

✓ Human Action Recognition in Videos using Keras (CNN + LSTM)



✓ Outline

- Step 1: Download and Visualize the Data with its Labels
- Step 2: Preprocess the Dataset
- Step 3: Split the Data into Train and Test Set
- Step 4: Implement the ConvLSTM Approach
 - Step 4.1: Construct the Model
 - Step 4.2: Compile & Train the Model
 - Step 4.3: Plot Model's Loss & Accuracy Curves
- Step 5: implement the LRCN Approach
 - Step 5.1: Construct the Model
 - Step 5.2: Compile & Train the Model
 - Step 5.3: Plot Model's Loss & Accuracy Curves
- Step 6: Test the Best Performing Model on YouTube videos

Alright, so without further ado, let's get started.

✓ Import the Libraries

We will start by installing and importing the required libraries.

```
# Discard the output of this cell.
%%capture

# Install the required libraries.
!pip install youtube-dl moviepy
!pip install git+https://github.com/TahaAnwar/pafy.git#egg=pafy
```

```
# Import the required libraries.
import os
import cv2
import pafy
import math
import random
import numpy as np
import datetime as dt
import tensorflow as tf
from collections import deque
import matplotlib.pyplot as plt

from moviepy.editor import *
%matplotlib inline

from sklearn.model_selection import train_test_split

from tensorflow.keras.layers import *
from tensorflow.keras.models import Sequential
from tensorflow.keras.utils import to_categorical
from tensorflow.keras.callbacks import EarlyStopping
from tensorflow.keras.utils import plot_model
```

And will set Numpy, Python, and Tensorflow seeds to get consistent results on every execution.

```
seed_constant = 27
np.random.seed(seed_constant)
random.seed(seed_constant)
tf.random.set_seed(seed_constant)
```

✓ Step 1: Download and Visualize the Data with its Labels

- 50 Action Categories
- 25 Groups of Videos per Action Category
- 133 Average Videos per Action Category
- 199 Average Number of Frames per Video
- 320 Average Frames Width per Video
- 240 Average Frames Height per Video
- 26 Average Frames Per Seconds per Video

```
# Discard the output of this cell.
%%capture

# Downlaod the UCF50 Dataset
!wget --no-check-certificate https://www.crcv.ucf.edu/data/UCF50.rar

#Extract the Dataset
!unrar x UCF50.rar
```

For visualization, we will pick 20 random categories from the dataset and a random video from each selected category and will visualize the first frame of the selected videos with their associated labels written. This way we'll be able to visualize a subset (20 random videos) of the dataset.

```
# Create a Matplotlib figure and specify the size of the figure.
plt.figure(figsize = (20, 20))

# Get the names of all classes/categories in UCF50.
all_classes_names = os.listdir('UCF50')

# Generate a list of 20 random values. The values will be between 0-50,
# where 50 is the total number of class in the dataset.
random_range = random.sample(range(len(all_classes_names)), 20)

# Iterating through all the generated random values.
for counter, random_index in enumerate(random_range, 1):

    # Retrieve a Class Name using the Random Index.
    selected_class_Name = all_classes_names[random_index]

    # Retrieve the list of all the video files present in the randomly selected Class Directory.
    video_files_names_list = os.listdir(f'UCF50/{selected_class_Name}')

    # Randomly select a video file from the list retrieved from the randomly selected Class Directory.
    selected_video_file_name = random.choice(video_files_names_list)

    # Initialize a VideoCapture object to read from the video File.
    video_reader = cv2.VideoCapture(f'UCF50/{selected_class_Name}/{selected_video_file_name}')

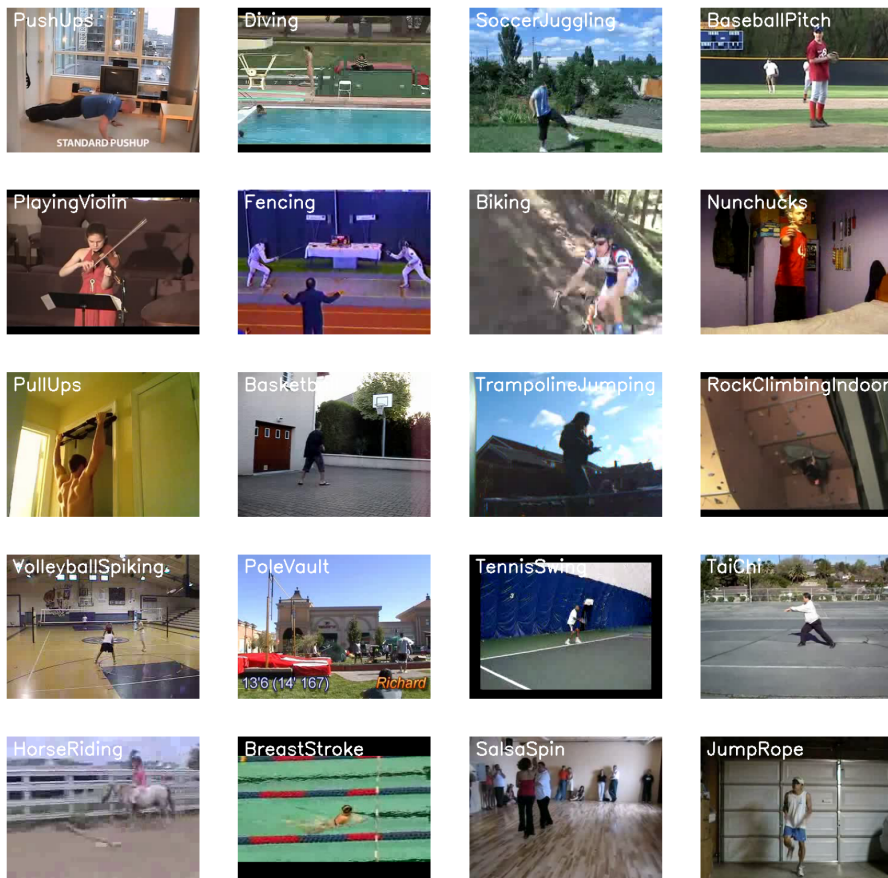
    # Read the first frame of the video file.
    _, bgr_frame = video_reader.read()

    # Release the VideoCapture object.
    video_reader.release()

    # Convert the frame from BGR into RGB format.
    rgb_frame = cv2.cvtColor(bgr_frame, cv2.COLOR_BGR2RGB)

    # Write the class name on the video frame.
    cv2.putText(rgb_frame, selected_class_Name, (10, 30), cv2.FONT_HERSHEY_SIMPLEX, 1, (255, 255, 255), 2)

    # Display the frame.
    plt.subplot(5, 4, counter);plt.imshow(rgb_frame);plt.axis('off')
```



✓ Step 2: Preprocess the Dataset

Next, we will perform some preprocessing on the dataset. First, we will read the video files from the dataset and resize the frames of the videos to a fixed width and height, to reduce the computations and normalized the data to range $[0-1]$ by dividing the pixel values with 255, which makes convergence faster while training the network.

But first, let's initialize some constants.

```
# Specify the height and width to which each video frame will be resized in our dataset.
IMAGE_HEIGHT , IMAGE_WIDTH = 64, 64

# Specify the number of frames of a video that will be fed to the model as one sequence.
SEQUENCE_LENGTH = 20

# Specify the directory containing the UCF50 dataset.
DATASET_DIR = "UCF50"

# Specify the list containing the names of the classes used for training. Feel free to choose any set of classes.
CLASSES_LIST = ["WalkingWithDog", "TaiChi", "Swing", "HorseRace"]
```

Note: The *IMAGE_HEIGHT*, *IMAGE_WIDTH* and *SEQUENCE_LENGTH* constants can be increased for better results, although increasing the sequence length is only effective to a certain point, and increasing the values will result in the process being more computationally expensive.

✓ Create a Function to Extract, Resize & Normalize Frames

We will create a function `frames_extraction()` that will create a list containing the resized and normalized frames of a video whose path is passed to it as an argument. The function will read the video file frame by frame, although not all frames are added to the list as we will only need an evenly distributed sequence length of frames.

```

def frames_extraction(video_path):
    '''
    This function will extract the required frames from a video after resizing and normalizing them.
    Args:
        video_path: The path of the video in the disk, whose frames are to be extracted.
    Returns:
        frames_list: A list containing the resized and normalized frames of the video.
    '''

    # Declare a list to store video frames.
    frames_list = []

    # Read the Video File using the VideoCapture object.
    video_reader = cv2.VideoCapture(video_path)

    # Get the total number of frames in the video.
    video_frames_count = int(video_reader.get(cv2.CAP_PROP_FRAME_COUNT))

    # Calculate the the interval after which frames will be added to the list.
    skip_frames_window = max(int(video_frames_count/SEQUENCE_LENGTH), 1)

    # Iterate through the Video Frames.
    for frame_counter in range(SEQUENCE_LENGTH):

        # Set the current frame position of the video.
        video_reader.set(cv2.CAP_PROP_POS_FRAMES, frame_counter * skip_frames_window)

        # Reading the frame from the video.
        success, frame = video_reader.read()

        # Check if Video frame is not successfully read then break the loop
        if not success:
            break

        # Resize the Frame to fixed height and width.
        resized_frame = cv2.resize(frame, (IMAGE_HEIGHT, IMAGE_WIDTH))

        # Normalize the resized frame by dividing it with 255 so that each pixel value then lies between 0 and 1
        normalized_frame = resized_frame / 255

        # Append the normalized frame into the frames list
        frames_list.append(normalized_frame)

    # Release the VideoCapture object.
    video_reader.release()

    # Return the frames list.
    return frames_list

```

✓ Create a Function for Dataset Creation

Now we will create a function `create_dataset()` that will iterate through all the classes specified in the `CLASSES_LIST` constant and will call the function `frame_extraction()` on every video file of the selected classes and return the frames (`features`), class index (`labels`), and video file path (`video_files_paths`).

```
def create_dataset():
    '''
    This function will extract the data of the selected classes and create the required dataset.
    Returns:
        features:      A list containing the extracted frames of the videos.
        labels:        A list containing the indexes of the classes associated with the videos.
        video_files_paths: A list containing the paths of the videos in the disk.
    '''

    # Declared Empty Lists to store the features, labels and video file path values.
    features = []
    labels = []
    video_files_paths = []

    # Iterating through all the classes mentioned in the classes list
    for class_index, class_name in enumerate(CLASSES_LIST):

        # Display the name of the class whose data is being extracted.
        print(f'Extracting Data of Class: {class_name}')

        # Get the list of video files present in the specific class name directory.
        files_list = os.listdir(os.path.join(DATASET_DIR, class_name))

        # Iterate through all the files present in the files list.
        for file_name in files_list:

            # Get the complete video path.
            video_file_path = os.path.join(DATASET_DIR, class_name, file_name)

            # Extract the frames of the video file.
            frames = frames_extraction(video_file_path)

            # Check if the extracted frames are equal to the SEQUENCE_LENGTH specified above.
            # So ignore the vides having frames less than the SEQUENCE_LENGTH.
            if len(frames) == SEQUENCE_LENGTH:

                # Append the data to their repective lists.
                features.append(frames)
                labels.append(class_index)
                video_files_paths.append(video_file_path)

    # Converting the list to numpy arrays
    features = np.asarray(features)
    labels = np.array(labels)

    # Return the frames, class index, and video file path.
    return features, labels, video_files_paths
```

Now we will utilize the function `create_dataset()` created above to extract the data of the selected classes and create the required dataset.

```
# Create the dataset.
features, labels, video_files_paths = create_dataset()

Extracting Data of Class: WalkingWithDog
Extracting Data of Class: TaiChi
Extracting Data of Class: Swing
Extracting Data of Class: HorseRace
```

Now we will convert `labels` (class indexes) into one-hot encoded vectors.

```
# Using Keras's to_categorical method to convert labels into one-hot-encoded vectors
one_hot_encoded_labels = to_categorical(labels)
```

✓ Step 3: Split the Data into Train and Test Set

As of now, we have the required `features` (a NumPy array containing all the extracted frames of the videos) and `one_hot_encoded_labels` (also a Numpy array containing all class labels in one hot encoded format). So now, we will split our data to create training and testing sets. We will also shuffle the dataset before the split to avoid any bias and get splits representing the overall distribution of the data.

```

# Split the Data into Train ( 75% ) and Test Set ( 25% ).
features_train, features_test, labels_train, labels_test = train_test_split(features, one_hot_encoded_labels,
                                                                              test_size = 0.25, shuffle = True,
                                                                              random_state = seed_constant)

def create_conv_lstm_model():
    """
    This function will construct the required conv_lstm model.
    Returns:
        model: It is the required constructed conv_lstm model.
    """

    # We will use a Sequential model for model construction
    model = Sequential()

    # Define the Model Architecture.
    #####

    model.add(ConvLSTM2D(filters = 4, kernel_size = (3, 3), activation = 'tanh', data_format = "channels_last",
                        recurrent_dropout=0.2, return_sequences=True, input_shape = (SEQUENCE_LENGTH,
                                                                                     IMAGE_HEIGHT, IMAGE_WIDTH, 3)))

    model.add(MaxPooling3D(pool_size=(1, 2, 2), padding='same', data_format='channels_last'))
    model.add(TimeDistributed(Dropout(0.2)))

    model.add(ConvLSTM2D(filters = 8, kernel_size = (3, 3), activation = 'tanh', data_format = "channels_last",
                        recurrent_dropout=0.2, return_sequences=True))

    model.add(MaxPooling3D(pool_size=(1, 2, 2), padding='same', data_format='channels_last'))
    model.add(TimeDistributed(Dropout(0.2)))

    model.add(ConvLSTM2D(filters = 14, kernel_size = (3, 3), activation = 'tanh', data_format = "channels_last",
                        recurrent_dropout=0.2, return_sequences=True))

    model.add(MaxPooling3D(pool_size=(1, 2, 2), padding='same', data_format='channels_last'))
    model.add(TimeDistributed(Dropout(0.2)))

    model.add(ConvLSTM2D(filters = 16, kernel_size = (3, 3), activation = 'tanh', data_format = "channels_last",
                        recurrent_dropout=0.2, return_sequences=True))

    model.add(MaxPooling3D(pool_size=(1, 2, 2), padding='same', data_format='channels_last'))
    #model.add(TimeDistributed(Dropout(0.2)))

    model.add(Flatten())

    model.add(Dense(len(CLASSES_LIST), activation = "softmax"))

    #####

    # Display the models summary.
    model.summary()

    # Return the constructed conv_lstm model.
    return model

```

Now we will utilize the function `create_conv_lstm_model()` created above, to construct the required `conv_lstm` model.

```

# Construct the required conv_lstm model.
conv_lstm_model = create_conv_lstm_model()

# Display the success message.
print("Model Created Successfully!")

```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
conv_lstm2d (ConvLSTM2D)	(None, 20, 62, 62, 4)	1024
max_pooling3d (MaxPooling3D)	(None, 20, 31, 31, 4)	0
time_distributed (TimeDistributed)	(None, 20, 31, 31, 4)	0
conv_lstm2d_1 (ConvLSTM2D)	(None, 20, 29, 29, 8)	3488


```

max_pooling3d_1 (MaxPooling (None, 20, 15, 15, 8) 0
3D)

time_distributed_1 (TimeDis (None, 20, 15, 15, 8) 0
tributed)

conv_lstm2d_2 (ConvLSTM2D) (None, 20, 13, 13, 14) 11144

max_pooling3d_2 (MaxPooling (None, 20, 7, 7, 14) 0
3D)

time_distributed_2 (TimeDis (None, 20, 7, 7, 14) 0
tributed)

conv_lstm2d_3 (ConvLSTM2D) (None, 20, 5, 5, 16) 17344

max_pooling3d_3 (MaxPooling (None, 20, 3, 3, 16) 0
3D)

flatten (Flatten) (None, 2880) 0

dense (Dense) (None, 4) 11524

```

```

=====
Total params: 44,524
Trainable params: 44,524
Non-trainable params: 0
=====
Model Created Successfully!

```

✓ Check Model's Structure:

Now we will use the `plot_model()` function, to check the structure of the constructed model, this is helpful while constructing a complex network and making that the network is created correctly.

```
# Create an Instance of Early Stopping Callback
```

```
early_stopping_callback = EarlyStopping(monitor = 'val_loss', patience = 10, mode = 'min', restore_best_weights = True)
```

```
# Compile the model and specify loss function, optimizer and metrics values to the model
```

```
convlstm_model.compile(loss = 'categorical_crossentropy', optimizer = 'Adam', metrics = ["accuracy"])
```