành biểu đồ quang phổ giúp chúng ta thực sự nhìn thấy âm thanh và cấu trúc cao độ của nó. Và lấy kết quả này là đặc trưng mà ta truyền vào mạng nơ-ron.



2. Giải pháp

Mô hình CNN: Convolutional neural network – mạng nơ-ron tích chập, là một trong những mô hình Deep Learning tiên tiến. giúp ta xây dựng được những hệ thống thông minh với độ chính xác cao. Mạng CNN là một tập hợp các lớp Convolution chồng lên nhau và sử dụng các hàm nonlinear activation như *ReLU*, *Tanh* để kích hoạt.

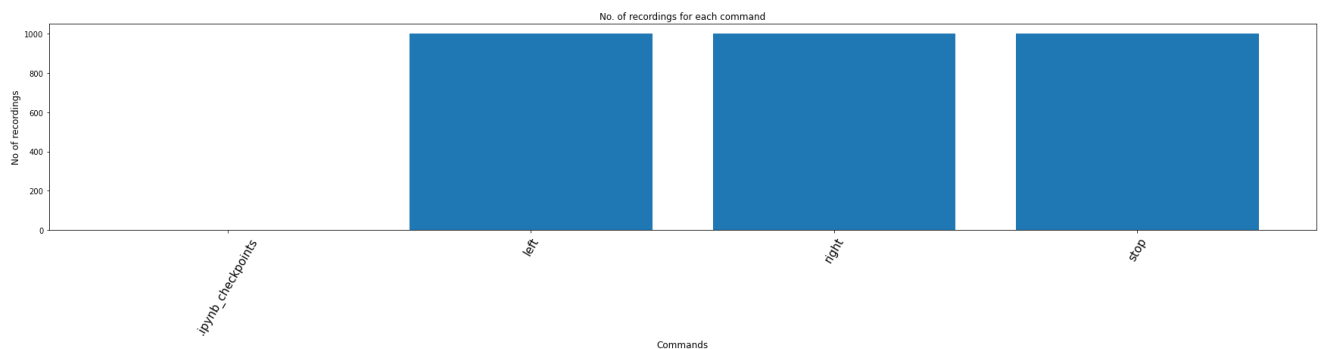


3. Source code

Tranning model sử dụng công cụ Jupyter Notebook.

3.1 Dataset

Sử dụng bộ dữ liệu giọng nói bao gồm 65.000 bản ghi của Tensorflow. Mỗi bản ghi âm có độ dài 1 giây bởi hàng ngàn người khác nhau. Trong khuôn khổ bản cáo cáo này, em sử dụng 3 từ bao gồm “left”, “right”, “stop”. Mỗi từ sử dụng 1000 bản ghi.



3.2 Tiền xử lý dữ liệu

Sử dụng thư viện **Librosa** và thư viện **SciPy** của Python để xử lý tín hiệu âm thanh.

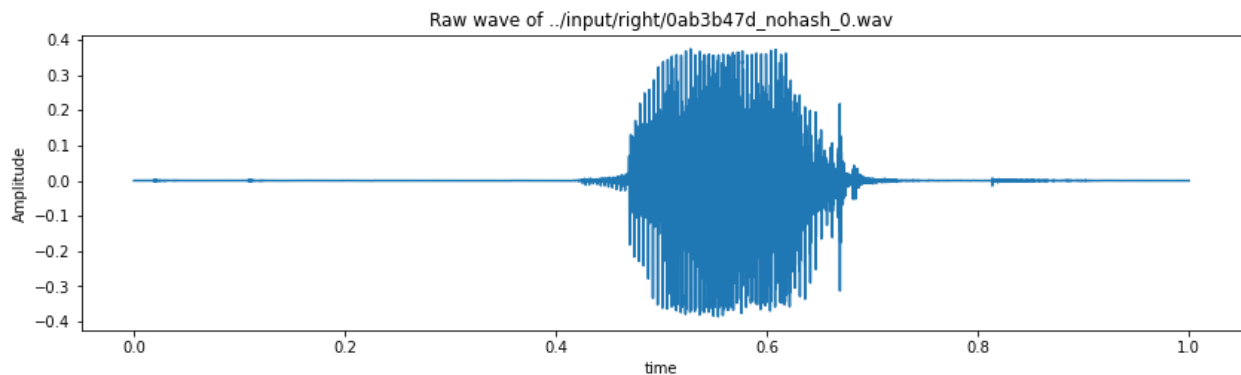
```
1 import librosa
2 import os
3 import IPython.display as idp
4 import matplotlib.pyplot as plt
5 import numpy as np
6 from scipy.io import wavfile
7 import warnings
8 warnings.filterwarnings("ignore")
```

Tiền xử lý dữ liệu.


```

1 # processing audio file
2 train_audio_path = "../input/"
3 samples, sample_rate = librosa.load(train_audio_path+"right/0ab3b47d_nohash_0.wav" , sr=16000)
4 fig = plt.figure(figsize=(14, 8))
5 ax1 = fig.add_subplot(211)
6 ax1.set_title('Raw wave of ' + '../input/right/0ab3b47d_nohash_0.wav')
7 ax1.set_xlabel('time')
8 ax1.set_ylabel('Amplitude')
9 ax1.plot(np.linspace(0, sample_rate/len(samples), sample_rate), samples)

```

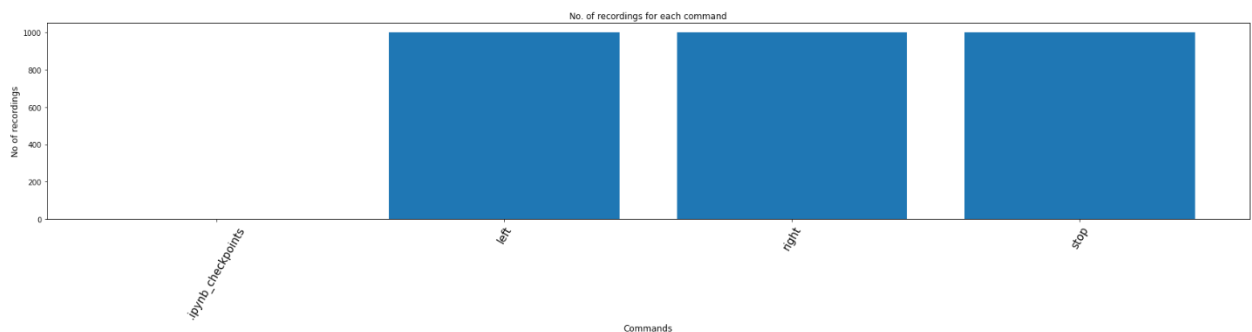


Lấy thông tin về bộ dữ liệu

```

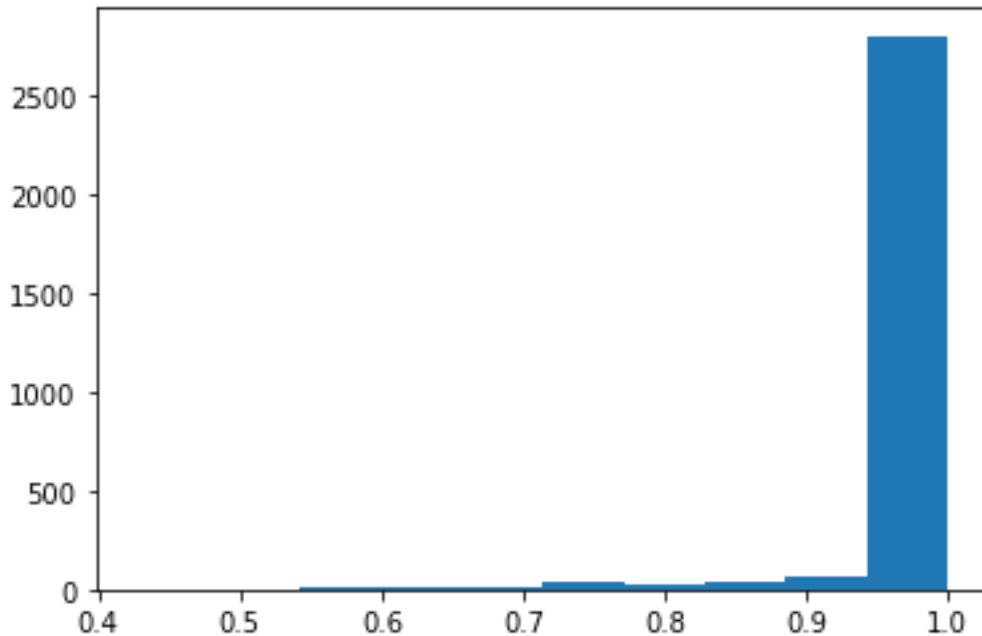
1 labels=os.listdir(train_audio_path)
2
3 #find count of each label and plot bar graph
4 no_of_recordings=[]
5 for label in labels:
6     print(label)
7     waves = [f for f in os.listdir(train_audio_path + label) if f.endswith('.wav')]
8     no_of_recordings.append(len(waves))
9
10
11 print(no_of_recordings)
12 # plot
13 plt.figure(figsize=(30,5))
14 index = np.arange(len(labels))
15 plt.bar(index, no_of_recordings)
16 plt.xlabel('Commands', fontsize=12)
17 plt.ylabel('No of recordings', fontsize=12)
18 plt.xticks(index, labels, fontsize=15, rotation=60)
19 plt.title('No. of recordings for each command')
20 plt.show()
21
22 labels=["left", "right", "stop"]

```



Tìm hiểu về các bản ghi âm

```
1 duration_of_recordings=[]
2 for label in labels:
3     waves = [f for f in os.listdir(train_audio_path + '/' + label) if f.endswith('.wav')]
4     for wav in waves:
5         sample_rate, samples = wavfile.read(train_audio_path + '/' + label + '/' + wav)
6         duration_of_recordings.append(float(len(samples)/sample_rate))
7     |
8 plt.hist(np.array(duration_of_recordings))
```



Tiền xử lý sóng âm thanh

Trong phần khám phá dữ liệu trước đó, chúng tôi đã thấy rằng thời lượng của một vài bản ghi âm nhỏ hơn 1 giây và tỷ lệ lấy mẫu (*sampling*) quá cao. Vì vậy, chúng ta hãy đọc sóng âm thanh và sử dụng các bước dưới đây tiền xử lý để đối phó với điều này.

Dưới đây là hai bước chúng ta sẽ làm theo:

- Lấy mẫu lại
- Loại bỏ các mẫu ngắn hơn dưới 1 giây

```

1 train_audio_path = './input/'
2
3 all_wave = []
4 all_label = []
5 for label in labels:
6     waves = [f for f in os.listdir(train_audio_path + '/' + label) if f.endswith('.wav')]
7     for wav in waves:
8         samples, sample_rate = librosa.load(train_audio_path + '/' + label + '/' + wav, sr = 16000)
9         samples = librosa.resample(samples, sample_rate, 8000)
10        if(len(samples)== 8000) :
11            all_wave.append(samples)
12            all_label.append(label)

```

Mã hóa nhãn

```

1 from sklearn.preprocessing import LabelEncoder
2 from keras.utils import np_utils
3 le = LabelEncoder()
4 y=le.fit_transform(all_label)
5 |
6 classes= list(le.classes_)
7 classes
8
9 y=np_utils.to_categorical(y, num_classes=len(labels))
10 y

```

Chia bộ dữ liệu training set và validation set: (80-20)

```

1 from sklearn.model_selection import train_test_split
2 x_tr, x_val, y_tr, y_val = train_test_split(np.array(all_wave),
3                                             np.array(y),
4                                             stratify=y,|
5                                             test_size = 0.2,
6                                             random_state=777,
7                                             shuffle=True)

```

3.3 Mô hình

```

1 from keras.layers import Dense, Dropout, Flatten, Conv1D, Input, MaxPooling1D
2 from keras.models import Model
3 from keras.callbacks import EarlyStopping, ModelCheckpoint
4 from keras import backend as K
5 K.clear_session()
6
7 inputs = Input(shape=(8000,1))
8
9 #First Conv1D layer
10 conv = Conv1D(8,13, padding='valid', activation='relu', strides=1)(inputs)
11 conv = MaxPooling1D(3)(conv)
12 conv = Dropout(0.3)(conv)
13
14 #Second Conv1D layer
15 conv = Conv1D(16, 11, padding='valid', activation='relu', strides=1)(conv)
16 conv = MaxPooling1D(3)(conv)
17 conv = Dropout(0.3)(conv)
18
19 #Third Conv1D layer
20 conv = Conv1D(32, 9, padding='valid', activation='relu', strides=1)(conv)
21 conv = MaxPooling1D(3)(conv)
22 conv = Dropout(0.3)(conv)
23
24 #Fourth Conv1D layer
25 conv = Conv1D(64, 7, padding='valid', activation='relu', strides=1)(conv)
26 conv = MaxPooling1D(3)(conv)
27 conv = Dropout(0.3)(conv)
28
29 #Flatten layer
30 conv = Flatten()(conv)
31
32 #Dense Layer 1
33 conv = Dense(256, activation='relu')(conv)
34 conv = Dropout(0.3)(conv)
35
36 #Dense Layer 2
37 conv = Dense(128, activation='relu')(conv)
38 conv = Dropout(0.3)(conv)
39
40 outputs = Dense(len(labels), activation='softmax')(conv)
41
42 model = Model(inputs, outputs)
43 model.summary()

```

Model: "functional_1"

Layer (type)	Output Shape	Param #
=====		
input_1 (InputLayer)	[(None, 8000, 1)]	0
conv1d (Conv1D)	(None, 7988, 8)	112
max_pooling1d (MaxPooling1D)	(None, 2662, 8)	0
dropout (Dropout)	(None, 2662, 8)	0
conv1d_1 (Conv1D)	(None, 2652, 16)	1424
max_pooling1d_1 (MaxPooling1D)	(None, 884, 16)	0
dropout_1 (Dropout)	(None, 884, 16)	0
conv1d_2 (Conv1D)	(None, 876, 32)	4640
max_pooling1d_2 (MaxPooling1D)	(None, 292, 32)	0
dropout_2 (Dropout)	(None, 292, 32)	0
conv1d_3 (Conv1D)	(None, 286, 64)	14400
max_pooling1d_3 (MaxPooling1D)	(None, 95, 64)	0
dropout_3 (Dropout)	(None, 95, 64)	0
flatten (Flatten)	(None, 6080)	0
dense (Dense)	(None, 256)	1556736
dropout_4 (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 128)	32896
dropout_5 (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 3)	387
=====		
Total params: 1,610,595		
Trainable params: 1,610,595		
Non-trainable params: 0		

Xác định loss function là: *categorical cross-entropy* cho vấn đề đa phân loại:

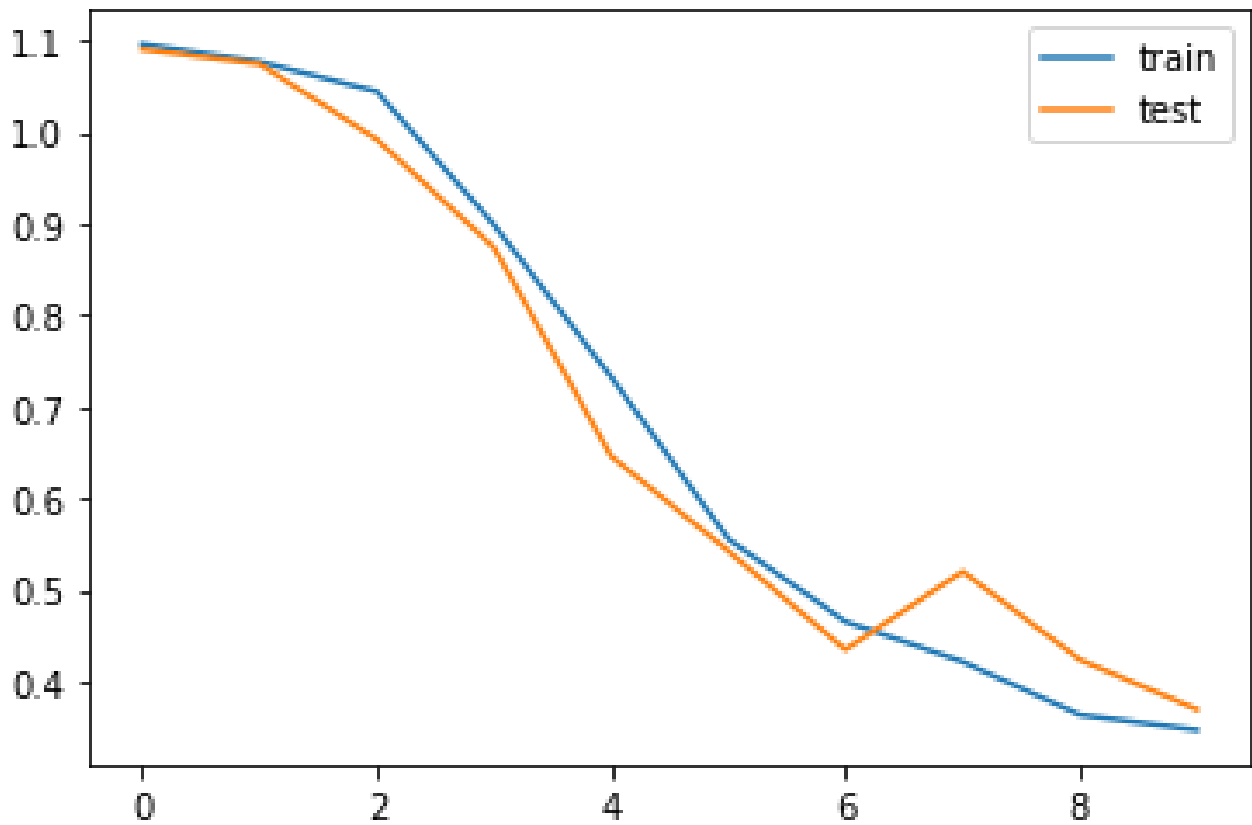
```
1 model.compile(loss='categorical_crossentropy',optimizer='adam',metrics=['accuracy'])
```

```
1 history=model.fit(x_tr, y_tr ,epochs=10, callbacks=[es,mc], batch_size=32, validation_data=(x_val,y_val))
```

3.4 Đánh giá

```
Epoch 10/10
69/69 [=====] - ETA: 0s - loss: 0.3482 - accuracy: 0.8657WARNING:tensorflow:Can save best model only w
ith val_acc available, skipping.
69/69 [=====] - 11s 161ms/step - loss: 0.3482 - accuracy: 0.8657 - val_loss: 0.4345 - val_accuracy: 0.
8597
```

```
1 from matplotlib import pyplot
2 pyplot.plot(history.history['loss'], label='train')
3 pyplot.plot(history.history['val_loss'], label='test')
4 pyplot.legend()
5 pyplot.show()
```



3.5 Nhận xét

Giá trị **accuracy** chưa được cao, đạt : **86,57%**

3.6 Phương hướng phát triển