

دانشگاه صنعتی خواجه نصیرالدین طوسی
دانشگاه صنعتی خواجه نصیرالدین طوسی
دانشکده مهندسی برق

درس مبانی مکاترونیک
استاد: دکتر مهدی دلربایی
گزارش کار پروژه پایانی

نام و نام خانوادگی	محمد امین محمدیون شبستری
شماره دانشجویی	۴۰۱۲۲۵۰۳
نام و نام خانوادگی	یکتا خلیلی
شماره دانشجویی	۴۰۱۱۷۹۰۳
نام و نام خانوادگی	سوگل سلامت
شماره دانشجویی	۴۰۱۱۹۳۷۳
تاریخ	شهریور ۱۴۰۴

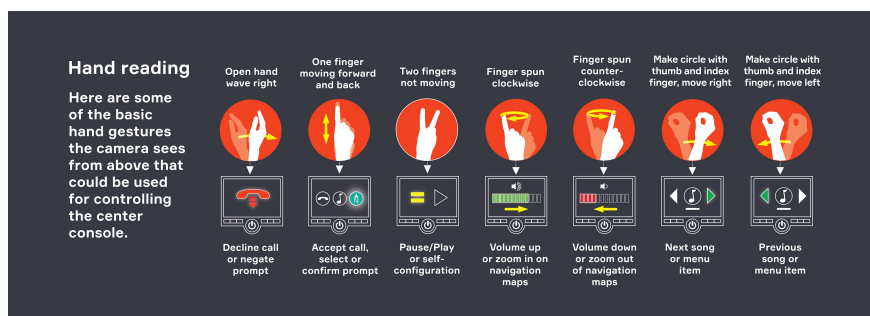
۱ مقدمه

۱.۱ کلیات و اهمیت موضوع

با گسترش روزافزون فناوری‌های نوین در حوزه هوش مصنوعی و بینایی ماشین، ارتباط و تعامل میان انسان و ماشین وارد مرحله‌ای تازه شده است. در گذشته، تعامل با سامانه‌های مکانیکی و الکترونیکی عمدتاً از طریق واسطه‌های فیزیکی مانند کلیدها، دکمه‌ها و صفحه‌کلیدها صورت می‌گرفت. با پیشرفت فناوری، روش‌های جدیدتری همچون صفحات لمسی، فرمان‌های صوتی و اخیراً تشخیص حرکات بدن و دست معرفی شده‌اند که می‌توانند تجربه کاربری را به طور چشمگیری بهبود بخشند.

۲.۱ کاربردهای تشخیص حرکات دست

یکی از مهم‌ترین شاخه‌های این حوزه، تشخیص حرکات دست (Hand Gesture Recognition) است که امکان صدور فرمان به دستگاه‌ها تنها با استفاده از حرکات طبیعی دست را فراهم می‌کند. این فناوری کاربردهای متنوعی در علوم و صنایع مختلف دارد؛ از کنترل ربات‌های صنعتی در محیط‌های حساس گرفته تا استفاده در سیستم‌های واقعیت مجازی (VR) و واقعیت افزوده (AR). در محیط‌های پزشکی، امکان کنترل تجهیزات بدون لمس می‌تواند از آلودگی متقاطع جلوگیری کند و در حوزه سرگرمی، به بازی‌های ویدئویی و سیستم‌های تعاملی سطح بالاتری از واقع‌گرایی می‌بخشد.



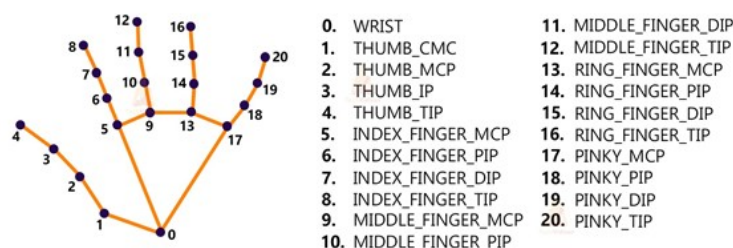
شکل ۱: نمونه‌ای از نگاشت حرکات دست به فرمان‌های کنترلی در یک سیستم تعاملی

۳.۱ نقش مکترونیک در پروژه

در مهندسی مکترونیک که ترکیبی از مکانیک، الکترونیک، کنترل و علوم رایانه است، چنین فناوری‌هایی به عنوان بخشی از سیستم‌های هوشمند نقش مهمی ایفا می‌کنند. این پروژه یک نمونه کاربردی از تلفیق بینایی ماشین و یادگیری ماشین است و نشان‌دهنده رویکرد میان‌رشته‌ای مکترونیک در ترکیب فناوری‌های نرم‌افزاری و سخت‌افزاری برای ایجاد سیستم‌های هوشمند و کارآمد می‌باشد. مزیت مهم آن این است که پیاده‌سازی به صورت نرم‌افزاری انجام شده و در عین حال امکان یکپارچه‌سازی با سخت‌افزارهای مکترونیکی نظیر ربات‌ها را نیز دارد.

۴.۱ معرفی ابزار کلیدی: MediaPipe

یکی از ابزارهای پیشرفته در این زمینه، کتابخانه MediaPipe است که توسط گوگل توسعه یافته و به صورت متن باز در اختیار پژوهشگران قرار گرفته است. این کتابخانه مجموعه‌ای از الگوریتم‌های از پیش ساخته را برای پردازش بلادرنگ تصویر و صوت ارائه می‌دهد. در بخش شناسایی دست، MediaPipe قادر است با دقت بالا ۲۱ نقطه کلیدی (Landmarks) هر دست را شناسایی کند. این نقاط شامل محل دقیق مفاصل، کف دست و نوک انگشتان بوده و ساختار هندسی کامل دست را بازسازی می‌کنند (تصویر ۲). مزیت اصلی این روش، عدم نیاز به تجهیزات گران‌قیمت و استفاده از یک دوربین معمولی برای استخراج داده است.



شکل ۲: ۲۱ نقطه کلیدی دست (Landmarks) که توسط MediaPipe شناسایی می‌شوند.

۵.۱ اهداف پروژه

هدف این پروژه، طراحی و پیاده‌سازی یک سامانه تشخیص حرکات دست با استفاده از MediaPipe و مدل‌های یادگیری ماشین است. فرآیند شامل مراحل زیر است:

- تهیه و بارگذاری یک مجموعه داده شامل تصاویر حرکات دست.
- استخراج مختصات نقاط کلیدی دست از هر تصویر با MediaPipe.
- آموزش یک مدل طبقه‌بندی‌کننده برای شناسایی الگوهای حرکتی.
- ارزیابی عملکرد مدل برای سنجش دقت و کارایی سیستم.
- پیاده‌سازی یک برنامه نهایی برای کنترل بی‌درنگ مدیاپایپر با استفاده از مدل آموزش دیده.

۲ فرآیند آموزش مدل

۱.۲ بارگذاری و آماده‌سازی داده‌ها

۱.۱.۲ فراخوانی کتابخانه‌ها و دانلود مجموعه داده

در گام نخست، کتابخانه‌های اصلی مورد نیاز پروژه شامل numpy، pandas و os فراخوانی شدند.

```
1 import numpy as np
2 import pandas as pd
3 import os
```

Code 1: اولیه کتابخانه‌های فراخوانی

سپس، با استفاده از کتابخانه‌های kagglehub و mediapipe، مجموعه داده hand-gesture-detection-system از Kaggle دانلود و در محیط اجرایی بارگذاری شد.

```
1 import kagglehub
2 import mediapipe as mp
3 path = kagglehub.dataset_download("marusagar/hand-gesture-detection-system")
4 print("Path to dataset files:", path)
```

Code 2: Kaggle از داده مجموعه دانلود

۲.۱.۲ تحلیل و اصلاح ساختار اولیه داده‌ها

پس از بارگذاری، مشخص شد که فایل train.csv ساختاری نامناسب دارد و تمام اطلاعات هر سطر در یک ستون با جداکننده «؛» قرار گرفته‌اند (تصویر ۳).

WIN_20180925_17_08_43_Pro_Left_Swipe_new;Left_Swipe_new;0	
0	WIN_20180925_17_18_28_Pro_Left_Swipe_new;Left...
1	WIN_20180925_17_18_56_Pro_Left_Swipe_new;Left...
2	WIN_20180925_17_19_51_Pro_Left_Swipe_new;Left...
3	WIN_20180925_17_20_14_Pro_Left_Swipe_new;Left...
4	WIN_20180925_17_21_28_Pro_Left_Swipe_new;Left...
...	...
657	WIN_20180907_16_38_23_Pro_Thumbs Up_new;Thumbs...
658	WIN_20180907_16_41_09_Pro_Thumbs Up_new;Thumbs...
659	WIN_20180907_16_42_05_Pro_Thumbs Up_new;Thumbs...
660	WIN_20180907_16_42_55_Pro_Thumbs Up_new;Thumbs...
661	WIN_20180907_16_43_39_Pro_Thumbs Up_new;Thumbs...
662 rows × 1 columns	

شکل ۳: ساختار اولیه و نامناسب فایل train.csv پیش از پردازش.

این مشکل با خواندن مجدد فایل و تعیین 'sep=';' و همچنین تعریف نام ستون‌ها (CLASS, LABEL, IMAGE_NAME) برطرف گردید.

```
1 df = pd.read_csv('.../train.csv', sep=';', names=['IMAGE_NAME', 'LABEL', 'CLASS'])
```

Code 3: صحیح ستون‌های نام و جداکننده با CSV فایل خواندن

	IMAGE_NAME	LABEL	CLASS
0	WIN_20180925_17_08_43_Pro_Left_Swipe_new	Left_Swipe_new	0
1	WIN_20180925_17_18_28_Pro_Left_Swipe_new	Left_Swipe_new	0
2	WIN_20180925_17_18_56_Pro_Left_Swipe_new	Left_Swipe_new	0
3	WIN_20180925_17_19_51_Pro_Left_Swipe_new	Left_Swipe_new	0
4	WIN_20180925_17_20_14_Pro_Left_Swipe_new	Left_Swipe_new	0
...
658	WIN_20180907_16_38_23_Pro_Thumbs Up_new	Thumbs Up_new	4
659	WIN_20180907_16_41_09_Pro_Thumbs Up_new	Thumbs Up_new	4
660	WIN_20180907_16_42_05_Pro_Thumbs Up_new	Thumbs Up_new	4
661	WIN_20180907_16_42_55_Pro_Thumbs Up_new	Thumbs Up_new	4
662	WIN_20180907_16_43_39_Pro_Thumbs Up_new	Thumbs Up_new	4
663 rows × 3 columns			

شکل ۴: ساختار اصلاح شده دیتافریم پس از بارگذاری صحیح.

در این دیتافریم، هر سطر به یک ویدیو کلیپ کوتاه از یک ژست خاص اشاره دارد. ستون LABEL شامل نام ژست مانند Left_Swipe_new، Thumbs_Up_new، Stop_Gesture_new و غیره است.

۲.۲ پاک‌سازی و پیش‌پردازش داده‌ها

۱.۲.۲ استانداردسازی برچسب‌های متنی

داده‌های خام معمولاً نیازمند پاک‌سازی هستند تا برای مدل‌های یادگیری ماشین قابل استفاده باشند. در اینجا، برچسب‌های موجود در ستون LABEL حاوی اطلاعات اضافی و فرمت‌های ناهمگون بودند. عباراتی مانند new_ و کاراکتر _ حذف شدند تا نام ژست‌ها خوانا تر و استاندارد شوند.

```

1 df['LABEL'] = df['LABEL'].str.replace('_', " ")
2 df['LABEL'] = df['LABEL'].str.replace('new', "")
3 df['LABEL'].value_counts()

```

برچسب‌ها از اضافی کاراکترهای و عبارات حذف: Code 4:



count	LABEL
137	Right Swipe
137	Thumbs Down
123	Thumbs Up
96	Left Swipe
93	Stop
40	Left Swipe Left Swipe
37	Stop Gesture

dtype: int64

شکل ۵: توزیع اولیه برچسب‌ها که شامل موارد تکراری است.

۲.۲.۲ یکپارچه‌سازی کلاس‌های مشابه

همانطور که در تصویر ۵ مشخص است، برخی ژست‌ها با نام‌های مختلفی ثبت شده بودند. برای مثال، "Left Swipe Left Swipe" و "Stop Gesture" وجود داشتند. این موارد به ترتیب به "Left Swipe" و "Stop" تبدیل شدند تا هر کلاس معنای منحصر به فردی داشته باشد. این کار از ایجاد کلاس‌های اضافی و سردرگمی مدل جلوگیری می‌کند.

```
1 df['LABEL'] = df['LABEL'].str.replace('Left Swipe Left Swipe ', "Left Swipe ")
2 df['LABEL'] = df['LABEL'].str.replace('Stop Gesture ', "Stop ")
```

مشابه برچسب‌های یکپارچه‌سازی: Code 5

count	
LABEL	
Right Swipe	137
Thumbs Down	137
Left Swipe	136
Stop	130
Thumbs Up	123
dtype: int64	

شکل ۶: توزیع نهایی و پاک‌سازی شده برچسب‌ها در مجموعه داده آموزشی.

این فرایند برای مجموعه داده اعتبارسنجی (val.csv) نیز تکرار شد.

count	
LABEL	
Right Swipe	23
Stop	22
Thumbs Down	21
Left Swipe	18
Thumbs Up	16
dtype: int64	

شکل ۷: توزیع پاک‌سازی شده برچسب‌ها در مجموعه داده اعتبارسنجی.

۳.۲.۲ مقایسه توزیع کلاس‌ها در داده‌های آموزش و اعتبارسنجی

یک مجموعه داده اعتبارسنجی خوب باید توزیع کلاس‌های مشابهی با مجموعه داده آموزشی داشته باشد. با مقایسه تصاویر ۶ و ۷ می‌توان دید که:



- در هر دو مجموعه، ۵ کلاس نهایی (Right Swipe, Thumbs Down, Left Swipe, Stop, Thumbs Up) وجود دارند.
- توزیع داده‌ها در هر دو مجموعه نسبتاً متوازن (Balanced) است، یعنی تعداد نمونه‌های هر کلاس تفاوت فاحشی با دیگری ندارد. این موضوع به مدل کمک می‌کند تا به تمام کلاس‌ها اهمیت یکسانی بدهد و از سوگیری (Bias) به سمت کلاس‌های پرتعدادتر جلوگیری شود.
- نسبت تعداد نمونه‌ها در داده اعتبارسنجی نیز منعکس‌کننده توزیع در داده آموزشی است که نشان می‌دهد این تقسیم‌بندی برای ارزیابی عملکرد مدل مناسب است.

۳.۲ استخراج ویژگی با MediaPipe

۱.۳.۲ نقاط کلیدی (Landmarks) چه هستند؟

به جای استفاده مستقیم از پیکسل‌های تصویر که ابعاد بسیار بزرگی دارند و شامل اطلاعات اضافی (مانند پس‌زمینه) هستند، از یک رویکرد مبتنی بر ویژگی استفاده می‌شود. Landmarks یا نقاط کلیدی، مجموعه‌ای از نقاط با مختصات مشخص روی یک شیء هستند که ساختار هندسی آن را توصیف می‌کنند. برای دست انسان، MediaPipe ۲۱ نقطه را تعریف کرده که متناظر با مفاصل انگشتان، نوک انگشتان و مچ است (همانطور که در تصویر ۲ نشان داده شد). هر نقطه دارای سه مختصات (x, y, z) است. این مختصات نرمال‌سازی شده‌اند:

- x و y : نسبت به عرض و ارتفاع تصویر نرمال شده و مقداری بین ۰ و ۱ دارند.
 - z : عمق نقطه را نسبت به مچ دست نشان می‌دهد. مقدار کوچکتر به معنای نزدیک‌تر بودن به دوربین است.
- با اتصال این ۲۱ نقطه، یک اسکلت سه‌بعدی از دست به دست می‌آید که اطلاعات بسیار غنی و فشرده‌ای از ژست دست را بدون توجه به پس‌زمینه یا رنگ پوست ارائه می‌دهد.

۲.۳.۲ پیکربندی ماژول MediaPipe Hands

ماژول `mp.solutions.hands.Hands` یک پایپ‌لاین پیچیده از مدل‌های یادگیری عمیق است. فرآیند کار آن به طور خلاصه به این صورت است: ابتدا یک مدل سبک به نام `Palm Detector` کل تصویر را برای پیدا کردن ناحیه کف دست جستجو می‌کند. پس از یافتن دست، ناحیه تصویر برش خورده و به یک مدل سنگین‌تر به نام `Hand Landmark Model` داده می‌شود که ۲۱ نقطه کلیدی را با دقت بالا روی آن مشخص می‌کند. این رویکرد دو مرحله‌ای باعث افزایش سرعت و دقت می‌شود.

```
1 import mediapipe as mp
2 hands = mp.solutions.hands.Hands(
3     static_image_mode=True,
4     max_num_hands=1,
5     min_detection_confidence=0.5
6 )
```

Code 6: دست تشخیص مدل پیکربندی

۳.۳.۲ تابع استخراج نقاط کلیدی و تبدیل فرمت

خروجی MediaPipe یک شیء پیچیده است. تابع `extract_landmarks` این شیء را به یک فرمت استاندارد و قابل استفاده برای مدل‌های یادگیری ماشین تبدیل می‌کند. این تابع، مختصات x, y, z هر ۲۱ نقطه را پشت سر هم قرار داده و یک بردار یک‌بعدی با طول $21 \times 3 = 63$ ایجاد می‌کند. برای مثال، خروجی برای یک نقطه به صورت $[x0, y0, z0, x1, y1, z1, \dots, x20, y20, z20]$ خواهد بود.

```
1 def extract_landmarks(result):
2     if result.multi_hand_landmarks:
3         hand_landmarks = result.multi_hand_landmarks[0]
4         landmarks_ = []
5         for lm in hand_landmarks.landmark:
6             landmarks_.extend([lm.x, lm.y, lm.z])
7         return np.array(landmarks_)
8     return None
```

دست کلیدی نقاط مسطح‌سازی و استخراج تابع: Code 7:

۴.۳.۲ اجرای فرآیند استخراج ویژگی روی داده‌ها

یک حلقه بر روی تمام تصاویر مجموعه داده اجرا شد. در هر تکرار، تصویر با OpenCV خوانده شده و به فرمت رنگی RGB تبدیل می‌شد. دلیل این تبدیل این است که OpenCV تصاویر را به صورت پیش‌فرض در فرمت BGR (آبی، سبز، قرمز) می‌خواند، در حالی که مدل‌های MediaPipe بر روی تصاویر با فرمت RGB (قرمز، سبز، آبی) آموزش دیده‌اند. عدم تطابق این فرمت‌ها می‌تواند منجر به تشخیص نادرست شود.

	Image_name	Folder	Label	Class	Landmarks
0	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180925_17_08_43_Pro_Left_Swipe_new	Left Swipe	0	[0.29479125142097473, 0.8378140330314636, 3.55...
1	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180925_17_08_43_Pro_Left_Swipe_new	Left Swipe	0	[0.26317331194877625, 0.7459539771080017, 2.65...
2	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180925_17_08_43_Pro_Left_Swipe_new	Left Swipe	0	[0.24926461279392242, 0.6994234323501587, 2.24...
3	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180925_17_08_43_Pro_Left_Swipe_new	Left Swipe	0	[0.21890145540237427, 0.6593414545059204, 2.04...
4	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180925_17_08_43_Pro_Left_Swipe_new	Left Swipe	0	[0.21465922892093658, 0.64450097080404541, 2.08...
...
19885	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180907_16_43_39_Pro_Thumbs Up_new	Thumbs Up	4	None
19886	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180907_16_43_39_Pro_Thumbs Up_new	Thumbs Up	4	None
19887	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180907_16_43_39_Pro_Thumbs Up_new	Thumbs Up	4	None
19888	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180907_16_43_39_Pro_Thumbs Up_new	Thumbs Up	4	None
19889	/root/.cache/kagglehub/datasets/marusagar/hand...	WIN_20180907_16_43_39_Pro_Thumbs Up_new	Thumbs Up	4	None

19890 rows × 5 columns

شکل ۸: دیتافریم نهایی شامل ستون Landmarks که ویژگی‌های استخراج شده را نگهداری می‌کند.

۵.۳.۲ مدیریت داده‌های فاقد ویژگی و دلیل آن

در این فرآیند، MediaPipe در ۵۰۴۶ تصویر آموزشی و ۶۷۳ تصویر اعتبارسنجی موفق به تشخیص دست نشد و مقدار None برگرداند. دلیل اصلی این اتفاق، پارامتر `min_detection_confidence=0.5` است. این بدان معناست که مدل تشخیص کف دست (Palm Detector) در این تصاویر، یک ناحیه را با اطمینان کمتر از ۵۰٪ به عنوان دست شناسایی کرده و در نتیجه آن را نادیده گرفته است. دلایل

دیگر می‌تواند شامل تاری شدید تصویر، نور بسیار کم یا زیاد، یا قرار گرفتن دست در زاویه‌ای نامتعارف باشد. از آنجایی که این سطرها فاقد ویژگی ورودی برای مدل بودند، با دستور dropna از دیتافریم‌ها حذف شدند. این کار برای جلوگیری از بروز خطا و آموزش مدل با داده‌های ناقص ضروری است.

```
1 df_images.dropna(subset=['Landmarks'], inplace=True)
2 df_valid.dropna(subset=['Landmarks'], inplace=True)
```

ویژگی فاقد سطرهای حذف: Code 8:

۳ تعریف، پیاده‌سازی و ارزیابی مدل‌ها

۱.۳ مدل ماشین بردار پشتیبان (SVM)

۱.۱.۳ دلیل انتخاب SVM به عنوان مدل اولیه

مدل SVM به چند دلیل به عنوان اولین انتخاب در نظر گرفته شد:

- کارایی در فضای با ابعاد بالا: داده‌های ما یک بردار ۶۳ بعدی هستند. SVM در چنین فضاهایی عملکرد بسیار خوبی دارد.
- سرعت و کارایی: برای مجموعه داده‌هایی با اندازه متوسط (مانند پروژه ما با حدود ۱۵۰۰۰ نمونه)، SVM معمولاً سریع‌تر از شبکه‌های عصبی عمیق آموزش می‌بیند و نیاز به منابع محاسباتی کمتری دارد.
- مقاومت در برابر بیش‌برازش: با تنظیم صحیح پارامتر C، می‌توان توازن خوبی بین دقت و تعمیم‌پذیری مدل ایجاد کرد.

۲.۱.۳ بهینه‌سازی هایپر پارامترها با GridSearchCV

برای یافتن بهترین ترکیب هایپر پارامترها، از GridSearchCV استفاده شد. پارامترهای بررسی شده در جدول زیر آمده‌اند:

پارامتر	مقادیر آزمایشی	توضیح
kernel	rbf, linear, poly	نوع هسته برای نگاشت داده‌ها به فضای با ابعاد بالاتر.
C	0.1, 1, 10, 100	پارامتر تنظیم (Regularization). مقدار بالاتر، خطای کمتری را در داده آموزشی مجاز می‌داند.
gamma	0.001, 0.01, 0.1, 1, 'scale'	ضریب هسته برای rbf و poly. تأثیر یک نمونه آموزشی را مشخص می‌کند.
degree	2, 3, 4	درجه چندجمله‌ای برای هسته poly.

جدول ۱: پارامترهای مورد استفاده در فرآیند جستجوی شبکه‌ای.

نتایج نشان داد که بهترین عملکرد با کرنل rbf، $C=100$ و $\gamma=0.001$ به دست آمد.

۳.۱.۳ اهمیت نرمال‌سازی داده‌ها

قبل از ورود داده‌ها به مدل SVM، از StandardScaler برای نرمال‌سازی استفاده شد. این کار مقیاس تمام ۶۳ ویژگی را یکسان می‌کند (میانگین صفر و واریانس یک). از آنجایی که SVM بر اساس محاسبه فاصله بین نمونه‌ها کار می‌کند، اگر ویژگی‌ها مقیاس‌های متفاوتی داشته باشند، ویژگی‌هایی با مقادیر بزرگتر تأثیر نامتناسبی بر مدل خواهند گذاشت. نرمال‌سازی این مشکل را حل کرده و به همگرایی سریع‌تر و عملکرد بهتر مدل کمک می‌کند.

۴.۱.۳ ارزیابی دقیق عملکرد مدل SVM

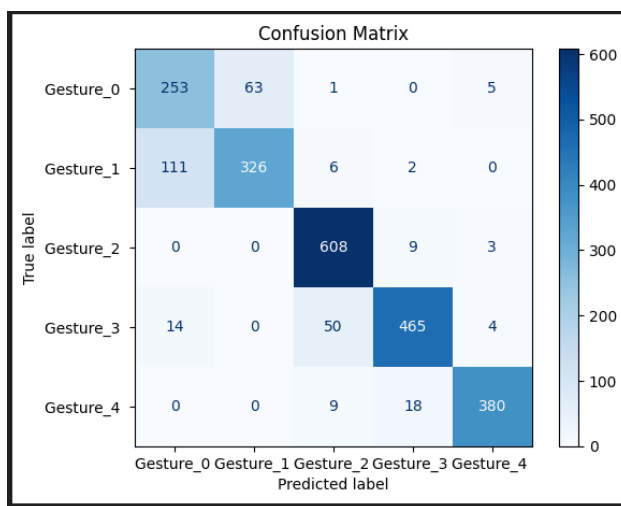
مدل نهایی دقت کلی ۸۷.۳۲٪ را کسب کرد. تحلیل دقیق‌تر معیارها در تصویر ۹ نشان می‌دهد:

- Precision (دقت پیش‌بینی): این معیار نشان می‌دهد که از میان نمونه‌هایی که مدل به عنوان یک کلاس خاص پیش‌بینی کرده، چه درصدی واقعاً درست بوده‌اند. برای مثال، Precision ۰.۹۷ برای کلاس ۴ (Thumbs Up) به این معناست که اگر مدل ژستی را به عنوان Thumbs Up تشخیص دهد، در ۹۷٪ موارد درست گفته است.
- Recall (بازخوانی): این معیار نشان می‌دهد که مدل چه درصدی از نمونه‌های واقعی یک کلاس را توانسته به درستی شناسایی کند. برای مثال، Recall ۰.۹۸ برای کلاس ۲ (Stop) یعنی مدل ۹۸٪ از تمام ژست‌های واقعی را پیدا کرده است.
- F1-Score: این معیار میانگین هماهنگ بین Precision و Recall است و یک دید کلی از توازن بین این دو ارائه می‌دهد.

تحلیل ماتریس درهم‌ریختگی (تصویر ۱۰) نشان می‌دهد که ۱۱۱ نمونه از کلاس ۱ (Right Swipe) به اشتباه کلاس ۰ (Left Swipe) تشخیص داده شده‌اند که بزرگترین منبع خطا است. این خطا به دلیل شباهت ساختاری و قرینه بودن این دو ژست کاملاً قابل انتظار است.

	precision	recall	f1-score	support
0	0.67	0.79	0.72	322
1	0.84	0.73	0.78	445
2	0.90	0.98	0.94	620
3	0.94	0.87	0.91	533
4	0.97	0.93	0.95	407
accuracy			0.87	2327
macro avg	0.86	0.86	0.86	2327
weighted avg	0.88	0.87	0.87	2327
Accuracy: 0.87322733132789				

شکل ۹: نتایج بهترین مدل svm



شکل ۱۰: ماتریس در هم ریختگی مربوط به svm

۲.۳ مدل پرسپترون چندلایه (MLP)

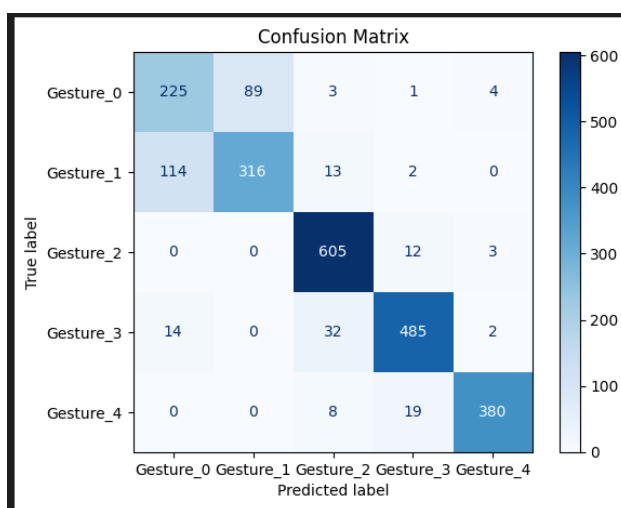
۱.۲.۳ معماری و هایپرپارامترهای مدل

یک مدل MLP به عنوان جایگزین SVM پیاده‌سازی شد. معماری آن شامل یک لایه ورودی با ۶۳ نورون، دو لایه مخفی با ۲۵۶ و ۱۲۸ نورون و تابع فعال‌سازی ReLU بود. برای جلوگیری از بیش‌برازش، پس از هر لایه مخفی یک لایه Dropout با نرخ ۰.۳ قرار داده شد که به طور تصادفی ۳۰٪ از نورون‌ها را در هر مرحله آموزش غیرفعال می‌کرد. هایپرپارامترهای آموزش به شرح زیر بودند:

- بهینه‌ساز: AdamW
- نرخ یادگیری: $1e-4$
- تابع هزینه: CrossEntropyLoss (مناسب برای مسائل طبقه‌بندی چندکلاسه)
- اندازه دسته: ۶۴
- مکانیسم توقف: توقف زودهنگام با $patience=5$

۲.۲.۳ ارزیابی عملکرد مدل MLP

این مدل به دقت کلی ۸۶.۴۲٪ دست یافت. تحلیل ماتریس درهم‌ریختگی آن (تصویر ۱۱) نیز نتایج مشابهی با SVM نشان داد؛ یعنی بیشترین خطا در تفکیک بین ژست‌های ۰ و ۱ بود. هرچند عملکرد کلی خوب بود، اما اندکی ضعیف‌تر از مدل SVM عمل کرد.



شکل ۱۱: ماتریس در هم ریختگی mlp

	precision	recall	f1-score	support
0	0.64	0.70	0.67	322
1	0.78	0.71	0.74	445
2	0.92	0.98	0.94	620
3	0.93	0.91	0.92	533
4	0.98	0.93	0.95	407
accuracy			0.86	2327
macro avg	0.85	0.85	0.85	2327
weighted avg	0.87	0.86	0.86	2327

شکل ۱۲: نتایج ساختار mlp

۳.۳ مقایسه و انتخاب مدل نهایی

جدول زیر دو مدل را بر اساس معیارهای مختلف مقایسه می‌کند:

مدل MLP	مدل SVM	معیار
۸۶.۴۲٪	۸۷.۳۲٪	دقت (Accuracy)
۰.۸۶	۰.۸۷	F1-Score (وزنی)
پیچیده‌تر	سبک و ساده	پیچیدگی مدل
بالا (اما کندتر از SVM)	بسیار بالا	سرعت پیش‌بینی (Inference)
نسبتاً سریع (با توقف زود هنگام)	طولانی (به دلیل GridSearch)	زمان آموزش
PyTorch (سنگین‌تر)	Scikit-learn	وابستگی‌ها

جدول ۲: مقایسه جامع بین مدل‌های SVM و MLP



با توجه به دقت بالاتر و همچنین سبک تر و سریع تر بودن در مرحله پیش بینی، مدل SVM به عنوان گزینه نهایی برای پیاده سازی در برنامه بی درنگ انتخاب شد.

۴ پیاده سازی نهایی و کاربرد بی درنگ

۱.۴ معماری سیستم بی درنگ و ابزارها

این سیستم با هدف پردازش لحظه ای تصویر وبکم، تشخیص ژست و ارسال فرمان به سیستم عامل طراحی شده است. کتابخانه های کلیدی در جدول قبل معرفی شدند.

۲.۴ بارگذاری مدل آموزش دیده

در ابتدای اسکریپت، مدل SVM ذخیره شده با استفاده از joblib.load بارگذاری می شود.

```
1 import joblib
2 try:
3     model = joblib.load('final_model_v2_svm (1).pkl')
4 except FileNotFoundError:
5     print("Error: Model file not found.")
6     exit()
```

Code 9: مدل بارگذاری SVM ذخیره شده

۳.۴ نگاشت ژست به فرمان (Gesture-to-Command Mapping)

دیکشنری GESTURE_MAP کلاس های خروجی مدل را به فرمان های قابل فهم و کلیدهای کیبورد نگاشت می دهد.

```
1 GESTURE_MAP = {
2     0: "Rewind",          # ->
3     1: "Fast Forward",   # ->
4     2: "Play/Pause",     # ->   Space
5     3: "Volume Down",    # ->
6     4: "Volume Up"       # ->
7 }
```

Code 10: کنترل فرمان های به پیش بینی کلاس های نگاشت

۴.۴ حلقه اصلی و پردازش فریم

برنامه وارد یک حلقه بی نهایت می شود که در هر تکرار، یک فریم از وبکم می خواند، آن را برمی گرداند و به فرمت RGB تبدیل می کند.



```
1 while True:
2     ret, frame = cap.read()
3     if not ret:
4         break
5     frame = cv2.flip(frame, 1)
6     img_rgb = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
7     result = hands.process(img_rgb)
8     # ...
```

ویدئویی فریم‌های پردازش حلقه اصلی بخش: Code 11

۵.۴ پیش‌بینی و اجرای فرمان با PyAutoGUI

پس از تشخیص دست، ویژگی‌ها استخراج و ژست توسط مدل SVM پیش‌بینی می‌شود. برای جلوگیری از ارسال فرمان‌های تکراری، یک تأخیر یک ثانیه‌ای (ACTION_DELAY) در نظر گرفته شده است.

```
1 current_time = time.time()
2 if current_time - last_action_time > ACTION_DELAY:
3     if action_name == "Play/Pause":
4         pyautogui.press('space')
5         video_is_playing = not video_is_playing
6         last_action_time = current_time
7     elif video_is_playing:
8         if action_name == "Fast Forward":
9             pyautogui.press('right')
10            last_action_time = current_time
11    # ...
```

تشخیص داده‌شده ژست اساس بر فرمان اجرای منطق: Code 12

۶.۴ بازخورد بصری و خاتمه برنامه

نقاط کلیدی دست، وضعیت پخش و نام ژست تشخیص داده‌شده به صورت زنده روی فریم ویدئو ترسیم می‌شوند. با فشردن کلید 'q'، برنامه خاتمه می‌یابد.

۵ نتیجه‌گیری و پیشنهادات

۱.۵ جمع‌بندی نهایی

این پروژه با موفقیت یک سیستم کامل از آموزش تا پیاده‌سازی برای تشخیص حرکات دست را اجرا کرد. نتایج نشان داد که استخراج ویژگی‌های ساختاری دست با استفاده از MediaPipe و طبقه‌بندی آن‌ها با یک مدل یادگیری ماشین کلاسیک و سبک مانند SVM، یک رویکرد بسیار مؤثر و کارآمد است. مدل نهایی با دستیابی به دقت ۸۷.۳۲٪، توانایی خود را در تفکیک پنج ژست مختلف به اثبات رساند و کاربرد عملی آن در کنترل مدیاپلیر، پتانسیل این فناوری را در تعاملات انسان و ماشین به خوبی نمایش داد.

۲.۵ پیشنهادات برای کارهای آتی

برای بهبود و گسترش این پروژه در آینده، موارد زیر پیشنهاد می‌شود:

- استفاده از مدل‌های زمانی (Temporal Models): برای تشخیص بهتر حرکات پویا که در طول زمان تعریف می‌شوند (مانند Swipe)، می‌توان از شبکه‌های عصبی بازگشتی (RNN) مانند LSTM یا GRU استفاده کرد که توالی فریم‌ها را تحلیل می‌کنند.
- افزایش داده (Data Augmentation): می‌توان با اعمال تغییرات جزئی مانند چرخش، تغییر مقیاس یا افزودن نویز به بردارهای ویژگی Landmarks، داده‌های آموزشی را غنی‌تر کرده و مدل را در برابر شرایط مختلف مقاوم‌تر نمود.
- گسترش مجموعه ژست‌ها: با جمع‌آوری داده‌های بیشتر، می‌توان تعداد ژست‌های قابل تشخیص را افزایش داد و سیستم را برای کاربردهای پیچیده‌تر، مانند کنترل یک بازوی رباتیک، توسعه داد.
- بهینه‌سازی برای محیط‌های نویزی: می‌توان مدل را با داده‌هایی که در شرایط نوری مختلف یا با پس‌زمینه‌های شلوغ گرفته شده‌اند، آموزش داد تا پایداری آن در محیط‌های واقعی افزایش یابد.
- یکپارچه‌سازی با سخت‌افزار و پیاده‌سازی: به عنوان فاز بعدی، پروژه بر روی یک برد Raspberry Pi 3 پیاده‌سازی خواهد شد. در این مرحله، عملگر (actuator) سیستم از کنترل نرم‌افزاری مدیاپلیر به کنترل سخت‌افزارهای فیزیکی تغییر خواهد یافت. اهداف این فاز شامل کنترل یک الگوی نوری با استفاده از آرایه‌ای از LEDها یا کنترل جهت و سرعت یک موتور الکتریکی بر اساس ژست‌های تشخیص داده شده می‌باشد. این گام، کاربرد عملی پروژه در سیستم‌های مکاترونیکی را به نمایش می‌گذارد.