

Nhận diện cử chỉ tay thời gian thực sử dụng MediaPipe và Random Forest

Nguyễn Thanh Hòa
SE183091

Nguyễn Quý Toàn
SE182785

Lê Minh Hùng
SE182706

July 2024

1 Giới Thiệu

Nhận diện cử chỉ tay đã thu hút được sự chú ý đáng kể do tiềm năng ứng dụng trong tương tác người với máy tính, giao tiếp bằng ngôn ngữ ký hiệu và thực tế ảo. Báo cáo này trình bày về hệ thống nhận diện cử chỉ tay theo thời gian thực, sử dụng MediaPipe để theo dõi bàn tay và bộ phân loại Random Forest cho việc phân loại cử chỉ, không dựa vào các kỹ thuật học sâu.

2 Phạm vi

2.1 Phạm vi của chủ đề

Phạm vi của báo cáo này là khám phá và triển khai hệ thống nhận dạng cử chỉ tay theo thời gian thực sử dụng MediaPipe để theo dõi bàn tay và Random Forest để phân loại cử chỉ.

2.2 Định nghĩa vấn đề

Nhận dạng cử chỉ tay là một lĩnh vực nghiên cứu quan trọng với nhiều ứng dụng thực tế khác nhau, bao gồm tương tác giữa người và máy tính, giao tiếp ngôn ngữ ký hiệu và thực tế ảo. Vấn đề chính được đề cập trong báo cáo này là phát triển một hệ thống đáng tin cậy và hiệu quả để nhận dạng cử chỉ tay trong thời gian thực.

3 Kiến Trúc Hệ Thống

Hệ thống của chúng tôi xây dựng theo kiến trúc hai giai đoạn: tích hợp khả năng theo dõi tay mạnh mẽ của MediaPipe với sức mạnh phân loại của bộ phân loại Random Forest.

3.1 Giai đoạn theo dõi bàn tay: Sử dụng MediaPipe Hands để theo dõi chuyển động của bàn tay. Framework của MediaPipe Hands

MediaPipe Hands là một công cụ hiệu quả và chính xác để theo dõi bàn tay, cung cấp 21 điểm mốc 2.5D đại diện cho các khớp xương của bàn tay. (sau khi MediaPipe Hands tìm thấy bàn tay sẽ tự tạo một khung xương 2.5D giúp mô tả chính xác cấu trúc và vị trí của bàn tay trong không gian). Sử dụng hai mô hình:

1. Bộ phát hiện lòng bàn tay: xác định vị trí ban đầu của tay và vẽ ra một hình hộp bao quanh tay và xoay theo hướng của bàn tay.
2. Mô hình dự đoán điểm khớp: 21 điểm khớp chi tiết của tay trong vùng lòng bàn tay đã phát hiện.

MediaPipe Hands cho phép theo dõi tay hiệu quả và chính xác và được ứng dụng rất nhiều như điều khiển máy tính bằng cử chỉ, chơi game và điều khiển robot bằng cử chỉ.

3.2 Random Forest Classifier

Random Forest là một phương pháp học tổng hợp xây dựng nhiều cây quyết định trong quá trình huấn luyện. Bộ phân loại sau đó dự đoán bằng cách xem xét số phiếu đa số được bầu từ tất cả các cây quyết định. Random Forest có độ bền, độ chính xác cao, và khả năng xử lý dữ liệu có chiều cao.

4 Thu thập và Xử lý Dữ liệu

4.1 Thu thập Dữ liệu

Để huấn luyện bộ phân loại Random Forest, nhóm đã thu thập một bộ dữ liệu các cử chỉ tay :

1. Capture Hand Gestures: Ghi lại các khung hình video của các cá nhân thực hiện các cử chỉ tay khác nhau.
2. Labeling: Gán nhãn thủ công cho từng khung hình với loại cử chỉ tương ứng.
3. Extract Landmarks: Sử dụng MediaPipe Hands để trích xuất các điểm khớp 2.5D của tay từ mỗi khung hình. Giải thích thuật ngữ: (Khung hình là từng ảnh tĩnh riêng được tạo ra từ video)

4.1.1 Code: 'create_dataset.py'

```
1 import os
2 import pickle
3
4 import mediapipe as mp
5 import cv2
6 import matplotlib.pyplot as plt
7
8
9 mp_hands = mp.solutions.hands
10 mp_drawing = mp.solutions.drawing_utils
11 mp_drawing_styles = mp.solutions.drawing_styles
12
13 hands = mp_hands.Hands(static_image_mode=True, min_detection_confidence=0.3)
14
15 DATA_DIR = './data'
16
17 data = []
18 labels = []
19 for dir_ in os.listdir(DATA_DIR):
20     for img_path in os.listdir(os.path.join(DATA_DIR, dir_)):
21         data_aux = []
22
23         x_ = []
24         y_ = []
25
26         img = cv2.imread(os.path.join(DATA_DIR, dir_, img_path))
27         img_rgb = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
28
29         results = hands.process(img_rgb)
30         if results.multi_hand_landmarks:
31             for hand_landmarks in results.multi_hand_landmarks:
32                 for i in range(len(hand_landmarks.landmark)):
33                     x = hand_landmarks.landmark[i].x
34                     y = hand_landmarks.landmark[i].y
35
36                     x_.append(x)
37                     y_.append(y)
38
39                 for i in range(len(hand_landmarks.landmark)):
40                     x = hand_landmarks.landmark[i].x
41                     y = hand_landmarks.landmark[i].y
42                     data_aux.append(x - min(x_))
43                     data_aux.append(y - min(y_))
44
45             data.append(data_aux)
46             labels.append(dir_)
47
48 f = open('data.pickle', 'wb')
49 pickle.dump({'data': data, 'labels': labels}, f)
50 f.close()
```

4.2 Xử lý Dữ liệu

Các điểm đặc trưng đã được trích xuất, sau đó được xử lý trước để chuẩn bị cho việc huấn luyện: (các điểm đặc trưng trên bàn tay, ví dụ như đầu ngón tay, cổ tay, khớp ngón tay,...)

1. Normalization: Chuẩn hóa các tọa độ điểm đặc trưng để loại bỏ sự khác biệt về kích thước và vị trí của tay.
2. Feature Vector: Chuyển đổi các điểm đặc trưng đã chuẩn hóa thành một vector đặc trưng, thường là một mảng 1D. Chúng tôi sử dụng các đặc trưng này để phân biệt giữa các cử chỉ khác nhau.

5 Huấn luyện phân loại (Classifier Training)

Dữ liệu được xử lý trước khi được sử dụng để huấn luyện bộ phân loại Random Forest:

1. Lựa chọn mô hình: Chọn các tham số thích hợp cho Random Forest, các tham số tối ưu cho Random Forest Classifier như số lượng cây, độ sâu tối đa của cây, ...
2. Huấn luyện: Huấn luyện bộ phân loại về các đặc điểm mốc được gắn nhãn bằng cách sử dụng kỹ thuật xác thực chéo để đánh giá mô hình.
3. Đánh giá: Đánh giá hiệu suất của mô hình bằng cách sử dụng các số liệu như độ chính xác (Tỷ lệ dự đoán chính xác trên tổng số dự đoán), Precision (Tỷ lệ dự đoán dương tính chính xác trong số tất cả các dự đoán)...

5.0.1 Code: 'collect_data.py'

```
1 import cv2
2 import os
3
4 def extract_frames(video_paths, output_dir, num_frames=100):
5     for idx, video_path in enumerate(video_paths, start=1):
6         # Tạo thư mục con trong thư mục data
7         video_folder = os.path.join(output_dir, f"{idx}")
8         if not os.path.exists(video_folder):
9             os.makedirs(video_folder)
10
11         # Đối tượng capture video từ file
12         cap = cv2.VideoCapture(video_path)
13
14         # Kiểm tra xem video có mở thành công hay không
15         if not cap.isOpened():
16             print(f"Không thể mở video: {video_path}")
17             continue
18
19         frame_count = 0
20         while frame_count < num_frames:
21             # Đọc từng khung hình
22             ret, frame = cap.read()
23
24             # Kiểm tra nếu không còn khung hình nào để đọc
25             if not ret:
26                 break
27
28             # Lưu khung hình thành file ảnh
29             frame_filename = os.path.join(video_folder, f"frame_{frame_count:04d}.png")
30             cv2.imwrite(frame_filename, frame)
31
32             frame_count += 1
33
34         # Giải phóng tài nguyên
35         cap.release()
36
37     # Đường dẫn đến các video đầu vào và thư mục đầu ra
38     video_paths = [
39         r'C:\Users\Junn\Desktop\FPT_study\CPV\Hand_gesture\bao.mp4',
40         r'C:\Users\Junn\Desktop\FPT_study\CPV\Hand_gesture\keo.mp4',
41         r'C:\Users\Junn\Desktop\FPT_study\CPV\Hand_gesture\bua.mp4'
42     ]
43     output_dir = 'data' # Thư mục để lưu các thư mục con của từng video
44
45     # Số lượng frame muốn lấy cho mỗi video
46     num_frames_per_video = 100
47
48     # Gọi hàm để tách video và lưu các frame vào các thư mục con
49     extract_frames(video_paths, output_dir, num_frames=num_frames_per_video)
```

6 Suy luận và nhận dạng cử chỉ

Sau khi huấn luyện bộ phân loại Rừng ngẫu nhiên, nó có thể được sử dụng để nhận dạng cử chỉ theo thời gian thực:

1. Theo dõi bàn tay (Hand Tracking): MediaPipe Hands liên tục theo dõi bàn tay của người dùng trong thời gian thực.
2. Trích xuất mốc (Landmark Extraction): MediaPipe trích xuất các mốc từ bàn tay được theo dõi.
3. Trích xuất đặc điểm (Feature Extraction): Các mốc được xử lý trước để tạo ra một vectơ đặc trưng.
4. Dự đoán cử chỉ (Gesture Prediction): Trình phân loại Random Forest dự đoán cử chỉ có khả năng xảy ra nhất dựa trên vectơ đặc trưng.

6.0.1 Code: ‘inference_classifier.py’

```
1  import pickle
2
3  import cv2
4  import mediapipe as mp
5  import numpy as np
6
7  model_dict = pickle.load(open('./model.p', 'rb'))
8  model = model_dict['model']
9
10 cap = cv2.VideoCapture(0)
11
12 mp_hands = mp.solutions.hands
13 mp_drawing = mp.solutions.drawing_utils
14 mp_drawing_styles = mp.solutions.drawing_styles
15
16 hands = mp_hands.Hands(static_image_mode=True, min_detection_confidence=0.3)
17
18 labels_dict = {1: 'bao', 2: 'keo', 3: 'bua'}
19
20 while True:
21     ret, frame = cap.read()
22
23     # Check if the frame was successfully read
24     if not ret:
25         print("Error: Unable to read frame from camera. Check camera connection or permissions.")
26         break # Exit the loop if there's an error
27
28     H, W, _ = frame.shape # Now frame.shape will work correctly
29
30     frame_rgb = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
31
32     results = hands.process(frame_rgb)
33     if results.multi_hand_landmarks:
34         for hand_landmarks in results.multi_hand_landmarks:
35             mp_drawing.draw_landmarks(
36                 frame, # image to draw
37                 hand_landmarks, # model output
38                 mp_hands.HAND_CONNECTIONS, # hand connections
39                 mp_drawing_styles.get_default_hand_landmarks_style(),
40                 mp_drawing_styles.get_default_hand_connections_style())
41
42     # Initialize x_ and y_ for each hand
```

```

42 # Initialize x_ and y_ for each hand
43 x_ = []
44 y_ = []
45
46 data_aux = []
47
48 for i in range(len(hand_landmarks.landmark)):
49     x = hand_landmarks.landmark[i].x
50     y = hand_landmarks.landmark[i].y
51
52     x_.append(x)
53     y_.append(y)
54
55 for i in range(len(hand_landmarks.landmark)):
56     x = hand_landmarks.landmark[i].x
57     y = hand_landmarks.landmark[i].y
58     data_aux.append(x - min(x_))
59     data_aux.append(y - min(y_))
60
61 x1 = int(min(x_) * W) - 10
62 y1 = int(min(y_) * H) - 10
63
64 x2 = int(max(x_) * W) - 10
65 y2 = int(max(y_) * H) - 10
66
67 prediction = model.predict([np.asarray(data_aux)])
68
69 predicted_character = labels_dict[int(prediction[0])]
70
71 cv2.rectangle(frame, (x1, y1), (x2, y2), (0, 0, 0), 4)
72 cv2.putText(frame, predicted_character, (x1, y1 - 10), cv2.FONT_HERSHEY_SIMPLEX, 1.3, (0, 0, 0), 3,
73             cv2.LINE_AA)
74
75 cv2.imshow('frame', frame)
76 cv2.waitKey(1)
77
78
79 cap.release()
80 cv2.destroyAllWindows()

```

7 Kết quả và thảo luận

- Các ưu điểm:
- Độ chính xác: Hệ thống đạt tỷ lệ chính xác cao trong việc nhận dạng cử chỉ tay, thể hiện tính hiệu quả của việc kết hợp MediaPipe với Random Forest.
- Hiệu suất thời gian thực: Hệ thống có thể xử lý các khung hình video trong thời gian thực, cho phép tương tác liền mạch với người dùng.
- Khả năng mở rộng: Khung có thể dễ dàng thích ứng với các bộ dữ liệu cử chỉ tay khác nhau và có thể được mở rộng để nhận dạng nhiều cử chỉ hơn.
- Nhược điểm:
- Hiệu suất của hệ thống có thể bị ảnh hưởng bởi các yếu tố như điều kiện ánh sáng, cách che tay và sự thay đổi kích thước bàn tay.

8 Kết luận và kế hoạch phát triển trong tương lai

Báo cáo này trình bày một hệ thống nhận dạng cử chỉ tay theo thời gian thực, tận dụng hiệu quả MediaPipe Hands và bộ phân loại Random Forest. Hệ thống này thể hiện độ chính xác cao và hiệu suất theo thời gian thực, khiến nó phù hợp với nhiều ứng dụng khác nhau.

Kế hoạch trong tương lai sẽ tập trung vào:

1. Mở rộng bộ dữ liệu cử chỉ: Thu thập bộ dữ liệu lớn hơn và đa dạng hơn về cử chỉ tay để cải thiện độ bền và tính tổng quát của mô hình.
2. Cải tiến về độ bền: Triển khai các kỹ thuật nhằm giải quyết tác động của các yếu tố môi trường như ánh sáng và chướng ngại vật.
3. Nhận dạng cử chỉ động: Kết hợp khả năng nhận dạng cử chỉ động để cho phép hệ thống nhận dạng các chuyển động tay liên tục.

9 Source Code

GitHub:https://github.com/Jikay-070203/ASS_CPV301.git

10 Nguồn

1. Fan Zhang và những người khác. MediaPipe Hands: On-device Real-time Hand Tracking. <https://arxiv.org/pdf/2006.10214v1.pdf>
2. Kumar, R. et al. Mediapipe and CNNs for Real-Time ASL Gesture Recognition. <https://arxiv.org/pdf/2305.05296.pdf>