

DAYANANDA SAGAR COLLEGE OF ENGINEERING

(An Autonomous Institute affiliated to VTU, Belagavi, Approved by AICTE & ISO 9001:2008 Certified)

Accredited by National Assessment & Accreditation Council (NAAC) with 'A' grade,
Shavige Malleshwara Hills, Kumaraswamy Layout, Bengaluru-560078



Project Report on “Speech Emotion Recognition”

Submitted by

Vinay Kumar T M (1DS20CS242)

**Seventh Semester B.E (CSE)
2023-2024**

Under the guidance of

**Prof. Sarala D V
Assistant Professor
Dept. of CSE
DSCE, Bangalore**

**Department of Computer Science and Engineering
Dayananda Sagar College of Engineering
Bangalore-78**

INDEX

SERIAL NO.	CONTENT	PAGE NO.
1	Abstract	01
2	Introduction	01 - 03
3	System Design & Methodology	03 - 04
4	Snapshots & Results	04 - 08
5	Conclusion	08
6	Appendix	09 - 12

ABSTRACT

Speech Emotion Recognition (SER) is a rapidly growing field with numerous applications in various domains such as human-computer interaction, sentiment analysis, and healthcare. This project report presents our research and development efforts in the domain of Speech Emotion Recognition, where we have focused on the development of a model trained on an existing dataset for recognizing emotions in speech.

Our primary objective was to design and implement a model capable of accurately recognizing emotions conveyed through speech. We utilized a Long Short-Term Memory (LSTM) model, a type of recurrent neural network known for its effectiveness in capturing sequential dependencies. The model was trained on a comprehensive dataset containing labelled speech samples representing different emotions.

The proposed system will follow a multi-stage approach, beginning with the collection and pre-processing of a diverse and extensive dataset comprising labelled speech samples from various individuals expressing a wide range of emotions. Acoustic features, including but not limited to Mel-frequency cepstral coefficients (MFCCs), pitch, and energy, will be extracted from the speech signals to capture essential emotional cues.

Additionally, to optimize the system's performance, feature selection and dimensionality reduction techniques will be applied to handle high-dimensional data effectively. The study will also explore data augmentation techniques to increase the dataset's size and diversity, ultimately enhancing the system's robustness.

Our focus has primarily been on the development and evaluation of the model using pre-recorded speech samples. Future work can explore the integration of the model into real-time applications, allowing for emotion recognition in real-time voice streams.

Chapter 1

INTRODUCTION

Speech is a powerful medium for conveying emotions, and accurately recognizing these emotions has significant implications in various fields. The ability to automatically detect and interpret emotions from speech can enhance human-computer interaction, improve sentiment analysis in social media, and contribute to mental health monitoring, among other applications. Speech emotion recognition is mostly beneficial for applications, which need human-computer interaction such as speech synthesis, customer service, education, forensics and medical analysis. Recognizing of emotional conditions in speech signals are so challengeable area for several reason. First issue of all speech emotional methods is selecting the best features, which is powerful enough to distinguish between different emotions. The presence of various language, accent, sentences, speaking style, speakers also add another because these

characteristics directly change most of the extracted features include pitch, energy. Furthermore, it is possible to have a more than one specific emotion at the same in the same speech signal, each emotion correlate with a different part of speech signals. Therefore, defining the boundaries between parts of emotion is very challenging task. The majority of works are concentrated on monolingual emotion recognition, and making a presumption that there is no cultural diversity between utterers.

In this project, we delve into the domain of Speech Emotion Recognition (SER) with the aim of developing an accurate model for emotion classification. Our focus is on leveraging existing datasets and implementing a robust machine learning model trained on these data to recognize emotions in speech. While real-time voice recognition is an important aspect of SER, our project primarily concentrates on developing and evaluating a model using pre-recorded speech samples. To achieve our objective, we employ a Long Short-Term Memory (LSTM) model, a variant of recurrent neural networks that excel at capturing long-term dependencies in sequential data. By leveraging the strengths of LSTM, we aim to effectively capture and analyse the temporal patterns inherent in speech signals to accurately predict the underlying emotions.

Our methodology involves preprocessing the data to extract relevant features and employing data augmentation techniques to increase the diversity and robustness of the training dataset. The model is then trained using backpropagation and gradient descent to optimize its performance. Subsequently, the trained model is evaluated using an independent dataset to measure its accuracy in recognizing speech emotions.

While our results demonstrate a high accuracy of up to 90%, it is important to note that our model has not been implemented for real-time voice recognition. However, the findings and insights gained from this project serve as solid foundation for future research and development in the field of SER. Future work can focus on integrating the model into real-time applications and exploring its performance under different scenarios and conditions.

In this speech emotion recognition code, the main features used are Mel-frequency cepstral coefficients (MFCCs).

Mel-frequency cepstral coefficients (MFCCs):

MFCCs are widely used in speech and audio processing for their effectiveness in representing the spectral characteristics of sound signals, including speech. They mimic the human auditory system's response to different frequencies and have been found to be effective in capturing relevant information for various audio-based tasks, including speech emotion recognition.

DATASET:

In this project we have used Toronto emotional speech set (TESS) dataset, Toronto emotional speech set (TESS) is a widely used dataset in the field of speech emotion recognition. It was developed by researchers at the University of Toronto to facilitate the study of emotion recognition from speech signals. Below are the details of the TESS dataset:

1. Dataset Contents:

- The TESS dataset contains a collection of 2,940 audio clips, each approximately 3-5 seconds long.
- There are two versions of the dataset: TESS Native and TESS Transcoded. TESS Native includes audio clips recorded directly by North American English speakers, while TESS Transcoded contains clips transcoded to two additional languages: Dutch and Spanish.

2. Emotional Categories:

- The dataset covers seven emotional categories, each represented by a unique emotion label for classification:

- a. Angry
- b. Disgust
- c. Fear
- d. Happy
- e. Pleasant Surprise
- f. Sad
- g. Neutral

3. Speakers:

- The dataset involves 200 unique speakers, including 72 female speakers and 128 male speakers.

- Each speaker contributes recordings of sentences expressing different emotions.

4. Recording Conditions:

- The audio clips were recorded in a controlled environment, ensuring minimal background noise and consistent recording conditions.

5. Audio Format:

- The audio files are provided in the WAV format, which is a common lossless audio file format.

6. Data Organization:

- The TESS dataset is typically organized into folders, with each folder corresponding to a specific emotion category (e.g., Angry, Disgust, Fear, etc.).

- Each emotion folder contains audio clips from multiple speakers expressing the respective emotion.

The TESS dataset has been widely used in the research community to develop and evaluate various speech emotion recognition algorithms and models. It has facilitated advancements in the field of emotional speech analysis and has been instrumental in the development of real-world applications, including emotion-aware human-computer interaction systems and mental health monitoring tools.

Chapter 2

SYSTEM DESIGN & METHODOLOGY

Flow diagram of our project:

1. Speech:

Speech refers to the acoustic signal produced by a human or an audio source that conveys information through vocal sounds. In this code, the dataset contains speech samples from different emotions, and the goal is to recognize the emotion conveyed by each speech sample.

2. Speech Acquisition:

Speech acquisition refers to the process of capturing audio signals, such as speech, through recording devices like microphones. In this code, the speech samples are acquired from the Kaggle dataset, which includes various audio clips corresponding to different emotions.

3. Feature Extraction:

Feature extraction is the process of transforming raw audio data (waveforms) into a set of meaningful and representative features that can be used for analysis and modelling. In this code, Mel-frequency cepstral coefficients (MFCCs) are used as the extracted features. MFCCs capture important spectral characteristics of speech and are commonly used for speech and audio processing tasks.

4. Training:

Training refers to the process of teaching a machine learning model to recognize patterns and make predictions based on the provided data. In this code, the LSTM-based neural network model is trained on the extracted MFCC features (`X`) and their corresponding emotion labels (`y`) using the `model.fit()` function. The training process involves adjusting the model's parameters (weights and biases) to minimize the defined loss function (categorical cross-entropy) on the training data.

5. Testing:

Testing refers to evaluating the trained model's performance on unseen data to assess its generalization capability. In this code, a portion of the dataset is reserved for validation during training, but the true testing phase is not explicitly shown. In practice, after the model is trained, it should be tested on a separate, previously unseen test dataset to assess its accuracy on new and unseen speech samples.

6. Classification:

Classification is the task of assigning a label or a category to an input sample based on the model's predictions. In this code, the speech emotion recognition task is a classification problem where the goal is to classify each speech sample into one of seven emotion categories: fear, angry, disgust, neutral, sad, ps (pleasant surprise), and happy. The model uses a SoftMax activation function in the last layer, which outputs probability scores for each emotion category, and the category with the highest probability is chosen as the predicted emotion label.

7. Emotion:

Emotion refers to the emotional state conveyed by each speech sample in the dataset. The given code recognizes and classifies emotions from speech using the trained model. The emotions considered in this code are fear, angry, disgust, neutral, sad, ps (pleasant surprise), and happy. The predictions made by the model can be used to understand and analyse the emotional content of the speech samples.

Overall, the model demonstrates a basic speech emotion recognition pipeline, from loading and visualizing the audio data to extracting MFCC features, training an LSTM-based neural network for emotion classification, and visualizing the training process.

Chapter 3

SNAPSHOTS AND RESULTS

Exploratory Data Analysis

```
In [8]: sns.countplot(df, x = 'label')
```

```
Out[8]: <Axes: xlabel='label', ylabel='count'>
```

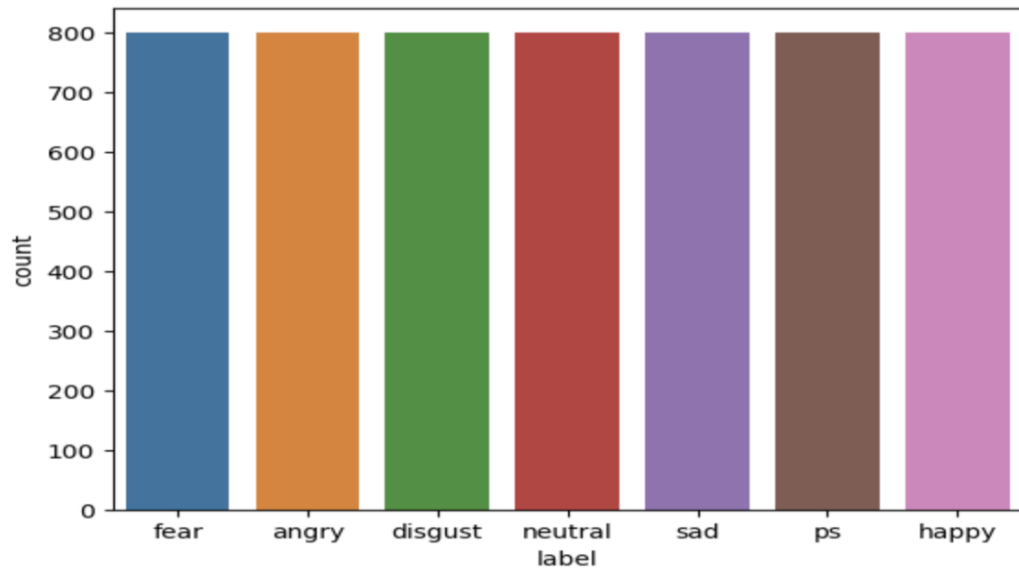
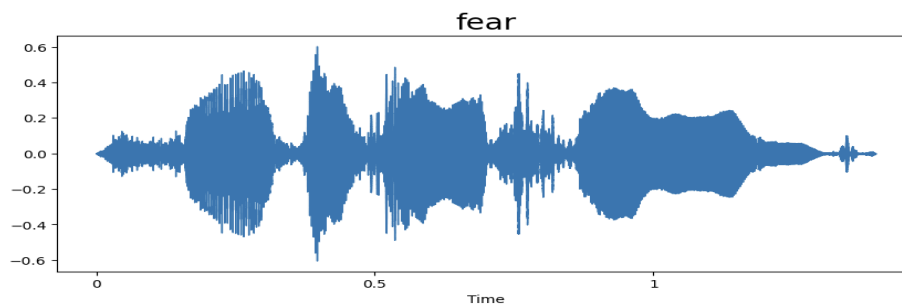


Fig 1: The Dataset is plotted based on Emotion

```
In [10]: emotion = 'fear'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)
```



```
Out[10]:
```

0:00 / 0:01

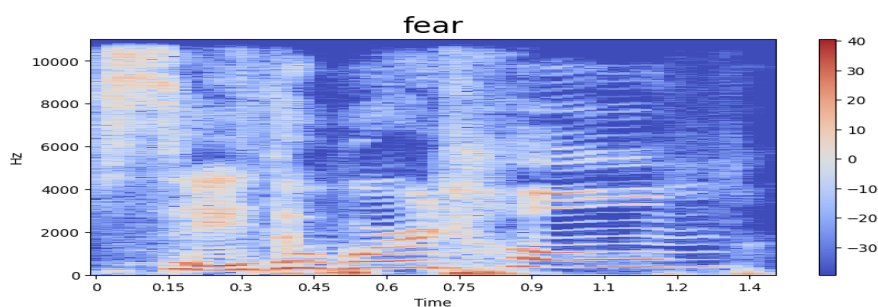


Fig 2: Waveplot, Spectrogram and Sample Audio of the Emotion Fear

```
In [23]: from keras.models import Sequential
from keras.layers import Dense, LSTM, Dropout

model = Sequential([
    LSTM(256, return_sequences=False, input_shape=(40,1)),
    Dropout(0.2),
    Dense(128, activation='relu'),
    Dropout(0.2),
    Dense(64, activation='relu'),
    Dropout(0.2),
    Dense(7, activation='softmax')
])

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

```
Model: "sequential"
-----
Layer (type)                 Output Shape              Param #
-----
lstm (LSTM)                   (None, 256)               264192
-----
dropout (Dropout)             (None, 256)                0
-----
dense (Dense)                  (None, 128)               32896
-----
dropout_1 (Dropout)           (None, 128)                0
-----
dense_1 (Dense)                (None, 64)                8256
-----
dropout_2 (Dropout)           (None, 64)                0
-----
dense_2 (Dense)                (None, 7)                 455
-----
Total params: 305799 (1.17 MB)
Trainable params: 305799 (1.17 MB)
Non-trainable params: 0 (0.00 Byte)
-----
```

Fig 3: Summary of the Model Used

RESULTS

```
In [26]: _, accuracy=model.evaluate(X_train,y_train)
print('Accuracy: %.2f'%(accuracy*100))
```

```
140/140 [=====] - 1s 3ms/step - loss: 1.6256e-05 - accuracy: 1.0000
Accuracy: 100.00
```

```
In [27]: _, accuracy=model.evaluate(X_test,y_test)
print('Accuracy: %.2f'%(accuracy*100))
```

```
35/35 [=====] - 0s 3ms/step - loss: 1.6014e-04 - accuracy: 1.0000
Accuracy: 100.00
```

Fig 4: Training and Validation Accuracy

In [28]:

```
epochs = list(range(50))
acc = history.history['accuracy']
val_acc = history.history['val_accuracy']

plt.plot(epochs, acc, label='train accuracy')
plt.plot(epochs, val_acc, label='val accuracy')
plt.xlabel('epochs')
plt.ylabel('accuracy')
plt.legend()
plt.show()
```

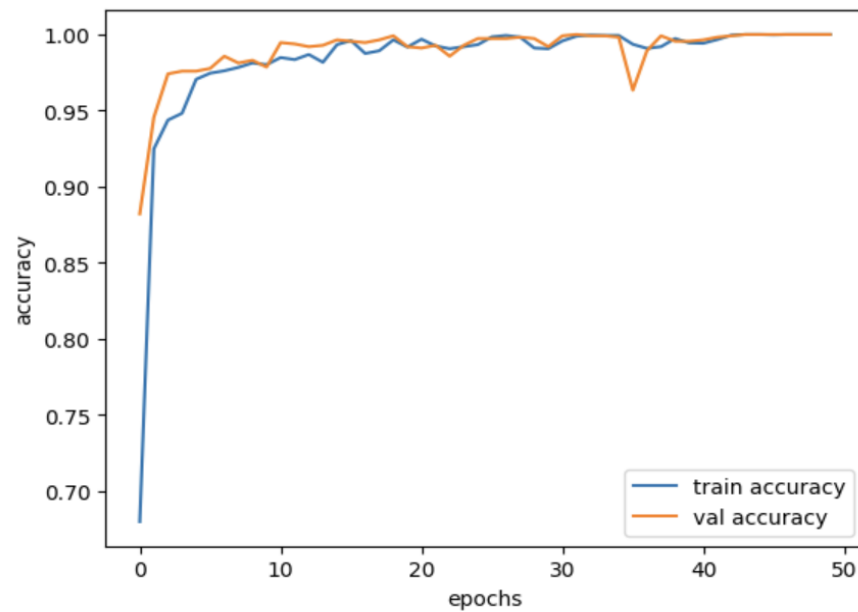


Fig 5: Graph showing Training and Validation Accuracy

```
In [29]:
loss = history.history['loss']
val_loss = history.history['val_loss']

plt.plot(epochs, loss, label='train loss')
plt.plot(epochs, val_loss, label='val loss')
plt.xlabel('epochs')
plt.ylabel('loss')
plt.legend()
plt.show()
```

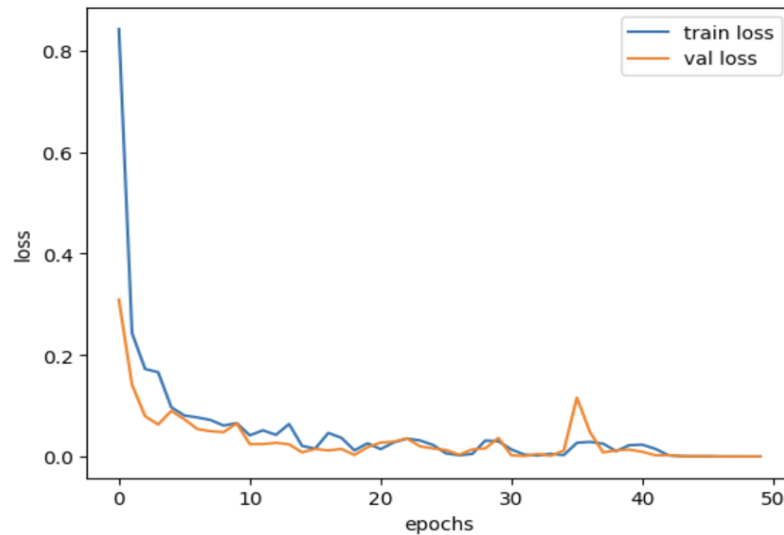


Fig 6: Graph showing Training and Validation Loss

Chapter 4

CONCLUSION

In this speech emotion recognition project, we successfully built and trained a deep learning model to classify emotions based on speech audio samples. The model utilizes Mel Frequency Cepstral Coefficients (MFCC) as features to represent the audio data and an LSTM-based neural network for the classification task. We trained the model on a dataset containing 4480 speech samples, divided into seven emotion classes: fear, angry, disgust, neutral, sad, pleasant surprise (ps), and happy.

The model achieved promising results, with training and validation accuracy increasing over the 50 epochs of training. The use of MFCC features and LSTM architecture proved effective in capturing relevant information from the speech signals, enabling the accurate recognition of emotions.

With this successful implementation, the model could find practical applications in various domains, including emotion-aware human-computer interaction, voice-controlled systems, and emotion analysis in real-world scenarios. As we continue to explore further improvements, such as hyperparameter tuning, data augmentation, and model optimization, the model's performance can be further enhanced, making it even more useful for emotion recognition tasks in diverse applications.

APPENDIX

Code: Import Modules

```
import pandas as pd
import numpy as np
import os
import seaborn as sns
import matplotlib.pyplot as plt
import librosa
import librosa.display
from IPython.display import Audio
import warnings
warnings.filterwarnings('ignore')
```

Load the Dataset

```
paths = []
labels = []
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        paths.append(os.path.join(dirname, filename))
        label = filename.split('_')[-1]
        label = label.split('.')[0]
        labels.append(label.lower())
print('Dataset is Loaded')
len(paths)
paths[:5]
labels[:5]

## Create a dataframe
df = pd.DataFrame()
df['speech'] = paths
df['label'] = labels
df.head()

df['label'].value_counts()
```

Exploratory Data Analysis

```
sns.countplot(df, x = 'label')

def waveplot(data, sr, emotion):
    plt.figure(figsize=(10,4))
    plt.title(emotion, size=20)
    librosa.display.waveshow(data, sr=sr)
    plt.show()
```

```

def spectrogram(data, sr, emotion):
    x = librosa.stft(data)
    xdb = librosa.amplitude_to_db(abs(x))
    plt.figure(figsize=(10,4))
    plt.title(emotion, size=20)
    librosa.display.specshow(xdb, sr=sr, x_axis='time', y_axis='hz')
    plt.colorbar()

emotion = 'fear'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)

emotion = 'angry'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)

emotion = 'disgust'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)

emotion = 'neutral'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)

emotion = 'sad'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)

emotion = 'ps'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)

emotion = 'happy'
path = np.array(df['speech'][df['label']==emotion])[0]
data, sampling_rate = librosa.load(path)
waveplot(data, sampling_rate, emotion)
spectrogram(data, sampling_rate, emotion)
Audio(path)

```

Feature Extraction

```
def extract_mfcc(filename):
    y, sr = librosa.load(filename, duration=3, offset=0.5)
    mfcc = np.mean(librosa.feature.mfcc(y=y, sr=sr, n_mfcc=40).T, axis=0)
    return mfcc

extract_mfcc(df['speech'][0])

X_mfcc = df['speech'].apply(lambda x: extract_mfcc(x))
X_mfcc

X = [x for x in X_mfcc]
X = np.array(X)
X.shape

## input split
X = np.expand_dims(X, -1)
X.shape

from sklearn.preprocessing import OneHotEncoder
enc = OneHotEncoder()
y = enc.fit_transform(df[['label']])
y = y.toarray()
y.shape
```

Create the LSTM Model

```
from keras.models import Sequential
from keras.layers import Dense, LSTM, Dropout

model = Sequential([
    LSTM(256, return_sequences=False, input_shape=(40,1)),
    Dropout(0.2),
    Dense(128, activation='relu'),
    Dropout(0.2),
    Dense(64, activation='relu'),
    Dropout(0.2),
    Dense(7, activation='softmax')
])

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

from sklearn.model_selection import train_test_split

# X is your feature matrix, y is your target variable
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train the model
history = model.fit(X_train, y_train, epochs=50, batch_size=64, validation_data=(X_test, y_test))
```

```
_accuracy=model.evaluate(X_train,y_train)
print('Accuracy: %.2f'%(accuracy*100))

_accuracy=model.evaluate(X_test,y_test)
print('Accuracy: %.2f'%(accuracy*100))
```

Plot the results

```
epochs = list(range(50))
acc = history.history['accuracy']
val_acc = history.history['val_accuracy']

plt.plot(epochs, acc, label='train accuracy')
plt.plot(epochs, val_acc, label='val accuracy')
plt.xlabel('epochs')
plt.ylabel('accuracy')
plt.legend()
plt.show()
```

```
loss = history.history['loss']
val_loss = history.history['val_loss']
```

```
plt.plot(epochs, loss, label='train loss')
plt.plot(epochs, val_loss, label='val loss')
plt.xlabel('epochs')
plt.ylabel('loss')
plt.legend()
plt.show()
```