

Natural Network Project Report

April 26, 2020

Yijun Zhang
629009140

Abstract

For smart home design, fall detection is very important. In this project, I study the vision-based solution using Convolutional Neural Network and Long Short-Term Memory Network to decide if a frame contains a person falling. The figure shows how the probability of falling changes according to the time of the video. The result shows that using the body landmark as input, the Convolutional Neural Network model can achieve accuracy of 0.9612 and Long Short-Term Memory Network could achieve accuracy of 0.9333. The final ensemble model has accuracy 0.9611.

1 Dataset

There are 3 different data sets namely – ‘UR’, ‘MC’, ‘YOUTUBE’. There are total 484 videos, including 252 videos with fall and 232 videos without fall.

1.1 UR Fall Detection Dataset^[1] (UR)

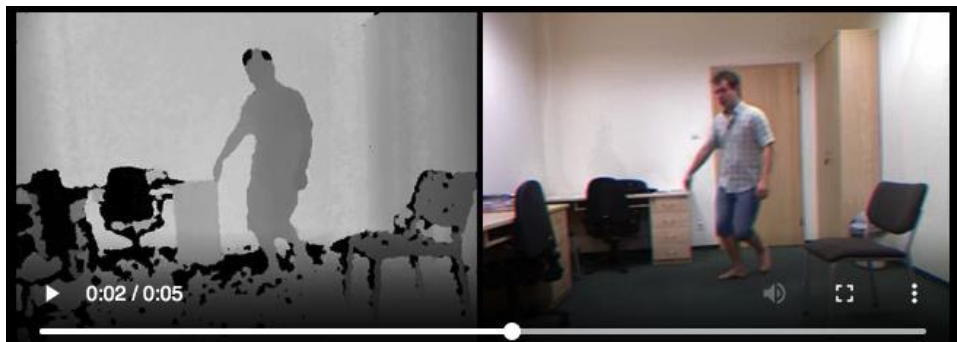


Figure 1: Example of video in UR Fall Detection Dataset

This dataset contains 70 (30 falls + 40 activities of daily living) sequences. Fall events are recorded with 2 Microsoft Kinect cameras(camera 0 from side angle and camera 1 from directly above angle) and corresponding accelerometric data. ADL events are recorded with only one device (camera 0). It also includes label files: [urfall-cam0-falls.csv](#) and [urfall-cam0-adls.csv](#). In these 2 csv files, each row contains one sample of data corresponding to one image frame of a video.

I use videos recorded by camera 0. There are 3 reasons:

- It is from the side angle, which is easier for later body landmark generation.
- Using videos recorded by camera 0 also because it is more similar to the camera set in smart home.
- ADL events are recorded with only one device (camera 0).

1.2 Multiple Camera Dataset_[2] (MC)



Figure 2: Example of video in Multiple Camera Dataset

This dataset contains 24 scenarios recorded with 8 IP video cameras. The first 22 first scenarios contain a fall and confounding events; the last 2 ones contain only confounding events. There is a total of 192 videos.

Each video may contain some staff and a person to perform falling action. For the accuracy of extracting body landmarks, I trim each video to only contain one person.

This dataset does not contain any label file. Therefore, I trim each video into 2 parts: one only contains frames of a person from falling to laying down, the other only contains frames of a person doing other things before falling. Later, I mark all frames of the first part to be fall(1) and marked all frames of the second part to be not-fall(0).

1.3 YOUTUBE Dataset

I download videos contains fall from YouTube. There are 23 videos including coach teaching old people how to fall safely, elderly people fall in the hospital, elderly people fall at home, people fall outdoor. For later labeling, I trim each video into 2 parts: one only contains frames of a person from falling to laying down, the other only contains frames of a person doing other things before falling. Later, I mark all frames of the first part to be fall(1) and marked all frames of the second part to be not-fall(0).

1.4 Variety of Videos

A person (such as elderly person, patient) may fall in various ways. There are many actions similar to fall actions. To correctly detect different kinds of fall and to reduce false positive, I consider the variety of videos in my dataset.

Videos in my datasets that contain fall action includes different kinds of fall. For daily activity videos that does not contains fall, there are different kinds of confusion actions.

Table 1: Types of fall videos

Type of fall	Number of video
Includes elderly person	10
Fall forward	187

Fall backward	25
Fall sideward	40
Outdoor	8
Indoor	244
Kick something then fall	10
Dizzy then fall	2

Table 2: Confusion Action of videos

Confusion Action	Number of video
Find or pick something up on the ground	31
Lay down in sofa	10
Squat	21
Lay down in bed	24
Sit in chair/sofa	35
Tie shoelaces	9
Do push up	8

2 Data Processing

2.1 Body Landmark Extraction

I use the open source tool^{[3][4][5][6]} from CMU to produce body landmark for frames of each video. The website of this tool: <https://github.com/CMU-Perceptual-Computing-Lab/openpose>.

All body landmarks for each frame are stored in a json file. The generated body landmarks for each frame contains 25 key points: {0, "Nose"}, {1, "Neck"}, {2, "RShoulder"}, {3, "RElbow"}, {4, "RWrist"}, {5, "LShoulder"}, {6, "LElbow"}, {7, "LWrist"}, {8, "MidHip"}, {9, "RHip"}, {10, "RKnee"}, {11, "RAnkle"}, {12, "LHip"}, {13, "LKnee"}, {14, "LAnkle"}, {15, "REye"}, {16, "LEye"}, {17, "REar"}, {18, "LEar"}, {19, "LBigToe"}, {20, "LSmallToe"}, {21, "LHeel"}, {22, "RBigToe"}, {23, "RSmallToe"}, {24, "RHeel"}, {25, "Background"}.

For fall detection, I think body landmarks of the face are unnecessary: {15, "REye"}(right eye), {16, "LEye"}(left eye), {17, "REar"}(right ear), {18, "LEar"}(left ear). After reading body landmarks from JSON files, I remove these 4 key points for each frame.

In case of unstable connection of Google Colab or crash of local machine. I store body landmarks into CSV files so that later I can read these data from CSV files, which will be faster than re-read all JSON files.

2.2 Labeling

UR dataset has its label files for each frame. Labels are stored in csv files. Attribute in csv files includes:

- sequence name - camera name is omitted, because all of the samples are from the front camera ('fall-01-cam0-d' is 'fall-01', 'adl-01-cam0-d' is 'adl-01' and so on).
- frame number - corresponding to number in sequence.
- label - describes human posture in the depth frame; '-1' means person is not lying, '1' means person is lying on the ground; '0' is temporary pose, when person "is falling", we don't use '0' frames in classification.

I transform the label using the rule:

- For the falling videos, original label '-1' means 0 (not falling), and original label '0' and '1' means 1 (falling).
- For the activities of daily living videos, all frames are labeled 0 (not falling).

Videos in MC dataset and YOUTUBE dataset are manually trimmed to 2 parts: one only contains frames of a person from falling to laying down, the other only contains frames of a person doing other things before falling. I mark all frames of the first part to be fall(1) and marked all frames of the second part to be not-fall(0).

2.3 Normalization and Up-Sampling

I normalize data.

```
1 # normalize data
2 for record in bodylandmark:
3     for i in range(0, 63, 3):
4         record[i] = float(record[i]) / 640
5         record[i + 1] = float(record[i + 1]) / 480
6
7 print(bodylandmark[0])
8 print(len(label))
9 print(len(bodylandmark))
```

```
[0.81401875, 0.2670375, 0.833728, 0.8508921875000001, 0.2724625, 0.739747, 0.844743
33678
33678
```

Figure 3: Annotated code of data normalization.

There are 20661 frames are not fall(0) and 13017 frames of fall(1). To make the data balance, I used up-sampling. After up-sampling, both fall frames and not fall frames are 20661.

Table 3: Frames before and after Up-sampling

	Fall frames(1)	Not fall frames (0)
Before up-sampling	13017	20661
After up-sampling	20661	20661

```

1 # up-sampling to make the data balance
2 # Separate majority and minority classes
3 df = pd.DataFrame.from_records(all_data)
4 header = ['label']
5 for i in range(63):
6     header.append(i)
7
8 df.columns = header
9 df_majority = df[df.label==0]
10 df_minority = df[df.label==1]
11 print("before re-sampling, not fall vs fall: ")
12 print(df['label'].value_counts())
13
14 from sklearn.utils import resample
15 print("begin to re sample...")
16 # Upsample minority class
17 df_minority_upsampled = resample(df_minority,
18                                 replace=True,          # sample with replacement
19                                 n_samples=20661)        # to match majority class
20
21 # Combine majority class with upsampled minority class
22 df_upsampled = pd.concat([df_majority, df_minority_upsampled])
23
24 print("after re-sampling...")
25 # Display new class counts
26 df_upsampled.label.value_counts()

```

```

before re-sampling, not fall vs fall:
0    20661
1    13017
Name: label, dtype: int64
begin to re sample...
after re-sampling...
1    20661
0    20661
Name: label, dtype: int64

```

Figure 4: Annotated code of up-sampling.

2.4 Train-Test set Split

I randomly shuffle the data to ensure that each data point creates an "independent" change on the model, without being biased by the same points before them.

Then, I split the dataset into training set and test set.

```
1 # split data into train and test set
2 np.random.shuffle(all_data)
3 split_point = len(all_data) // 9
4 print("split point:", split_point)
5 test_data = all_data[:split_point]
6 train_data = all_data[split_point:]
7
8 print("number of test data:", len(test_data))
9 print("number of train data:", len(train_data))
10
11 print(train_data[0])
```

Figure 5: Annotated code of train-test set split.

3 Model Training

3.1 LSTM model Training

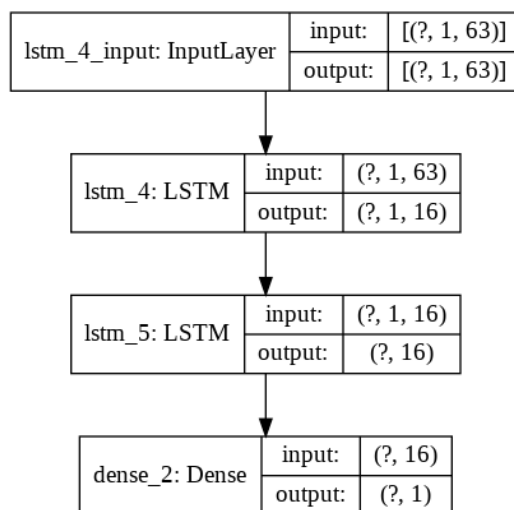


Figure 6: Shape of tensor of each layer of LSTM model

Table 4: Shape of tensor of the first and last layer of LSTM model

Input shape of tensor	x_train: (36731, 1, 63)
	x_test: (4591, 1, 63)
Output shape of tensor	y_train shape is (36731,16)
	y_test shape is (4591,1)

The LSTM model is built on a linear stack of layers with the sequential model. There are two hidden layers in total.

Model: "sequential_3"

Layer (type)	Output Shape	Param #
lstm_4 (LSTM)	(None, 1, 16)	5120
lstm_5 (LSTM)	(None, 16)	2112
dense_4 (Dense)	(None, 1)	17
Total params: 7,249		
Trainable params: 7,249		
Non-trainable params: 0		

Figure 7: Summary of LSTM model.

For learning rate, I define a scheduler. When epoch is less than 10, the learning rate will be 0.01. When epoch is more than 10, it will decrease as epoch increase.

I apply early stopping with a patience 5 epochs, which means if the loss is not improving for 5 epochs, the training will stop. This policy can avoid unnecessary full training.

I try different hyperparameters, and combination of the best performance is Batch size 16, epochs 52 and without dropout.

Table 5: Performance comparison of different batch size

Batch size	16	32	64
Loss	0.2491	0.2668	0.2962
Accuracy	0.8883	0.8789	0.8624

Best batch size is 16.

Table 6: Performance comparison of different drop out

Drop out	none	0.1	0.5
Loss	0.2491	0.2884	0.2544
Accuracy	0.8883	0.8789	0.8801

The performance does not improve, so I removed dropout.

According to the above discussion, hyperparameters of LSTM model are as follows:

Table 7: Hyperparameters of LSTM model

Batch Size	16
Epochs	52
Dropout	(unnecessary)

A snippet of code of constructing the model/network is shown as follows:

```

8 def scheduler(epoch):
9     if epoch < 10:
10         return 0.001
11     else:
12         return 0.001 * np.exp(0.1 * (10 - epoch))
13
14 x_train = []
15 y_train = []
16
17 for record in train_data:
18     x_train.append(record[1:])
19     label = int(record[0])
20     y_train.append(label)
21
22 def train_lstm_model(x_train, y_train):
23     x_train = array(x_train)
24     x_train = x_train.reshape((len(x_train), 1, len(x_train[0])))
25     print("x_train.shape", x_train.shape)
26     print(x_train[0])
27
28     y_train = array(y_train)
29     print("y_train.shape", y_train.shape)
30
31     # improve log: use batch size 16 and add one more lstm layer
32
33     lstm_model = Sequential()
34     lstm_model.add(LSTM(16,
35                         input_shape=(1, 63),
36                         return_sequences=True))
37     lstm_model.add(LSTM(16, ))
38     lstm_model.add(layers.Dense(1, activation='sigmoid'))
39     lstm_model.compile(optimizer='rmsprop',
40                       loss='binary_crossentropy',
41                       metrics=['acc',
42                               metrics.AUC(),
43                               metrics.FalseNegatives(),
44                               metrics.Recall(),
45                               metrics.Precision(),
46                               metrics.FalseNegatives(),
47                               metrics.TrueNegatives(),
48                               metrics.FalsePositives(),
49                               metrics.TruePositives()])
50     lstm_history = lstm_model.fit(x_train, y_train,
51                                 epochs=100,
52                                 batch_size=16,
53                                 validation_split=0.2,
54                                 callbacks=[callbacks.EarlyStopping(monitor='val_loss', patience=5),
55                                           callbacks.LearningRateScheduler(scheduler)])
56     print("finish training lstm model")
57     return lstm_model, lstm_history
58
59 lstm_model, lstm_history = train_lstm_model(x_train, y_train)

```

Figure 8: Annotated code of LSTM model training

Training and Testing performance of LSTM model are as follows:

```
Epoch 1/100
1837/1837 [=====] - 23s 13ms/step - loss: 0.3716 - acc: 0.8316 - auc_4: 0.8498 .
Epoch 2/100
1837/1837 [=====] - 23s 13ms/step - loss: 0.2633 - acc: 0.8903 - auc_4: 0.9312 .
Epoch 3/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.2438 - acc: 0.8981 - auc_4: 0.9444 .
Epoch 4/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.2304 - acc: 0.9052 - auc_4: 0.9502 .
Epoch 5/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.2219 - acc: 0.9104 - auc_4: 0.9544 .
Epoch 6/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.2139 - acc: 0.9124 - auc_4: 0.9576 .
Epoch 7/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.2084 - acc: 0.9158 - auc_4: 0.9599 .
Epoch 8/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.2023 - acc: 0.9180 - auc_4: 0.9613 .
Epoch 9/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1979 - acc: 0.9190 - auc_4: 0.9629 .
Epoch 10/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1931 - acc: 0.9221 - auc_4: 0.9643 .
Epoch 11/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1879 - acc: 0.9228 - auc_4: 0.9656 .
Epoch 12/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1850 - acc: 0.9248 - auc_4: 0.9667 .
Epoch 13/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1805 - acc: 0.9264 - auc_4: 0.9677 .
Epoch 14/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1773 - acc: 0.9278 - auc_4: 0.9686 .
Epoch 15/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1728 - acc: 0.9291 - auc_4: 0.9694 .
Epoch 16/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1707 - acc: 0.9302 - auc_4: 0.9701 .
Epoch 17/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1678 - acc: 0.9315 - auc_4: 0.9709 .
Epoch 18/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1654 - acc: 0.9326 - auc_4: 0.9715 .
Epoch 19/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1636 - acc: 0.9337 - auc_4: 0.9721 .
Epoch 20/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1611 - acc: 0.9334 - auc_4: 0.9727 .
Epoch 21/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1598 - acc: 0.9348 - auc_4: 0.9732 .
Epoch 22/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1583 - acc: 0.9344 - auc_4: 0.9737 .
Epoch 23/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1570 - acc: 0.9355 - auc_4: 0.9742 .
Epoch 24/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1554 - acc: 0.9363 - auc_4: 0.9746 .
Epoch 25/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1540 - acc: 0.9371 - auc_4: 0.9750 .
Epoch 26/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1536 - acc: 0.9373 - auc_4: 0.9754 .
Epoch 27/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1528 - acc: 0.9378 - auc_4: 0.9757 .
Epoch 28/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1512 - acc: 0.9384 - auc_4: 0.9761 .
Epoch 29/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1511 - acc: 0.9378 - auc_4: 0.9764 .
Epoch 30/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1502 - acc: 0.9386 - auc_4: 0.9767 .
Epoch 31/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1498 - acc: 0.9385 - auc_4: 0.9769 .
Epoch 32/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1490 - acc: 0.9394 - auc_4: 0.9772 .
Epoch 33/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1488 - acc: 0.9393 - auc_4: 0.9775 .
Epoch 34/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1483 - acc: 0.9393 - auc_4: 0.9777 .
Epoch 35/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1476 - acc: 0.9394 - auc_4: 0.9779 .
Epoch 36/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1473 - acc: 0.9398 - auc_4: 0.9781 .
Epoch 37/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1471 - acc: 0.9402 - auc_4: 0.9783 .
Epoch 38/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1469 - acc: 0.9399 - auc_4: 0.9785 .
Epoch 39/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1464 - acc: 0.9401 - auc_4: 0.9787 .
Epoch 40/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1465 - acc: 0.9397 - auc_4: 0.9789
```

```

Epoch 41/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1462 - acc: 0.9401 - auc_4: 0.9790
Epoch 42/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1459 - acc: 0.9404 - auc_4: 0.9792
Epoch 43/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1458 - acc: 0.9406 - auc_4: 0.9793
Epoch 44/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1456 - acc: 0.9401 - auc_4: 0.9795
Epoch 45/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1458 - acc: 0.9404 - auc_4: 0.9796
Epoch 46/100
1837/1837 [=====] - 21s 12ms/step - loss: 0.1454 - acc: 0.9402 - auc_4: 0.9797
Epoch 47/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1452 - acc: 0.9404 - auc_4: 0.9799
Epoch 48/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1450 - acc: 0.9410 - auc_4: 0.9800
Epoch 49/100
1837/1837 [=====] - 22s 12ms/step - loss: 0.1450 - acc: 0.9405 - auc_4: 0.9801
Epoch 50/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1449 - acc: 0.9407 - auc_4: 0.9802
Epoch 51/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1448 - acc: 0.9408 - auc_4: 0.9803
Epoch 52/100
1837/1837 [=====] - 21s 11ms/step - loss: 0.1447 - acc: 0.9407 - auc_4: 0.9804
finish training lstm model

```

Figure 9: Training log of LSTM model training

```

1148/1148 [=====] - 7s 6ms/step - loss: 0.1497 - acc: 0.9395 - auc_4: 0.9805 - false_negatives_7: 49462.9883 - recall_4: 0.9486 - precision_4: 0.9125 - false_negatives_8: 49462.9883
[0.14972610771656036, 0.9394516944885254, 0.9805395603179932, 49462.98828125, 0.9485976696014404, 0.9125279188156128, 49462.98828125, 878713.0625, 87488.15625, 912731.1875]

```

Figure 10: Testing result of LSTM model

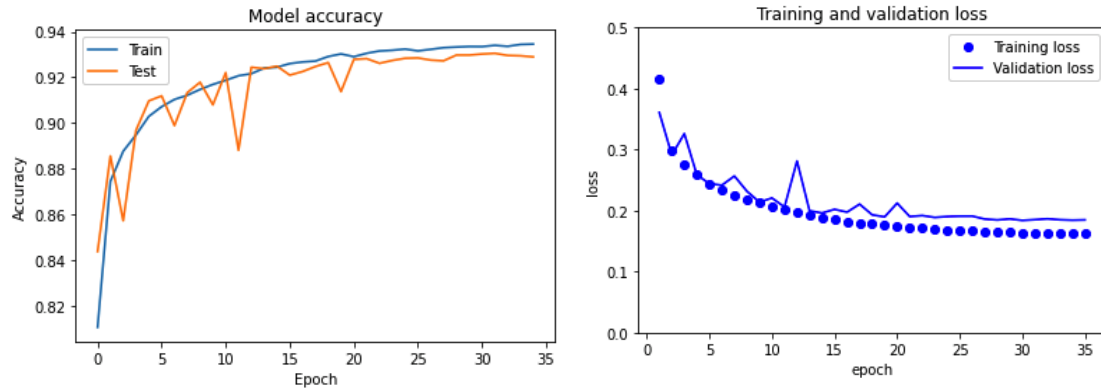


Figure 11: Performance of LSTM model training

Table 8: Test result of LSTM model

Testing Loss	0.1497
Testing Accuracy	0.9395
AUC	0.9805
Sensitivity /Recall ($Sensitivity = \frac{TP}{TP+FN}$)	0.9486
Specificity ($Specificity = \frac{TN}{TN+FP}$)	0.9095

After training the model and use test dataset to get the test score, I trained the final LSTM model using all the data.

```

1 # train the final lstm model with all data (include test data)
2
3 x_train_final = []
4 y_train_final = []
5
6 for record in all_data:
7     x_train_final.append(record[1:])
8     label = int(record[0])
9     y_train_final.append(label)
10
11 final_lstm_model, final_lstm_history = train_lstm_model(x_train_final, y_train_final)
12
13 final_lstm_model.save("/content/drive/My Drive/636project/final_lstm_model.h5")
14

```

Figure 12: Annotated code of LSTM model training using all data

3.2 CNN model Training

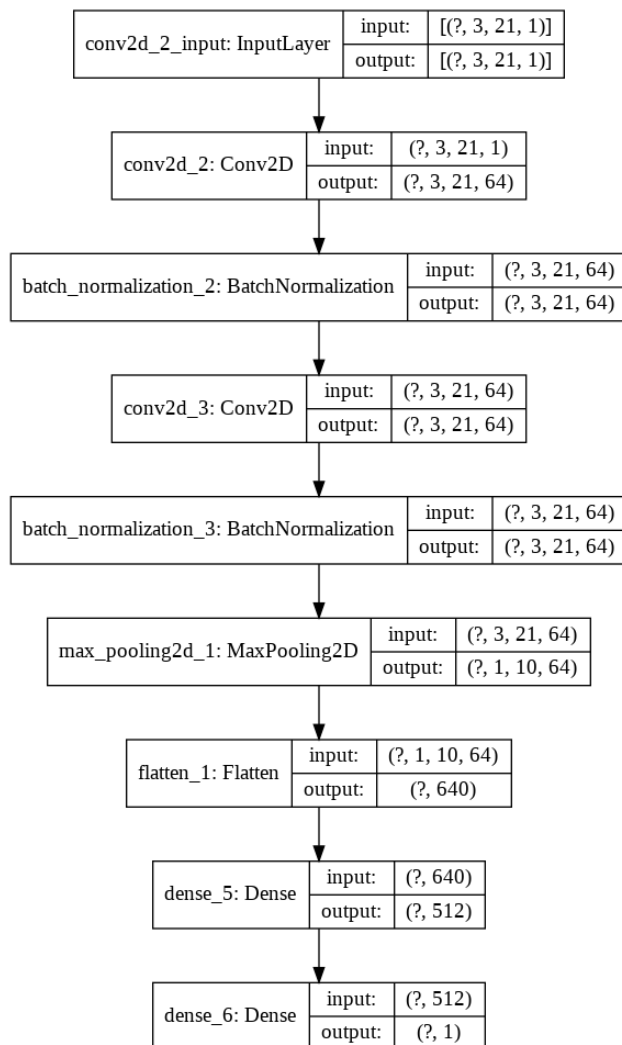


Figure 13: Shape of tensor of each layer of CNN model

Table 9: Shape of tensor of the first and last layer of CNN model

Input shape of tensor	x_train: (36731, 3, 21, 63)
	x_test: (4591, 3, 21, 1)
Output shape of tensor	y_train shape is (36731,512)
	y_test shape is (4591,1)

The CNN model is built on a linear stack of layers with the sequential model. There are 2 2D-convolution layers, 1 max pooling layers and 2 dense layers.

I added batch normalization between 2 convolution layers to adaptively normalize data as the mean and variance change over time during training.

Model: "sequential_4"

Layer (type)	Output Shape	Param #
conv2d_2 (Conv2D)	(None, 3, 21, 64)	640
batch_normalization_2 (Batch Normalization)	(None, 3, 21, 64)	12
conv2d_3 (Conv2D)	(None, 3, 21, 64)	36928
batch_normalization_3 (Batch Normalization)	(None, 3, 21, 64)	12
max_pooling2d_1 (MaxPooling2D)	(None, 1, 10, 64)	0
flatten_1 (Flatten)	(None, 640)	0
dense_5 (Dense)	(None, 512)	328192
dense_6 (Dense)	(None, 1)	513
Total params: 366,297		
Trainable params: 366,285		
Non-trainable params: 12		

Figure 14: Summary of CNN model.

According to the above discussion, Hyperparameters of CNN model are as follows:

Table 10: Hyperparameters of CNN model

Batch Size	16
Epochs	18
Dropout	(unnecessary)

For the learning rate, I use the same scheduler as LSTM model. When the epoch is less than 10, the learning rate will be 0.01. When the epoch is more than 10, it will decrease as epoch increases.

I apply early stopping with a patience 5 epochs, which means if the loss is not improving for 5 epochs, the training will stop. This policy can avoid unnecessary full training.

I try different hyperparameters, and the combination of the best performance is Batch size 16, epochs 22 and without dropout.

Table 11: Performance comparison of different batch size

Batch size	16	32	64
Loss	0.0667	0.0713	0.0812
Accuracy	0.9692	0.9621	0.9343

Best batch size is 16.

Table 12: Performance comparison of different drop out

Drop out	none	0.1	0.5
Loss	0.0667	0.0877	0.0670
Accuracy	0.9692	0.9601	0.9652

The performance does not improve, so I remove dropout.

Regarding the optimization of the network parameters, I found Adam is better than RMSprop here.

Table 13: Performance comparison of different optimizer

Optimizer	Adam	RMSprop
Loss	0.0667	0.0877
Accuracy	0.9692	0.9601

A snippet of code of constructing the model/network is shown as follows:

```

19 def train_cnn_model(x_train, y_train):
20     x_train = array(x_train)
21     x_train = x_train.reshape((len(x_train), 3, int(len(x_train[0])/3), 1))
22
23     y_train = array(y_train)
24
25     #create model
26     cnn_model = Sequential()
27     cnn_model.add(Conv2D(64,
28                         kernel_size=3,
29                         activation='relu',
30                         input_shape=(3,21,1),
31                         padding='same'))
32     cnn_model.add(layers.BatchNormalization(1))
33     cnn_model.add(Conv2D(64,
34                         kernel_size=3,
35                         activation='relu',
36                         padding='same'))
37     cnn_model.add(layers.BatchNormalization(1))
38     cnn_model.add(MaxPooling2D(2,2))
39     cnn_model.add(Flatten())
40     cnn_model.add(Dense(512, activation = 'relu'))
41     cnn_model.add(Dense(1, activation='sigmoid'))
42

```

```

43 # compile and fit
44 cnn_model.compile(optimizer='Adam',
45                   loss='binary_crossentropy',
46                   metrics=['acc',
47                           metrics.AUC(),
48                           metrics.FalseNegatives(),
49                           metrics.Recall(),
50                           metrics.Precision(),
51                           metrics.FalseNegatives(),
52                           metrics.TrueNegatives(),
53                           metrics.FalsePositives(),
54                           metrics.TruePositives()])
55 cnn_history = cnn_model.fit(x_train, y_train,
56                             epochs=100,
57                             batch_size=16,
58                             validation_split=0.2,
59                             callbacks=[callbacks.EarlyStopping(monitor='val_loss', patience=5),
60                                       callbacks.LearningRateScheduler(scheduler)])
61
62 print("finish training cnn model")
63 return cnn_model, cnn_history
64
65
66 cnn_model, cnn_history = train_cnn_model(x_train, y_train)

```

Figure 15: Annotated code of CNN model training

Training and Testing performance of CNN models are as follows:

```

Epoch 1/100
1837/1837 [=====] - 17s 9ms/step - loss: 0.2590 - acc: 0.8909 - auc_5: 0.9174
Epoch 2/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.1693 - acc: 0.9295 - auc_5: 0.9688
Epoch 3/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.1472 - acc: 0.9385 - auc_5: 0.9749
Epoch 4/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.1283 - acc: 0.9454 - auc_5: 0.9790
Epoch 5/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.1174 - acc: 0.9499 - auc_5: 0.9816
Epoch 6/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.1068 - acc: 0.9536 - auc_5: 0.9835
Epoch 7/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.1014 - acc: 0.9570 - auc_5: 0.9850
Epoch 8/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0946 - acc: 0.9586 - auc_5: 0.9862
Epoch 9/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0897 - acc: 0.9600 - auc_5: 0.9872
Epoch 10/100
1837/1837 [=====] - 17s 9ms/step - loss: 0.0869 - acc: 0.9618 - auc_5: 0.9880
Epoch 11/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0836 - acc: 0.9621 - auc_5: 0.9887
Epoch 12/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0814 - acc: 0.9630 - auc_5: 0.9893
Epoch 13/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0751 - acc: 0.9654 - auc_5: 0.9898
Epoch 14/100
1837/1837 [=====] - 17s 9ms/step - loss: 0.0721 - acc: 0.9659 - auc_5: 0.9902
Epoch 15/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0685 - acc: 0.9679 - auc_5: 0.9907
Epoch 16/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0674 - acc: 0.9677 - auc_5: 0.9910
Epoch 17/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0660 - acc: 0.9687 - auc_5: 0.9914
Epoch 18/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0629 - acc: 0.9693 - auc_5: 0.9917
Epoch 19/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0612 - acc: 0.9701 - auc_5: 0.9920

```

```

Epoch 20/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0605 - acc: 0.9702 - auc_5: 0.9922
Epoch 21/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0609 - acc: 0.9703 - auc_5: 0.9924
Epoch 22/100
1837/1837 [=====] - 16s 9ms/step - loss: 0.0582 - acc: 0.9710 - auc_5: 0.9926
finish training cnn model

```

Figure 16: Training log of CNN model training

```

1148/1148 [=====] - 7s 6ms/step - loss: 0.0667 - acc: 0.9692 - auc_5: 0.9928 - false_negatives_9: 8458.3076 - recall_5: 0.9795 - precision_5: 0.9361 - false_negatives_10: 8458.3076 .
[0.0666593685746193, 0.9692357778549194, 0.9928162693977356, 8458.3076171875, 0.9794872999191284, 0.9361050128936768, 8458.3076171875, 386510.875, 27570.66796875, 403926.0625]

```

Figure 17: Testing result of CNN model

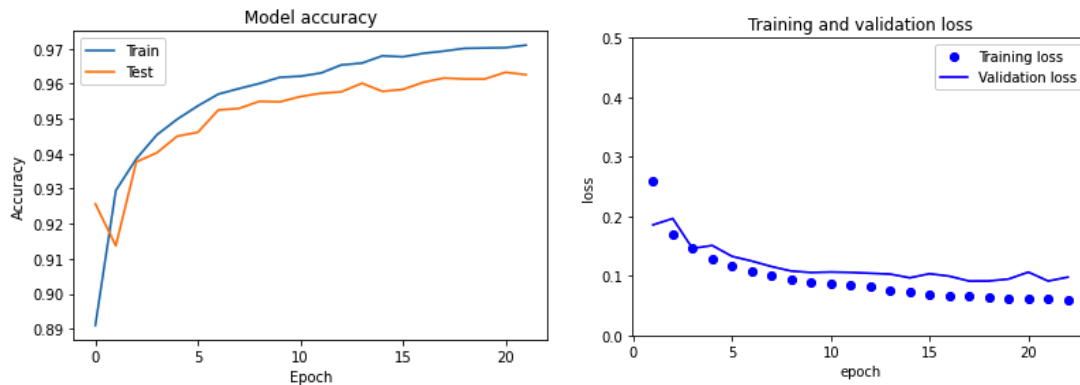


Figure 18: Performance of CNN model training

Table 14: Test result of CNN model

Testing Loss	0.0667
Testing Accuracy	0.9692
AUC	0.9928
Sensitivity / Recall ($Sensitivity = \frac{TP}{TP+FN}$)	0.9795
Specificity ($Specificity = \frac{TN}{TN+FP}$)	0.9332

After training the model and use test dataset to get the test score, I trained the final CNN model using all the data.

```

1 # train the final cnn model with all data (include test data)
2 x_train_final = []
3 y_train_final = []
4
5 for record in all_data:
6     x_train_final.append(record[1:])
7     label = int(record[0])
8     y_train_final.append(label)
9
10 final_cnn_model, final_cnn_history = train_cnn_model(x_train_final, y_train_final)
11
12 final_cnn_model.save("/content/drive/My Drive/636project/final_cnn_model.h5")

```

Figure 19: Annotated code of CNN model training using all data

3.3 Ensemble Model Training

The LSTM model and CNN model are good models and they prediction from different aspects. I ensemble these two models with different ratio($\text{prediction of ensemble model} = \text{prediction of CNN model} * \text{fraction1} + \text{prediction of LSTM model} * \text{fraction2}$, $\text{fraction1} + \text{fraction2} = 1$) and compare the performance($\text{Accuracy} = \frac{(\text{TruePositive} + \text{TrueNegative})}{\text{ALLframes}}$) of ensemble model and CNN model and LSTM model.

Table 15: Ratio of CNN model to LSTM model, Accuracy comparison

CNN vs. LSTM	Ensembled model	CNN model	LSTM model
0.1 : 0.9	0.9362491834	0.9611615891	0.9332501930
0.2 : 0.8	0.9399311122		
0.3 : 0.7	0.9440287427		
0.4 : 0.6	0.9496704079		
0.5 : 0.5	0.9591127739		
0.6 : 0.4	0.9606271156		
0.7 : 0.3	0.9609240453		
0.8 : 0.2	0.9611318961		
0.9 : 0.1	0.9611318961		

It is shown that the ensemble model does not have improved the accuracy of prediction. It seems unnecessary to do model ensembling. But in the last experiment(models trained without YOUTUBE dataset), with ratio 0.8 : 0.2(CNN vs. LSTM) the ensemble model have better performance(0.961132) compare to CNN model(0.961162) and LSTM model(0.933250). Also, for predictions of frames for a whole video, the ensemble model has better performance.

These are two test case for two of the test videos:

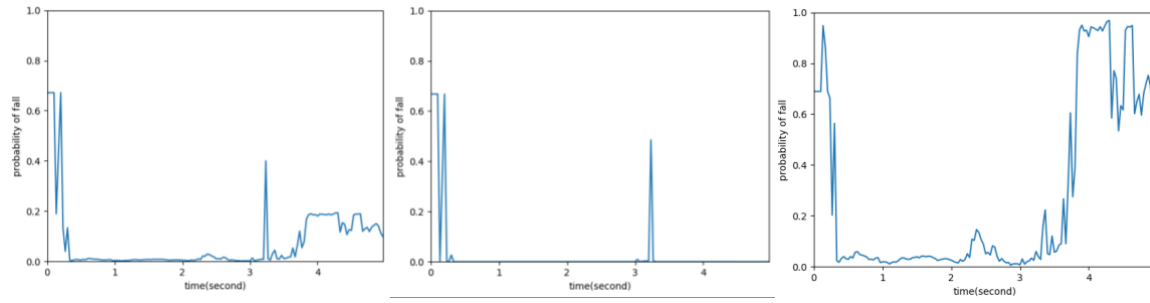


Figure 20: Prediction for video [UR-adl-01-cam0](#) (no fall) of different models

from left to right: ensembled model, CNN model, LSTM model

The above video does not contain fall action, we can see that in the third second, the ensembled model has a prediction of less than 0.5, while CNN model has a prediction of 0.5. In the fourth and fifth seconds, the ensembled model performs better than LSTM model.

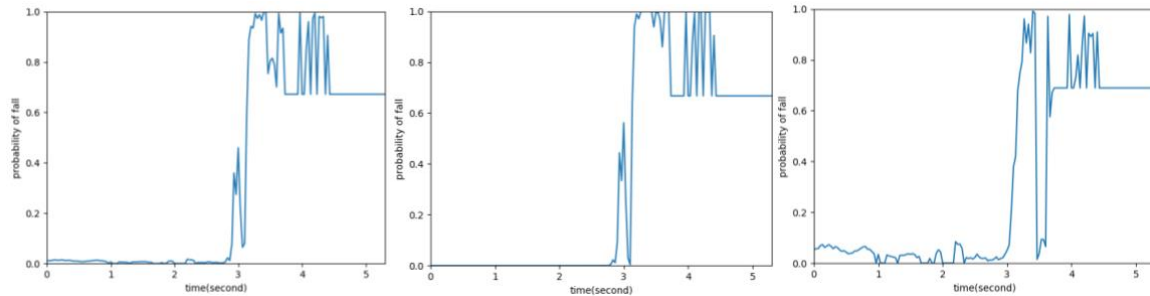


Figure 21: Prediction for video [UR-fall-01-cam0](#) (contains fall) of different models

from left to right: ensembled model, CNN model, LSTM model

In this video, in the first 2.5 seconds, the person does not fall, then he falls and lays down until the end. The prediction of the ensembled model is much more accurate than LSTM model.

Therefore, I decide to make the final model for prediction as the ensembled model.

This is the code of model ensembling:

```
1 # df = pd.DataFrame.from_records(all_data)
2 test_cnn = []
3 test_lstm = []
4 y_test = []
5 correct = 0
6 incorrect = 0
7 cnn_correct = 0
8 cnn_incorrect = 0
9 lstm_correct = 0
10 lstm_incorrect = 0
11
12 for record in all_data:
13     test_cnn.append(record[1:])
14     test_lstm.append(record[1:])
15     label = int(record[0])
16     y_test.append(label)
```

```

17
18
19 # use model and write the json file for a video
20 print("loading model and computing probability for each frame...")
21 cnn_model = keras.models.load_model('/content/drive/My Drive/636project/model_improved_cnn.h5')
22 lstm_model = keras.models.load_model('/content/drive/My Drive/636project/model_improved_lstm.h5')
23
24 # for cnn model, reshape and predict
25 test_cnn = array(test_cnn)
26 test_cnn = test_cnn.reshape((len(test_cnn), 3, int(len(test_cnn[0])/3), 1))
27 probability_cnn = cnn_model.predict_proba(test_cnn)
28
29 # for lstm model, reshape and predict
30 test_lstm = array(test_lstm)
31 test_lstm = test_lstm.reshape((len(test_lstm), 1, len(test_lstm[0])))
32 probability_lstm = lstm_model.predict_proba(test_lstm)
33
34 # ensemble the prediction of cnn and lstm model
35 frame_num = 0
36
37 # fp: false positive
38 # fn: false negative
39 fp = []
40 fn = []
41
42 for i in range(len(probability_cnn)):
43     timestamp = frame_num / 30
44     probability_fall = float(probability_cnn[i][0]) * 0.9 + float(probability_lstm[i][0]) * 0.1
45     # emsembled model
46     if probability_fall < 0.5:
47         prediction = 0.0
48     else:
49         prediction = 1.0
50
51     if prediction == float(y_test[i]):
52         correct += 1
53     else:
54         incorrect += 1
55         if prediction == 1.0:
56             fp.append(i)
57         else:
58             fn.append(i)
59
60     # cnn model
61     if probability_cnn[i][0] < 0.5:
62         prediction_cnn = 0.0
63     else:
64         prediction_cnn = 1.0
65
66     if prediction_cnn == float(y_test[i]):
67         cnn_correct += 1
68     else:
69         cnn_incorrect += 1
70
71     # lstm model
72     if probability_lstm[i][0] < 0.5:
73         prediction_lstm = 0.0
74     else:
75         prediction_lstm = 1.0
76
77     if prediction_lstm == float(y_test[i]):
78         lstm_correct += 1
79     else:
80         lstm_incorrect += 1
81
82 print('finish prediction')

```

Figure 22: Annotated code of model ensembling

4 Problems and Fail Cases

For all data I used, I use the final model(ensembled model) to get all fail cases. For a total of 33678 video frames, there are 1309 fail cases. In these fails frames, there are 1185 false-positive frames(actual fall but predicts not fall) 124 false-negative frames(actual not fall but predicts fall). For all these fail cases, I got a CSV file of their video number and frame number. I could analyze them by checking each frame.

```
1 # get name of fail frames
2 false_positive_frames = []
3 false_negative_frames = []
4
5 for fail in fp:
6     false_positive_frames.append(frame_name[fail])
7
8
9 for fail in fn:
10    false_negative_frames.append(frame_name[fail])
11
12 print(len(false_positive_frames))
13 print(len(false_negative_frames))

1185
124

1 # save names of fail frames
2 df_false_negative = pd.DataFrame.from_records(false_negative_frames)
3 df_false_negative.to_csv (r'/content/drive/My Drive/636project/false_negative_frames.csv', index = False, header=False)
4 df_false_positive = pd.DataFrame.from_records(false_positive_frames)
5 df_false_positive.to_csv (r'/content/drive/My Drive/636project/false_positive_frames.csv', index = False, header=False)
```

Figure 23: Annotated code of getting fail cases

By checking the image of fail video frames. I summarize main reasons of fail prediction.

Table 16: Reasons of fail prediction

Fail reason	Number of fail frames	Percentage
Uncertainty of falling action	610	46.6%
Confusing actions	431	32.9%
Uncertain reason	268	20.5%

First, the uncertainty of falling action. These fail cases happen before the fall and are consider as fall. The main reason for this failure is because the trimming and labeling of videos.

UR dataset has a label for each frame, but the Multiple Camera Dataset does not have. Then I trimmed each video into 2 parts: the first part is before fall(the person may walk, sit or stand) and the second part is falling and laying on the floor. I labeled all frames of the first part to be 0(not fall) and all frames of the second part to be 1(fall). The trimming is judge by a human, and the fall action is a consecutive action. Therefore, there is no exact time to trim the video.

For example, the frame image below shows a person is leaning and about to fall. I trimmed the video at this moment. However, these two frames will have every similar body landmark data.



Figure 24: Video 59 from Multiple Camera Dataset

Left frame: marked as fall;

Right frame: marked as not fall

Second, not enough data for confusing actions.

The actor performs confusing actions. For example, laying on the sofa bending to seek objects under the sofa or grabbing an object on the floor. There are some videos contains confusing actions for training, but not enough. Then the network does not have enough samples to learn.

Videos from YOUTUBE dataset improved the performance by adding videos of real video of old people falling in the hospital or people falling outdoor. However, it is still not enough videos. Fall is a dangerous and not fun action, there are not many videos could be found from the internet and it also hurts to trying make fall video by myself.

If there are enough video data of different kinds of fall and confusion actions, my training structure will get a more precise model.

Third, other uncertain reasons.

I did not know the exact reason of other fail predictions. It could be the limitations of the cameras and the body landmark extraction algorithm.

The UR dataset only has two camera positions (parallel to the floor and ceiling-mounted, respectively). The video from ceiling mounted cannot generate body landmark, so I only use video recorded by camera parallel to the floor. Though the Multiple Camera Dataset contains 8 cameras from different positions, there are only 22 videos contains fall. Videos from YOUTUBE dataset has different angles but there are only 46 videos of this dataset. Therefore, there is not enough video from different positions for training. The body landmark extraction also has its limitations and cannot be 100% accurate.

5 Smart Home Solution

I found that for every 30 frames of each video, there are less than 15 consecutive fail frames. Therefore, for a real application, the problem of fall prediction can be solved by for every 30 frames(for video, 30 frames would be 1 second), the prediction of fall or not fall could be decided by the mode of predictions of these 30 frames.

6 Annotated Code

Github repository:

https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model

Model Training:

(The code of training CNN model and LSTM model)

[model_training_Apr25.ipynb](#)

[model_training_Apr25.py](#)

Model Executing code:

Final model(ensembled model): [execute_model_ensembled.py](#)

LSTM model: [execute_model_lstm.py](#)

CNN model: [execute_model_cnn.py](#)

Other code:

Model ensembling and test accuracy of ensembled model:

[test_acc_emsembled_model.py](#)

[test_acc_emsembled_model.ipynb](#)

Plot the architecture of models:

[model_plot.py](#)

7 Instruction of Using Models

1. Clone this repository:

https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model

2. Turn to the directory:

```
```cd ../YJZFlora/Fall_Detection_Deep_Learning_Model```
```

### 3. Install necessary libraries:

- \* python 3
- \* Keras
- \* Tensorflow
- \* pandas
- \* Numpy

**4. Generate body landmark json files, copy the directory of it.(or use the json files provided in sample directory**

I used the open source tool by CMU-Perceptual-Computing-Lab(<https://github.com/CMU-Perceptual-Computing-Lab/openpose>) to generate body landmark.

Sample body landmark files have been included in samples directory.

eg, Fall\_Detection\_Deep\_Learning\_Model/samples/bodylandmark/ test\_case1

## 5. Run one of the following scripts:

Execute new LSTM model:

```
```python3 execute_model_lstm.py <directory of body landmarks for a video>```
```

Execute CNN model:

```
```python3 execute_model_cnn.py <directory of body landmarks for a video>```
```

Execute emsembled model:

```
``` python3 execute_model_ensembled.py <directory of body landmarks for a video>```
```

eg:

```
```python3 execute_model_ensembled.py ./samples/bodylandmark/test_case1```
```

## 6. Result files

Ensembled model:

./results/timeLabel.json

./results/timeLabel.png

LSTM model:

./results/timeLabel\_lstm.json

./results/timeLabel\_lstm.png

CNN model:

./results/timeLabel\_cnn.json

./results/timeLabel\_cnn.png

And in each pair in the file, such as [2.721365,0.23445825932504438], the first number is the time, and the second number is the value of the label your binary-classification model predicts. So the above example shows that at 2.721365 second in the video, the label predicted by your binary-classification model changes to 0.23445825932504438.)

## 7. Video demo

<https://youtu.be/r2CNC9QNPMg>

## 8. Test samples

1)

Video: <https://youtu.be/U8xwla8tQzA>

Body landmark position:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/bodylandmark/test\\_case1](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/bodylandmark/test_case1)

Results:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/test\\_case1](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/test_case1)

2)

Video: [https://youtu.be/5KCNIL\\_NP18](https://youtu.be/5KCNIL_NP18)

Body landmark position:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/bodylandmark/test\\_case2](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/bodylandmark/test_case2)

Results:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/test\\_case2](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/test_case2)

3)

Video: [https://youtu.be/553VQJozd\\_c](https://youtu.be/553VQJozd_c)

Body landmark position:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/bodylandmark/test\\_case3](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/bodylandmark/test_case3)

Results:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/test\\_case3](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/test_case3)

4)

Video: <https://youtu.be/55pul8OFEMs>

Body landmark position:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/bodylandmark/test\\_case4](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/bodylandmark/test_case4)

Results:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/test\\_case4](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/test_case4)

5)

Video: <https://youtu.be/3WK1YGBYQDc>

Body landmark position:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/bodylandmark/test\\_case5](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/bodylandmark/test_case5)

Results:

[https://github.com/YJZFlora/Fall\\_Detection\\_Deep\\_Learning\\_Model/tree/master/samples/test\\_case5](https://github.com/YJZFlora/Fall_Detection_Deep_Learning_Model/tree/master/samples/test_case5)

## Reference

- [1] Bogdan Kwalek, Michal Kepski, Human fall detection on embedded platform using depth maps and wireless accelerometer, *Computer Methods and Programs in Biomedicine*, Volume 117, Issue 3, December 2014, Pages 489-501, ISSN 0169-2607.
- [2] E. Auvinet, C. Rougier, J.Meunier, A. St-Arnaud, J. Rousseau, Multiple cameras fall dataset, Technical report 1350, DIRO - Université de Montréal, July 2010.
- [3] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, Y. A. Sheikh, OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [4] Tomas Simon, Hanbyul Joo, Iain Matthews, Yaser Sheikh, Hand Keypoint Detection in Single Images using Multiview Bootstrapping, 2017.
- [5] Zhe Cao, Tomas Simon, Shih-En Wei, Yaser Sheikh, Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields, 2017
- [6] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, Yaser Sheikh, Convolutional pose machines, 2016.