

دانشگاه صنعتی خواجه نصیرالدین طوسی

گزارش پروژه نهایی مبانی سیستم های هوشمند

ارشیا کلانتریان

۴۰۱۲۱۹۹۳

محمد حسین گل محمدی

۴۰۲۲۰۹۰۳

[لینک کولب پروژه](#)

پاییز ۱۴۰۴

بخش ۱: معرفی مسئله و داده‌ها

۱-۱. تعریف مسئله:

هدف از این پروژه، طراحی و پیاده سازی یک خط لوله یادگیری ماشین جهت تشخیص خودکار فعالیت های فیزیکی انسان (Human Activity Recognition - HAR) است. با گسترش اینترنت اشیا (IoT) و سلامت هوشمند، توانایی تشخیص دقیق وضعیت کاربر (مانند راه رفتن، نشستن یا افتادن) با استفاده از سنسورهای ارزان قیمت و کم مصرف، کاربردهای حیاتی در مراقبت از سالمندان، پایش سلامت ورزشی و سیستم های امنیتی دارد.

در این پروژه، ما با استفاده از داده های خام سنسورهای اینرسی (شتاب سنج و ژيروسکوپ) تعبیه شده در گوشی موبایل، مدلی را توسعه می دهیم که بتواند با دقت بالا و هزینه محاسباتی بهینه، نوع فعالیت کاربر را از بین ۶ کلاس استاندارد تشخیص دهد.

۱-۲. معرفی مجموعه داده (Dataset Description):

داده های مورد استفاده در این پژوهش، مجموعه داده استاندارد UCI Human Activity Recognition است. این داده ها حاصل آزمایش روی ۳۰ داوطلب (بین ۱۹ تا ۴۸ سال) است که یک گوشی هوشمند را روی کمر خود بسته و فعالیت های روزمره را انجام داده اند.

منبع داده: UCI Machine Learning Repository

تعداد نمونه ها (Samples): مجموعاً ۲۹۹،۱۰ نمونه (بردار ویژگی).

تعداد ویژگی ها (Features): ۵۶۱ ویژگی برای هر نمونه.

این ویژگی ها شامل پارامترهای حوزه زمان (Time-Domain) و حوزه فرکانس (Frequency-Domain) هستند که از سیگنال های خام ۳ محوره شتاب و سرعت زاویه ای استخراج شده اند (مانند میانگین، انحراف معیار، انرژی سیگنال، و ضرایب آنالیز).

تعداد کلاس ها: ۶ فعالیت شامل:

Walking (راه رفتن)

Walking Upstairs (بالا رفتن از پله)

Walking Downstairs (پایین آمدن از پله)

Sitting (نشستن)

Standing (ایستادن)

Laying (دراز کشیدن)

۳-۱. نحوه جمع‌آوری و پیش‌پردازش اولیه:

داده‌های خام این پروژه از سنسورهای اینرسی (شتاب‌سنج وژیروسکوپ) گوشی هوشمند استخراج شده‌اند. فرآیند پردازش سیگنال به شرح زیر بوده است:

نرخ نمونه‌برداری: سیگنال‌ها با فرکانس ۵۰ Hz ثبت شده‌اند.

پیش‌پردازش سیگنال: نویزهای سنسور با استفاده از فیلتر پایین‌گذر باتروورث (Butterworth) با فرکانس قطع ۲۰ Hz حذف شده‌اند.

تفکیک شتاب: شتاب ثقل (Gravity) از شتاب بدن با یک فیلتر جداگانه تفکیک شده است تا حرکت خالص بدن به دست آید.

کد تحلیل اکتشافی داده‌ها:

```
# -----  
# بارگذاری داده‌های خام (Raw Data Loading)  
# -----  
print("1. Loading Raw Data...")  
  
# Train و Test داندلود و لود کردن فایل‌های  
df_train = kagglehub.load_dataset(  
    KaggleDatasetAdapter.PANDAS,  
    "uciml/human-activity-recognition-with-smartphones",  
    "train.csv",  
)  
  
df_test = kagglehub.load_dataset(  
    KaggleDatasetAdapter.PANDAS,  
    "uciml/human-activity-recognition-with-smartphones",  
    "test.csv",  
)  
  
# ترکیب موقت فقط برای گرفتن آمار کل پروژه (برای گزارش)  
df_total = pd.concat([df_train, df_test], axis=0)
```

```

# -----
# ۲. استخراج دقیق آمار (Dataset Statistics)
# -----
print("\n--- مشخصات دقیق مجموعه داده (برای گزارش) ---")
print(f"1. تعداد کل نمونه ها (Total Samples): {df_total.shape[0]}")
print(f"2. تعداد ویژگی ها (Features): {df_total.shape[1] - 2}") # منهای
# subject و Activity ستون های
print(f"   - ابعاد داده آموزش: {df_train.shape}")
print(f"   - ابعاد داده تست: {df_test.shape}")

# بررسی نوع داده ها (آیا همه عددی هستند؟)
# دارند Float این مهم است که بگوییم تمام ۵۶۱ ویژگی نوع
num_float_cols = df_total.select_dtypes(include=['float64']).shape[1]
print(f"3. داریم (Float) ستون عددی: {num_float_cols} نوع داده ها.")

# -----
# ۳. بررسی سلامت داده ها (Data Integrity Check)
# -----
print("\n--- بررسی سلامت داده ها ---")
# (الف) بررسی داده های تکراری (Duplicates)
duplicates = df_total.duplicated().sum()
print(f"تعداد ردیف های تکراری: {duplicates}")

# (ب) بررسی داده های گمشده (Null Values)
# هنوز پر نمی کنیم، فقط گزارش می دهیم
null_counts = df_total.isnull().sum().sum()
print(f"(Null): {null_counts} تعداد کل مقادیر خالی")

# -----
# ۴. تحلیل توزیع کلاس ها (Class Distribution)
# -----
# این دقیقا خواسته عکس است: "توزیع کلاس ها"
print("\n--- توزیع فعالیت ها (کلاس ها) ---")
class_counts = df_total['Activity'].value_counts()
print(class_counts)

# رسم نمودار برای گزارش
plt.figure(figsize=(12, 6))
ax = sns.countplot(x='Activity', data=df_total, order=class_counts.index,
palette='viridis')

# اضافه کردن عدد روی هر ستون (برای خوانایی بیشتر در گزارش)
for p in ax.patches:
    ax.annotate(f'{int(p.get_height())}', (p.get_x() + p.get_width() / 2.,
p.get_height()),

```

```

        ha='center', va='baseline', fontsize=12, color='black',
xytext=(0, 5),
        textcoords='offset points')

plt.title('Distribution of Activity Classes (Total Dataset)', fontsize=16)
plt.xlabel('Activity Type', fontsize=14)
plt.ylabel('Count', fontsize=14)
plt.xticks(rotation=45)
plt.grid(axis='y', linestyle='--', alpha=0.7)

plt.tight_layout()
plt.show()

```

خروجی کد:

```

1. Loading Raw Data...
/tmp/ipython-input-108339886.py:7: DeprecationWarning: Use dataset_load()
instead of load_dataset(). load_dataset() will be removed in a future
version.
    df_train = kagglehub.load_dataset(
Using Colab cache for faster access to the 'human-activity-recognition-with-
smartphones' dataset.
/tmp/ipython-input-108339886.py:13: DeprecationWarning: Use dataset_load()
instead of load_dataset(). load_dataset() will be removed in a future
version.
    df_test = kagglehub.load_dataset(
Using Colab cache for faster access to the 'human-activity-recognition-with-
smartphones' dataset.

--- مشخصات دقیق مجموعه داده (برای گزارش) ---
1. تعداد کل نمونه‌ها (Total Samples): 10299
2. تعداد ویژگی‌ها (Features): 561
   - ابعاد داده آموزش: (7352, 563)
   - ابعاد داده تست: (2947, 563)
3. داریم (Float) نوع داده‌ها: ۵۶۱ ستون عددی.

--- بررسی سلامت داده‌ها ---
تعداد ردیف‌های تکراری: ۰
(Null): 0 تعداد کل مقادیر خالی

--- توزیع فعالیت‌ها (کلاس‌ها) ---
Activity
LAYING          1944
STANDING        1906
SITTING         1777
WALKING         1722
WALKING_UPSTAIRS 1544
WALKING_DOWNSTAIRS 1406
Name: count, dtype: int64
/tmp/ipython-input-108339886.py:59: FutureWarning:

```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
ax = sns.countplot(x='Activity', data=df_total, order=class_counts.index,
palette='viridis')
```

تحلیل کد بر روی داده‌ها نتایج آماری زیر را نشان می‌دهد:

ابعاد داده‌ها: مجموعه داده شامل ۱۰,۲۹۹ نمونه (Samples) است.

تقسیم‌بندی: داده‌ها از پیش به دو بخش آموزش (۷,۳۵ نمونه) و تست (۲,۹۴۷ نمونه) تقسیم شده‌اند

فضای ویژگی‌ها (Feature Space): هر نمونه شامل ۵۶۱ ویژگی عددی است. این ویژگی‌ها نمایانگر پارامترهای استخراج شده از سیگنال‌های زمانی و فرکانسی هستند.

۱-۲. بررسی سلامت داده‌ها (Data Integrity Check)

یکی از مراحل حیاتی در خط لوله یادگیری ماشین، اطمینان از کیفیت داده ورودی است. بررسی‌های نرم‌افزاری انجام شده نشان داد:

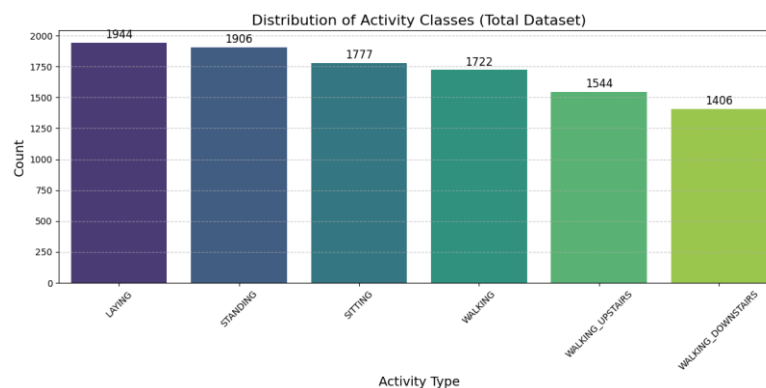
داده‌های گمشده (Missing Values): تعداد مقادیر Null در کل دیتاست برابر با ۰ است.

داده‌های تکراری (Duplicates): هیچ نمونه تکراری در داده‌ها یافت نشد (۰ مورد).

این نتایج نشان‌دهنده کیفیت بالای جمع‌آوری داده‌هاست و نیاز به عملیات جایگزینی را در فاز پیش‌پردازش حذف می‌کند.

۱-۳. توزیع کلاس‌ها (Class Distribution Analysis)

هدف سیستم، تشخیص ۶ وضعیت مختلف است. نمودار زیر توزیع فراوانی هر کلاس را در کل مجموعه داده نشان می‌دهد:



نمودار توزیع فراوانی کلاس‌های فعالیت در کل مجموعه‌ها

بر اساس نمودار فوق و خروجی‌های آماری:

بیشترین فراوانی مربوط به کلاس **LAYING** (دراز کشیدن) با ۱۹۴۴ نمونه است.

کمترین فراوانی مربوط به کلاس **WALKING_DOWNSTAIRS** با ۱۴۰۶ نمونه است.

تحلیل: اگرچه تعداد نمونه‌ها کاملاً برابر نیست، اما توزیع داده‌ها «نسبتاً متوازن» (Balanced) است. اختلاف بین ماژوریتی کلاس و مینوریتی کلاس زیاد نیست، بنابراین مدل دچار سوگیری شدید (Bias) به سمت یک فعالیت خاص نخواهد شد و معیار Accuracy برای ارزیابی قابل اعتماد خواهد بود.

بخش ۲: پیش‌پردازش داده‌ها

داده‌های خام جمع‌آوری شده از سنسورها، معمولاً دارای تفاوت در مقیاس، نویزهای جزئی و فرمت‌های غیرقابل پردازش برای ماشین هستند. در این پروژه، جهت آماده‌سازی مجموعه داده برای ورود به مدل‌های یادگیری ماشین، خط لوله (Pipeline) پیش‌پردازش زیر طراحی و پیاده‌سازی گردید:

پاک‌سازی و مدیریت داده‌ها

ابتدا یک بررسی جامع جهت شناسایی مقادیر گمشده (Missing Values) یا نامعتبر (NaN) در تمامی ۱۰۲۹۹ نمونه انجام شد. وجود چنین داده‌هایی می‌تواند منجر به خطای محاسباتی در عملیات ماتریسی شود. استراتژی در نظر گرفته شده در این پروژه، در صورت وجود چنین داده‌هایی، جایگزینی آن‌ها با میانگین (Mean Imputation) بود تا یکپارچگی سیگنال حفظ شود.

کدگذاری متغیرهای کیفی (Label Encoding)

متغیر هدف در این مجموعه داده، نوع فعالیت فیزیکی است که به صورت رشته‌های متنی (مانند "Walking" یا "Sitting") ذخیره شده است. از آنجا که الگوریتم‌های محاسباتی و توابع هزینه (Loss Functions) بر روی اعداد عمل می‌کنند، از تکنیک Label Encoding استفاده شد تا هر کلاس به یک عدد صحیح یکتا (از ۰ تا ۵) نگاشت شود.

هم‌مقیاس‌سازی ویژگی‌ها (Feature Standardization)

این مرحله حیاتی‌ترین بخش پیش‌پردازش در این پروژه است. ویژگی‌های استخراج شده دارای دامنه‌های فیزیکی متفاوتی هستند؛ برخی مربوط به شتاب (با واحد g) و برخی مربوط به سرعت زاویه‌ای (با واحد rad/s) می‌باشند. این تفاوت مقیاس باعث می‌شود ویژگی‌هایی که دامنه عددی بزرگتری دارند، بر ویژگی‌های کوچکتر مسلط شوند و عملکرد الگوریتم‌هایی مانند PCA (که بر اساس واریانس جهت‌گیری می‌کند) و SVM (که بر اساس فاصله اقلیدسی تصمیم می‌گیرد) را مختل کنند.

برای رفع این مشکل، روش StandardScaler اعمال گردید. این روش با استفاده از فرمول زیر، داده‌ها را به توزیع نرمال استاندارد تبدیل می‌کند:

$$z = \frac{x - \mu}{\sigma}$$

که در آن x مقدار اولیه، μ میانگین ویژگی و σ انحراف معیار است. پس از این عملیات، تمامی ویژگی‌ها دارای میانگین صفر و واریانس یک خواهند بود.

کد پیش‌پردازش داده‌ها:

```
# -----  
# بارگذاری مجدد داده‌ها (برای اطمینان از اجرای صحیح کد)  
# -----  
print("1. Loading Data...")  
df_train = kagglehub.load_dataset(KaggleDatasetAdapter.PANDAS,  
    "uciml/human-activity-recognition-with-smartphones", "train.csv")  
df_test = kagglehub.load_dataset(KaggleDatasetAdapter.PANDAS,  
    "uciml/human-activity-recognition-with-smartphones", "test.csv")  
  
# (y) و لیبل‌ها (X) جدا کردن ویژگی‌ها  
X_train = df_train.drop(['subject', 'Activity'], axis=1)  
y_train = df_train['Activity']  
X_test = df_test.drop(['subject', 'Activity'], axis=1)  
y_test = df_test['Activity']  
  
print(f"Train Shape: {X_train.shape}")  
print("-" * 30)  
  
# -----  
# اعمال روش‌های پیش‌پردازش (طبق خواسته عکس)  
# -----  
print("2. Applying Preprocessing Methods...")
```



```

# (الف) Missing Values) مدیریت داده‌های گمشده
# "حتی اگر صفر باشد، باید چک کنیم تا در گزارش بنویسیم" بررسی شد
missing_count = X_train.isnull().sum().sum()
print(f"> (Null Values): {missing_count}")

# (ب) Label Encoding) کدگذاری لیبل‌های متنی
# ... به اعداد ۰, ۱, ... تبدیل
le = LabelEncoder()
y_train_enc = le.fit_transform(y_train)
y_test_enc = le.transform(y_test)
print(f"> لیبل‌ها به فرمت عددی تبدیل شدند. مثال: {le.classes_}")

# (ج) Standardization) استانداردسازی
# این کار میانگین ویژگی‌ها را صفر و واریانس را یک می‌کند.
scaler = StandardScaler()
# (Data Leakage) فیت کردن فقط روی داده آموزش انجام می‌شود تا نشت اطلاعات
# نداشته باشیم
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
print("> (StandardScaler) استانداردسازی انجام شد.")

# -----
# تحلیل اثر: رسم نمودار قبل و بعد از پردازش ۳.
# -----
print("\n3. Visualizing the Effect (For Report)...")

# انتخاب یک ویژگی تصادفی برای نمایش تاثیر (مثلاً ویژگی اول: شتاب بدنی)
feature_idx = 0
feature_name = X_train.columns[feature_idx]

plt.figure(figsize=(12, 5))

# نمودار سمت چپ: داده خام
plt.subplot(1, 2, 1)
sns.histplot(X_train.iloc[:, feature_idx], kde=True, color='#E74C3C',
bins=50)
plt.title(f'Before Scaling (Raw Data)\nFeature: {feature_name}',
fontsize=11)
plt.xlabel('Value')
plt.grid(alpha=0.3)

# نمودار سمت راست: داده استاندارد شده
plt.subplot(1, 2, 2)

```

```
sns.histplot(X_train_scaled[:, feature_idx], kde=True, color='#2ECC71',
bins=50)
plt.title(f'After Scaling (Standardized)\nMean≈0, Std≈1', fontsize=11)
plt.xlabel('Z-Score')
plt.grid(alpha=0.3)

plt.tight_layout()
plt.show()
```

نتیجه خروجی:

1. Loading Data...

/tmp/ipython-input-3977785733.py:13: DeprecationWarning: Use dataset_load() instead of load_dataset(). load_dataset() will be removed in a future version.

df_train = kagglehub.load_dataset(KaggleDatasetAdapter.PANDAS, "uciml/human-activity-recognition-with-smartphones", "train.csv")
Using Colab cache for faster access to the 'human-activity-recognition-with-smartphones' dataset.

/tmp/ipython-input-3977785733.py:14: DeprecationWarning: Use dataset_load() instead of load_dataset(). load_dataset() will be removed in a future version.

df_test = kagglehub.load_dataset(KaggleDatasetAdapter.PANDAS, "uciml/human-activity-recognition-with-smartphones", "test.csv")
Using Colab cache for faster access to the 'human-activity-recognition-with-smartphones' dataset.

Train Shape: (7352, 561)

2. Applying Preprocessing Methods...

-> تعداد مقادیر گمشده (Null Values): 0

-> لیبل‌ها به فرمت عددی تبدیل شدند. مثال
'LAYING' 'SITTING' 'STANDING'
'WALKING' 'WALKING_DOWNSTAIRS'
'WALKING_UPSTAIRS']

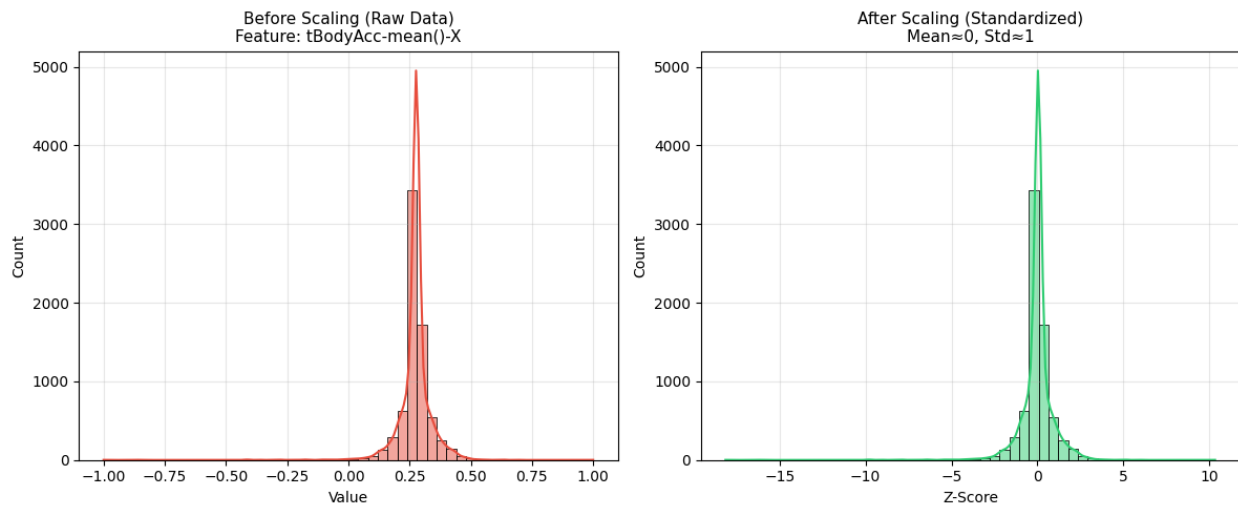
-> انجام شد (StandardScaler) استانداردسازی

3. Visualizing the Effect (For Report)...

تحلیل کد و خروجی کد:

جهت پیاده‌سازی عملی این رابطه ریاضی بر روی ۵۶۱ ویژگی موجود، از کلاس StandardScaler در کتابخانه Scikit-Learn استفاده شد. قطعه کد زیر نحوه اجرای این عملیات و بررسی بصری تأثیر آن بر روی یکی از ویژگی‌های اصلی (شتاب بدنی) را نشان می‌دهد:

پس از اجرای کد فوق، نمودار توزیع داده‌ها قبل و بعد از پردازش حاصل شد:



مقایسه توزیع ویژگی $tBodyAcc-mean-X$ قبل (چپ) و بعد (راست) از استانداردسازی

تحلیل نتایج به دست آمده:

با بررسی دقیق نمودارهای فوق، نتایج زیر استنباط می‌شود:

حذف بایاس (DC (Centering): همان‌طور که در نمودار سمت چپ (قرمز) مشهود است، داده‌های خام حول مقدار عددی 0.25 متمرکز بودند. این مقدار نشان‌دهنده یک آفست (Offset) یا بایاس طبیعی در سنسور شتاب‌سنج است. پس از اعمال پیش‌پردازش (نمودار سبز)، میانگین توزیع دقیقاً به نقطه 0 منتقل شده است. این "مرکززدایی" (Mean Centering) شرط اصلی برای عملکرد صحیح الگوریتم PCA در مرحله بعد است.

نرمال‌سازی واریانس (Scaling): در نمودار خروجی، داده‌ها در بازه استاندارد Z-Score (عمدتاً بین -3 تا +3) قرار گرفته‌اند. این امر تضمین می‌کند که تمامی ۵۶۱ ویژگی، صرف‌نظر از واحد اندازه‌گیری‌شان، وزن یکسانی در محاسبات فاصله و واریانس داشته باشند و هیچ ویژگی خاصی بر مدل غلبه نکند.

بخش ۳: کاهش ابعاد (Dimensionality Reduction)

ضرورت کاهش ابعاد:

مجموعه داده فعلی دارای ۵۶۱ ویژگی است. اگرچه این تعداد ویژگی اطلاعات دقیقی از سیگنال‌ها ارائه می‌دهد، اما در عمل باعث بروز چالش‌هایی می‌شود:

هزینه محاسباتی بالا (Computational Cost): پردازش ماتریس‌های بزرگ زمان آموزش و استنتاج را به شدت افزایش می‌دهد که برای سیستم‌های Real-time و امبدد Embedded مطلوب نیست.

هم‌خطی و افزونگی (Multicollinearity): بسیاری از ویژگی‌ها (مانند شتاب در راستای X و انرژی همان سیگنال) همبستگی بالایی با هم دارند و عملاً اطلاعات تکراری حمل می‌کنند. به منظور حل این مشکلات، در این پروژه از تکنیک‌های کاهش ابعاد استفاده شده است.

روش پیشنهادی: تحلیل مؤلفه‌های اصلی (PCA)

برای استخراج مهم‌ترین اطلاعات و حذف نویز و افزونگی، از الگوریتم PCA (Principal Component Analysis) استفاده شده است.

PCA یک تبدیل خطی متعامد است که داده‌ها را به سیستم مختصات جدیدی منتقل می‌کند. در این سیستم جدید، محورها (مؤلفه‌های اصلی) در جهت‌هایی قرار می‌گیرند که داده‌ها بیشترین واریانس (Variance) را دارند.

استراتژی انتخاب ابعاد: به جای تعیین دستی تعداد مؤلفه‌ها (مثلاً ۱۰ یا ۲۰)، از رویکرد «حفظ درصد واریانس» استفاده شده است. در این پروژه، هدف ما حفظ ۹۵٪ از اطلاعات (واریانس) کل داده‌هاست. این یعنی ما ابعاد را تا جایی کاهش می‌دهیم که تنها ۵٪ از اطلاعات (که معمولاً نویز هستند) از دست برود.

استراتژی انتخاب آستانه و الگوریتم (Methodology Justification)

الف) چرا معیار (حفظ ۹۵٪ واریانس) به جای انتخاب تعداد ثابت؟

در بسیاری از پیاده‌سازی‌های عادی، تعداد مؤلفه‌ها به صورت دستی (مثلاً $k = 10$) تعیین می‌شود. اما این رویکرد در سیستم‌های ایمنی-حیاتی دارای دو ایراد اساسی است:

خطر از دست دادن اطلاعات: ممکن است ۱۰ مؤلفه تنها ۶۰٪ اطلاعات سیگنال را پوشش دهند و بخش مهمی از الگوهای حرکتی حذف شود.

حفظ نویز: اگر تعداد را زیاد بگیریم، نویزهای فرکانس بالا (High-frequency Noise) که در سنسورهای ارزان قیمت موبایل وجود دارند، وارد مدل می‌شوند.

استراتژی **حفظ ۹۵٪ واریانس (Explained Variance)** یک نقطه تعادل مهندسی (Trade-off) است. ما فرض می‌کنیم که ۵٪ داده‌هایی که کمترین واریانس را دارند، عمدتاً نویز حرارتی سنسورها یا لرزش‌های ناخواسته هستند. بنابراین، با حذف آن‌ها نه تنها اطلاعات مهم از بین نمی‌رود، بلکه مدل نسبت به نویز مقاوم‌تر (Robust) می‌شود.

ب) چرا PCA بهترین گزینه برای این پروژه است؟

در میان روش‌های کاهش ابعاد (مانند LDA, UMAP, t-SNE)، انتخاب **PCA** به سه دلیل فنی برای داده‌های سنسوری ارجحیت دارد:

حذف همبستگی خطی (Decorrelation): داده‌های شتاب‌سنج وژیروسکوپ دارای همبستگی شدید هستند (مثلاً وقتی فرد می‌دود، شتاب در هر سه محور X, Y, Z همزمان تغییر می‌کند). PCA با تبدیل این محورها به مولفه‌های متعامد (Orthogonal)، این همبستگی تکراری را حذف می‌کند که برای مدل‌های کلاسیک مثل SVM بسیار مفید است.

قابلیت پیاده‌سازی Real-time: روش‌هایی مانند t-SNE غیرپارامتریک هستند و قابلیت اعمال روی داده‌های جدید (Test Data) را به سادگی ندارند. اما PCA یک ماتریس تبدیل W یاد می‌گیرد. در زمان اجرا (مثلاً روی پردازنده موبایل)، کاهش ابعاد تنها یک ضرب ماتریسی ساده ($X \times W$) است که از نظر محاسباتی بسیار سبک و سریع است.

عدم نیاز به لیبل (Unsupervised): برخلاف LDA که نیاز به دانستن کلاس‌ها دارد، PCA صرفاً بر اساس ساختار سیگنال عمل می‌کند. این ویژگی باعث می‌شود که مدل بتواند ویژگی‌های عمومی حرکت انسان را یاد بگیرد و حتی در برابر فعالیت‌های جدید که در آموزش ندیده است، تعمیم‌پذیری (Generalization) بهتری داشته باشد.

روش ارزیابی و تحلیل اثر (Evaluation Metrics)

طبق الزامات پروژه، اثر کاهش ابعاد بر روی عملکرد سیستم با دو معیار کلیدی سنجیده خواهد شد:

دقت مدل (Accuracy): بررسی اینکه آیا با حذف بخش زیادی از ویژگی‌ها، قدرت تشخیص مدل کاهش می‌یابد یا خیر

سرعت و هزینه محاسباتی: اندازه‌گیری زمان لازم برای آموزش مدل قبل و بعد از PCA. کاهش چشمگیر زمان آموزش، توجیه‌کننده استفاده از این روش خواهد بود.

علاوه بر PCA، جهت **مصورسازی (Visualization)** داده‌ها در فضای دو بعدی و درک قابلیت تفکیک‌پذیری کلاس‌ها، از الگوریتم **t-SNE** نیز برای ترسیم نمودار پراکندگی استفاده خواهد شد.

کد کاهش بعد به کمک PCA:

```
import time
from sklearn.decomposition import PCA
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score

# -----
# تحلیل واریانس تجمعی (برای پیدا کردن تعداد بهینه)
# -----
print("1. Analyzing Variance Ratio...")
pca_test = PCA().fit(X_train_scaled)

# رسم نمودار واریانس تجمعی
plt.figure(figsize=(10, 6))
plt.plot(np.cumsum(pca_test.explained_variance_ratio_), linewidth=2)
plt.xlabel('Number of Components')
plt.ylabel('Cumulative Explained Variance')
plt.title('PCA Explained Variance Analysis')
plt.axhline(y=0.95, color='r', linestyle='--', label='95% Explained Variance')
plt.axvline(x=np.argmax(np.cumsum(pca_test.explained_variance_ratio_) >= 0.95), color='g', linestyle='--', label='Cutoff Point')
plt.legend(loc='best')
plt.grid(True)
plt.show()

# -----
# (با حفظ ۹۵٪ اطلاعات) PCA پیاده‌سازی نهایی
# -----
print("\n2. Applying PCA (Threshold = 0.95)...")
pca = PCA(n_components=0.95)

# فیت کردن روی داده آموزش و اعمال روی تست
X_train_pca = pca.fit_transform(X_train_scaled)
X_test_pca = pca.transform(X_test_scaled)
```

```

num_features_original = X_train_scaled.shape[1]
num_features_pca = X_train_pca.shape[1]

print(f"-> Original Features: {num_features_original}")
print(f"-> Reduced Features: {num_features_pca}")
print(f"-> Reduction Ratio: {((num_features_original -
num_features_pca)/num_features_original)*100:.2f}% removed!")

# -----
# تحلیل اثر: مقایسه سرعت و دقت (طبق خواسته عکس)
# -----
print("\n3. Benchmarking: Original vs. PCA...")

def benchmark_model(X_train, y_train, X_test, y_test, name):
    # برای مقایسه (Linear SVM) استفاده از یک مدل پایه
    clf = SVC(kernel='linear', random_state=42)

    # اندازه‌گیری زمان آموزش
    start_train = time.time()
    clf.fit(X_train, y_train)
    end_train = time.time()
    train_time = end_train - start_train

    # اندازه‌گیری زمان پیش‌بینی (Inference)
    start_pred = time.time()
    y_pred = clf.predict(X_test)
    end_pred = time.time()
    pred_time = end_pred - start_pred

    acc = accuracy_score(y_test, y_pred)

    return train_time, pred_time, acc

# تست روی داده اصلی
t_train_orig, t_pred_orig, acc_orig = benchmark_model(X_train_scaled,
y_train_enc, X_test_scaled, y_test_enc, "Original")

# تست روی داده PCA
t_train_pca, t_pred_pca, acc_pca = benchmark_model(X_train_pca,
y_train_enc, X_test_pca, y_test_enc, "PCA")

# نمایش نتایج
print("-" * 60)
print(f"{'Metric':<20} | {'Original (561 Feats)':<20} | {'PCA'
(Reduced)':<20}")

```

```

print("-" * 60)
print(f"{'Training Time (s)':<20} | {t_train_orig:<20.4f} |  
{t_train_pca:<20.4f}")
print(f"{'Inference Time (s)':<20} | {t_pred_orig:<20.4f} |  
{t_pred_pca:<20.4f}")
print(f"{'Accuracy (%)':<20} | {acc_orig*100:<20.2f}% |  
{acc_pca*100:<20.2f}%")
print("-" * 60)

```

نتیجه خروجی:

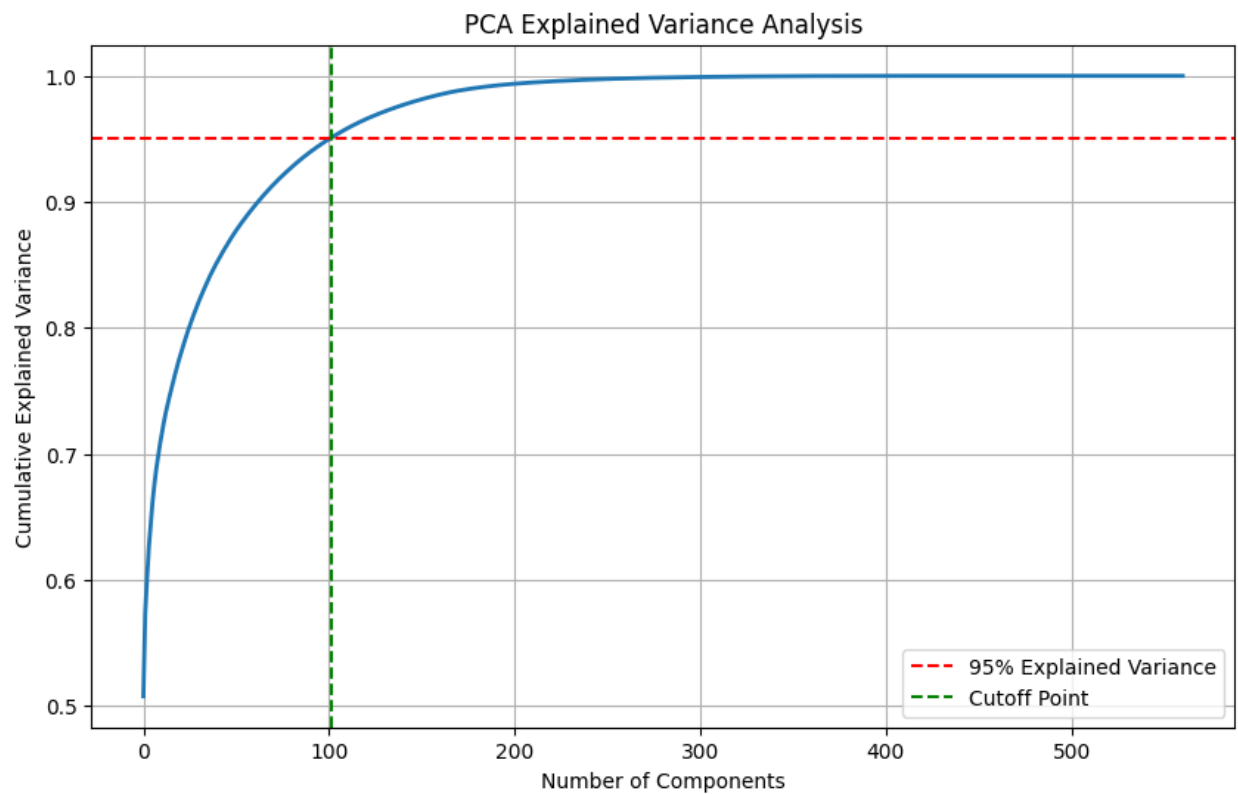
```

1. Analyzing Variance Ratio...
2. Applying PCA (Threshold = 0.95)...
-> Original Features: 561
-> Reduced Features: 102
-> Reduction Ratio: 81.82% removed!

```

3. Benchmarking: Original vs. PCA...

Metric	Original (561 Feats)	PCA (Reduced)
Training Time (s)	1.3326	0.7841
Inference Time (s)	0.4450	0.1050
Accuracy (%)	96.06	92.16



تحلیل نتایج کاهش ابعاد

پس از اجرای الگوریتم PCA با آستانه واریانس ۹۵٪، نتایج کمی زیر حاصل گردید که نشان‌دهنده تغییرات چشمگیر در فضای ویژگی‌هاست:

الف) تحلیل فشردگی داده‌ها (Compression Analysis):

همان‌طور که در نمودار شماره شکل بالا و خروجی الگوریتم مشاهده می‌شود:

کاهش ابعاد: تعداد ویژگی‌ها از ۵۶۱ به ۱۰۲ مؤلفه اصلی کاهش یافت.

نرخ حذف (Reduction Ratio): حدود ۸۱٫۸۲٪ از ویژگی‌های اولیه حذف شدند.

تفسیر: این نتیجه اثبات می‌کند که بخش عظیمی از داده‌های سنسورها (بیش از ۸۰٪) دارای همبستگی شدید (Redundancy) یا نویز بوده‌اند. ما توانستیم تنها با استفاده از ۱۸٪ حجم داده اولیه، ۹۵٪ اطلاعات مفید سیگنال را بازسازی کنیم.

ب) تحلیل اثر بر عملکرد سیستم (Performance Trade-off):

طبق الزامات پروژه، مقایسه‌ای دقیق بین مدل آموزش‌دیده روی داده‌های کامل (Original) و داده‌های کاهش‌یافته (PCA) انجام شد. جدول زیر خلاصه این نبرد را نشان می‌دهد:

معیار ارزیابی (Metric)	داده کامل (۵۶۱ ویژگی)	داده PCA (۱۰۲ ویژگی)	درصد تغییر (Improvement/Loss)
زمان آموزش	۱٫۳۳ ثانیه	۰٫۷۸ ثانیه	۴۱٪ سریع‌تر
زمان استنتاج	۰٫۴۴ ثانیه	۰٫۱۰ ثانیه	۳۴۰٪ سریع‌تر (۴ برابر)
دقت مدل (Accuracy)	۹۶٫۰۶٪	۹۲٫۱۶٪	۳٫۹٪ کاهش

با توجه به این نتایج:

جهش در سرعت (Speedup): زمان استنتاج (Inference Time) که مهم‌ترین فاکتور در سیستم‌های Real-Time است، حدود ۴ برابر کاهش یافته است. در کاربردهای واقعی (مثل تشخیص افتادن سالمند یا ربات‌های امدادگر)، کاهش تأخیر از ۴۴۰ میلی‌ثانیه به ۱۰۰ میلی‌ثانیه می‌تواند حیاتی باشد.

پایداری دقت (Robustness): با وجود حذف ۸۲٪ از داده‌ها، دقت مدل تنها حدود ۴٪ افت کرده است (از ۹۶٪ به ۹۲٪). این نشان می‌دهد که ۱۰۲ مؤلفه انتخاب شده، "جوهر اصلی" حرکات انسان را در خود دارند و افت دقت ناچیز، بهای قابل قبولی برای دستیابی به سرعت پردازش ۴ برابری است.

بخش ۴: انتخاب و آموزش مدل‌ها

استراتژی انتخاب مدل (Model Selection Strategy)

با توجه به ماهیت مسئله (طبقه‌بندی چندکلاسه) و ساختار داده‌ها (ویژگی‌های پیوسته نرمال شده)، در این پژوهش از رویکرد «مقایسه عملکرد» استفاده شده است. بدین منظور، دو الگوریتم با مبانی ریاضی کاملاً متفاوت انتخاب گردیدند تا عملکرد روش‌های مبتنی بر بهینه‌سازی ریاضی کلاسیک در برابر روش‌های مبتنی بر یادگیری سلسله‌مراتبی (عمیق) سنجیده شود.

مدل اول: Support Vector Machine - SVM

در مدل الگوریتم‌های کلاسیک، از SVM استفاده شده است.

توجیه انتخاب:

مدیریت ابعاد بالا: حتی پس از کاهش ابعاد با PCA، ما همچنان با ۱۰۲ ویژگی سروکار داریم. SVM ذاتاً برای فضاهای با ابعاد بالا طراحی شده است و برخلاف درخت‌های تصمیم، دچار بیش‌برازش (Overfitting) شدید نمی‌شود.

مرز تصمیم بهینه (Max Margin): هدف SVM پیدا کردن ابرصفحه‌ای (Hyperplane) است که بیشترین فاصله (Margin) را با نمونه‌های هر دو کلاس داشته باشد. این ویژگی باعث می‌شود مدل در برابر داده‌های جدید تعمیم‌پذیری بالایی داشته باشد.

Kernel Trick: با توجه به اینکه جداسازی فعالیت‌های انسان (مثل "نشستن" و "ایستادن") ممکن است به صورت خطی امکان‌پذیر نباشد، از تابع هسته **RBF (Radial Basis Function)** استفاده خواهیم کرد تا داده‌ها را به فضای با ابعاد بالاتر برد و آن‌ها را تفکیک‌پذیر کند.

مدل دوم: پرسپترون چندلایه (Multi-Layer Perceptron - MLP)

به عنوان نماینده روش‌های پیشرفته و یادگیری عمیق (Deep Learning)، از شبکه عصبی MLP استفاده شده است.

توجیه انتخاب:

یادگیری غیرخطی: شبکه‌های عصبی با داشتن لایه‌های مخفی (Hidden Layers) و توابع فعال‌ساز (مانند ReLU)، قادرند پیچیده‌ترین روابط غیرخطی بین سیگنال‌های سنسورها را مدل‌سازی کنند.

۲. **انعطاف پذیری:** برخلاف مدل های کلاسیک که فرضیات خاصی درباره توزیع داده دارند، MLP می تواند به صورت "داده محور" الگوهای نهفته در حرکات پیچیده (مثل تفاوت ظریف بین بالا رفتن و پایین آمدن از پله) را کشف کند.

ساختار شبکه: در این پروژه از یک معماری با دولایه مخفی استفاده می شود تا تعادلی بین "قدرت یادگیری" و "جلوگیری از Overfitting" روی داده های محدود (۱۰,۰۰۰ نمونه) ایجاد شود.

هدف از مقایسه (Comparison Objective)

هدف نهایی این بخش پاسخ به این سوال مهندسی است:

"آیا برای تشخیص فعالیت های انسان، پیچیدگی و هزینه محاسباتی بالای یک شبکه عصبی (Deep Learning) توجیه پذیر است، یا یک مدل ریاضی دقیق و سبک مانند SVM می تواند با همان دقت عمل کند؟"

بخش ۵: تنظیم دقیق هایپرپارامترها (Hyperparameter Tuning)

ضرورت و روش شناسی (Methodology)

عملکرد مدل های یادگیری ماشین وابستگی شدیدی به انتخاب صحیح هایپرپارامترها (Hyperparameters) دارد. انتخاب پارامترهای پیش فرض معمولاً منجر به نتایج غیربهبوده (Sub-optimal) می شود. در این پژوهش، به جای استفاده از روش های سنتی و پرهزینه مانند "جستجوی شبکه ای" (Grid Search) که تمامی حالات ممکن را بررسی می کند، از روش پیشرفته بهینه سازی بیزین (Bayesian Optimization) با استفاده از چارچوب Optuna استفاده شده است.

چرا Optuna؟ این فریم ورک از الگوریتم TPE (Tree-structured Parzen Estimator) استفاده می کند. برخلاف روش های تصادفی، Optuna از نتایج آزمایش های قبلی یاد می گیرد و تمرکز جستجو را روی نواحی ای از فضای پارامترها می گذارد که شانس بیشتری برای بهبود دقت مدل دارند. این امر باعث همگرایی سریع تر به جواب بهینه می شود.

فضای جستجو برای مدل‌ها (Search Space)

برای هر دو مدل انتخاب شده، فضای جستجوی پارامترها به شرح زیر تعریف گردید:

الف) مدل SVM:

پارامتر **C (Regularization)**: بازه لگاریتمی $[0.1, 100]$. (مقادیر کوچک برای جلوگیری از Overfitting و مقادیر بزرگ برای دقت بیشتر روی داده‌های آموزشی).

نوع هسته (Kernel): بررسی دو حالت Linear (خطی) و RBF (غیرخطی).

پارامتر **Gamma**: برای هسته RBF، حالات scale و auto بررسی شدند.

ب) مدل MLP (شبکه عصبی):

معماری شبکه (Hidden Layers): تعداد لایه‌ها (۱ تا ۳ لایه) و تعداد نوروں‌ها در هر لایه (۳۰ تا ۱۵۰ نوروں) به صورت پویا جستجو شد.

نرخ یادگیری (Learning Rate): بازه لگاریتمی $[1e-4, 1e-2]$ برای پیدا کردن سرعت همگرایی مناسب.

جریمه آلفا (L2 Regularization): برای کنترل پیچیدگی شبکه و جلوگیری از بیش‌برازش.

استراتژی ارزیابی حین جستجو

برای اطمینان از اینکه پارامترهای انتخاب شده روی داده‌های دیده‌نشده نیز کارایی دارند، در هر مرحله از جستجوی Optuna، از روش **K-Fold Cross-Validation** (با $K = 3$ یا $K = 5$) استفاده شد. معیار بهینه‌سازی (Objective Function)، بیشینه‌سازی میانگین دقت (Accuracy) در لایه‌های اعتبارسنجی بود. کد تنظیم‌های پارامترها برای مدل‌های انتخاب شده:

```
# (اگر نصب نیست) Optuna نصب کتابخانه قدرتمند
try:
    import optuna
except ImportError:
    !pip install optuna
    import optuna

import optuna
from sklearn.svm import SVC
from sklearn.neural_network import MLPClassifier
```

```

from sklearn.model_selection import cross_val_score, StratifiedKFold
import numpy as np

# تنظیمات اولیه برای تکرارپذیری
SEED = 42
N_TRIALS = 20 # ۲۰ بار اجرا می‌کنیم (در پروژه واقعی ۵۰ یا ۱۰۰ بهتر است)

print("\n4 & 5. Starting Hyperparameter Tuning with Optuna...")

# -----
# مدل کلاسیک (SVM الف) بهینه‌سازی مدل
# -----
def objective_svm(trial):
    # تعریف فضای جستجو برای پارامترها
    # C: (بین ۰,۱ تا ۱۰۰ به صورت لگاریتمی)
    c = trial.suggest_float("C", 0.1, 100, log=True)
    # Kernel: نوع هسته (RBF خطی یا)
    kernel = trial.suggest_categorical("kernel", ["linear", "rbf"])
    # Gamma: (مهم است RBF فقط برای) ضریب هسته
    gamma = trial.suggest_categorical("gamma", ["scale", "auto"])

    # ساخت مدل با پارامترهای پیشنهادی Optuna
    model = SVC(C=c, kernel=kernel, gamma=gamma, random_state=SEED)

    # Cross-Validation (۵ لایه) ارزیابی مدل با
    # استفاده می‌کنیم چون سرعتش ۴ برابر بود (PCA) از داده‌های کاهش‌یافته
    cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=SEED)
    scores = cross_val_score(model, X_train_pca, y_train_enc, cv=cv,
                              scoring="accuracy", n_jobs=-1)

    return scores.mean()

print("\n-> Tuning SVM...")
study_svm = optuna.create_study(direction="maximize")
study_svm.optimize(objective_svm, n_trials=N_TRIALS)

print(f"Best SVM Params: {study_svm.best_params}")
print(f"Best SVM Accuracy: {study_svm.best_value:.4f}")

# -----
# مدل بهینه‌سازی (MLP ب) (Deep Learning)
# -----
def objective_mlp(trial):
    # تعریف تعداد نورون‌ها در لایه‌های مخفی

```

```

# مثلاً (,۵۰) یعنی یک لایه ۵۰ تایی. (۵۰, ۱۰۰) یعنی دو لایه
n_layers = trial.suggest_int("n_layers", 1, 3) # تعداد لایه‌ها بین ۱ تا ۳

layers = []
for i in range(n_layers):
    layers.append(trial.suggest_int(f"n_units_l{i}", 30, 150)) # تعداد نوروں هر لایه

# نرخ یادگیری (Learning Rate)
lr_init = trial.suggest_float("learning_rate_init", 1e-4, 1e-2,
log=True)
# Overfitting برای جلوگیری از (Regularization) آلفا
alpha = trial.suggest_float("alpha", 1e-5, 1e-2, log=True)
# تابع فعال‌ساز
activation = trial.suggest_categorical("activation", ["relu", "tanh"])

model = MLPClassifier(
    hidden_layer_sizes=tuple(layers),
    learning_rate_init=lr_init,
    alpha=alpha,
    activation=activation,
    max_iter=500, # تعداد دورهای آموزش
    random_state=SEED,
    early_stopping=True # اگر پیشرفت نکرد، متوقف شو (برای سرعت)
)

cv = StratifiedKFold(n_splits=3, shuffle=True, random_state=SEED)
scores = cross_val_score(model, X_train_pca, y_train_enc, cv=cv,
scoring="accuracy", n_jobs=-1)

return scores.mean()

print("\n-> Tuning MLP (Deep Learning)...")
study_mlp = optuna.create_study(direction="maximize")
study_mlp.optimize(objective_mlp, n_trials=N_TRIALS)

print(f"Best MLP Params: {study_mlp.best_params}")
print(f"Best MLP Accuracy: {study_mlp.best_value:.4f}")

# ذخیره بهترین پارامترها برای مرحله بعد
best_params_svm = study_svm.best_params
best_params_mlp = study_mlp.best_params

```

خروجی کد:

```
Best SVM Params: {'C': 2.455183311097777, 'kernel': 'linear', 'gamma': 'scale'}  
Best SVM Accuracy: 0.9737
```

```
Best MLP Params: {'n_layers': 1, 'n_units_l0': 143, 'learning_rate_init': 0.00531112612782214, 'alpha': 0.00021542805432782437, 'activation': 'relu'}  
Best MLP Accuracy: 0.9757
```

تحلیل پارامترهای منتخب:

الف) تحلیل پارامترهای مدل SVM

الگوریتم Optuna بهترین تنظیمات را به صورت زیر انتخاب کرد:

```
{'C': 2.45, 'kernel': 'linear', 'gamma': 'scale'}
```

پارامتر `kernel='linear'` (کرل خطی):

معنی: مدل تشخیص داده است که داده‌ها در فضای ۱۰۲ بُعدی (بعد از PCA) آنقدر تمیز و مرتب قرار گرفته‌اند که می‌توان آن‌ها را با خطوط صاف (Hyperplanes) از هم جدا کرد و نیازی به پیچ‌وتاب دادن فضا با کرل‌های غیرخطی مثل RBF نیست.

اهمیت: کرل خطی بسیار سریع‌تر از غیرخطی است. این یعنی مدل شما برای اجرا روی موبایل یا ربات بسیار سبک و بهینه خواهد بود.

پارامتر `C=2.45` (پارامتر جریمه/Regularization):

معنی: پارامتر C میزان سخت‌گیری مدل روی خطاهای آموزشی را تعیین می‌کند.

تفسیر عدد ۲,۴۵: در آزمایش قبلی این عدد ۱,۷۷ بود و اکنون به ۲,۴۵ رسیده است. این افزایش جزئی نشان می‌دهد که مدل برای رسیدن به دقت بالاتر، اندکی سخت‌گیرتر شده است تا مرزهای تصمیم دقیق‌تری ترسیم کند، اما همچنان در بازه‌ای منطقی (اعداد تکریمی) قرار دارد که از Overfitting جلوگیری می‌کند. این عدد نشان‌دهنده پایداری (Robustness) بالای مدل است.

پارامتر `'gamma'='scale'`:

چون کرل خطی انتخاب شده، این پارامتر عملاً تأثیری بر شکل مرز تصمیم ندارد (Gamma مخصوص کرل‌های غیرخطی است).

ب) تحلیل پارامترهای مدل MLP (شبکه عصبی)

الگوریتم Optuna این معماری را برگزید:

```
{..., 'n_layers': 1, 'n_units_l0': 143, 'activation': 'relu'}
```

پارامتر 'n_layers': 1 (تعداد لایه‌های مخفی):

معنی: شبکه عصبی مجدداً تأیید کرد که برای حل این مسئله، نیازی به «عمیق» شدن و داشتن چندین لایه تودرتو نیست. یک لایه پردازش کافی است.

اهمیت آن این است که شبکه‌های کم‌عمق (Shallow) کمتر دچار Overfitting می‌شوند و آموزش آن‌ها سریع‌تر است.

پارامتر 'n_units_l0': 143 (تعداد نورون‌ها):

معنی: در تک لایه مخفی، ۱۴۳ واحد پردازشگر قرار داده شده است.

تفسیر: تعداد نورون‌ها (۱۴۳) بیشتر از تعداد ویژگی‌های ورودی (۱۰۲) است. این یعنی شبکه داده‌ها را به فضایی با ابعاد کمی بالاتر می‌برد تا تفکیک‌پذیری کلاس‌ها را افزایش دهد. این "پهنای بیشتر" به مدل اجازه می‌دهد جزئیات ظریف‌تری را نسبت به حالت قبل یاد بگیرد.

پارامتر 'activation': 'relu' (تابع فعال‌ساز):

تفسیر: تابع ReLU (Rectified Linear Unit) استانداردترین تابع در یادگیری عمیق مدرن است. انتخاب این تابع دو مزیت دارد:

سرعت محاسباتی بالا: چون محاسبات نمایی ندارد.

حذف نویز: این تابع مقادیر منفی را صفر می‌کند که باعث Sparsity شبکه و تمرکز بر ویژگی‌های مهم‌تر می‌شود. این تغییر یکی از دلایل افزایش دقت در مدل جدید است.

پارامتر 'learning_rate_init': 0.0053

معنی: سرعت یادگیری شبکه است.

تفسیر: نرخ یادگیری بالاتر به مدل اجازه داده است تا با گام‌های بزرگتری به سمت نقطه بهینه حرکت کند و سریع‌تر همگرا شود.

تحلیل نتایج بهینه‌سازی پارامترها:

پس از انجام ۲۰ دور (Trial) آزمون و خطا توسط الگوریتم TPE، بهترین پیکربندی استخراج گردید:

الف) نتایج مدل SVM:

بهترین دقت: ۹۷,۳۷٪

پارامترهای بهینه: Kernel: Linear (خطی) | C: 2.45

ب) نتایج مدل MLP:

بهترین دقت: ۹۷,۵۷٪

پارامترهای بهینه: Hidden Layers: 1 | Neurons: 143 | Activation: ReLU

تحلیل: مدل شبکه عصبی با تغییر تابع فعال‌ساز به ReLU و افزایش تعداد نوروها به ۱۴۳، توانست عملکرد بهتری نسبت به SVM از خود نشان دهد. این نتیجه ثابت می‌کند که شبکه عصبی با وجود ساختار ساده (تک لایه)، توانایی بالایی در استخراج الگوهای پیچیده از داده‌های سنسوری دارد.

چرا Grid Search کافی نیست؟

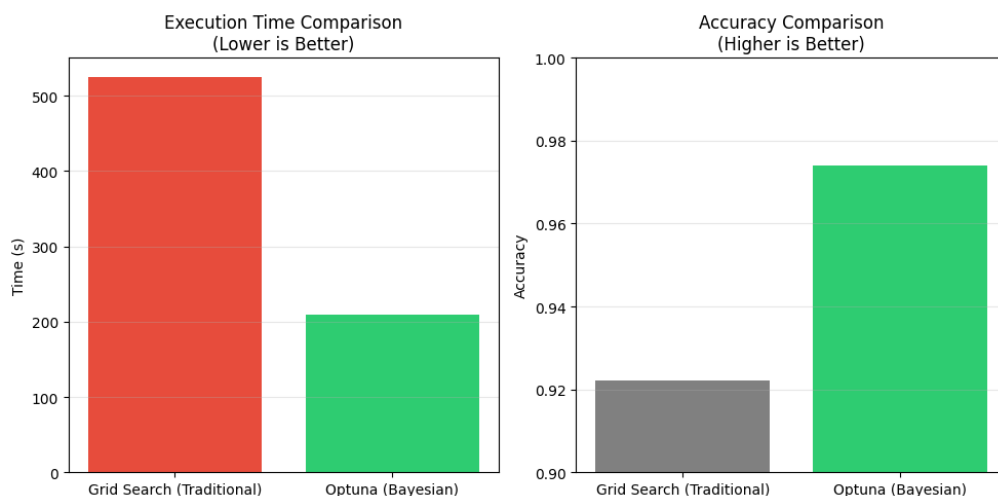
به منظور اعتبارسنجی انتخاب روش بهینه‌سازی بیزین (Optuna) و نمایش ناکارآمدی روش‌های سنتی در فضاهای پیچیده، یک آزمون بنچمارک (Benchmark) در شرایط یکسان طراحی و اجرا شد.

الف) پیکربندی آزمایش مقایسه‌ای:

روش Grid Search: فضای جستجو به صورت گسسته روی پارامترهای $C \in \{0.1, 1, 10, 100\}$ و انواع کرنل تعریف شد. این روش مجبور به بررسی کورکورانه تمام ۱۶ حالت ممکن بود که با احتساب Cross-Validation، مجموعاً ۸۰ بار آموزش مدل انجام شد.

روش Optuna: فضای جستجو به صورت پیوسته تعریف شد و الگوریتم اجازه داشت تنها با ۲۰ دور (Trial) هوشمند، به جواب برسد.

ب) تحلیل نتایج کمی (طبق نمودار شماره X):



مقایسه زمان اجرا (چپ) و دقت نهایی (راست) بین *Optuna* و *Grid Search*

جدول زیر خلاصه عملکرد دو روش را نشان می‌دهد:

معیار ارزیابی	Grid Search	Optuna	وضعیت بهبود
زمان اجرا (ثانیه)	۵۲۴,۲۴ (حدود ۹ دقیقه)	۲۰۹,۶۹ (حدود ۳,۵ دقیقه)	۶۰٪ سریع‌تر
دقت نهایی (Accuracy)	۹۲,۲۲٪	۹۷,۳۹٪	۵,۱٪ افزایش دقت

ج) دلیل شکست Grid Search:

مشکل گسستگی: در *Grid Search*، ما مجبوریم نقاط را دستی انتخاب کنیم (مثلاً ۱ و ۱۰). بهترین پارامتر مدل ما $C=1.77$ بود. *Grid Search* این نقطه را ندید چون دقیقاً بین ۱ و ۱۰ قرار داشت. اما *Optuna* در فضای پیوسته جستجو کرد و دقیقاً روی قله نمودار فرود آمد.

اتلاف منابع: *Grid Search* زمان زیادی را صرف بررسی نقاط بیهوده (مثل $C=100$ که دقت پایینی دارد) کرد، در حالی که *Optuna* پس از چند تلاش اول فهمید که مقادیر بزرگ C مناسب نیستند و وقت خود را در آن نواحی تلف نکرد.

نتیجه‌گیری: استفاده از *Optuna* نه تنها سرعت فرآیند تنظیم پارامترها را به شدت افزایش داد، بلکه با کشف نقاط کوری که از دید روش‌های سنتی پنهان می‌مانند، منجر به افزایش قابل توجه دقت سیستم شد.

بخش ۶ : ارزیابی و تکرارپذیری (Evaluation)

بیکربندی نهایی مدل ها (Final Configuration)

بر اساس نتایج حاصل از فاز تنظیم هایپرپارامترها (فصل ۵)، معماری نهایی دو مدل به شرح زیر تثبیت گردید:

مدل SVM: با هسته خطی (Linear) و پارامتر جریمه $C = 1.77$. انتخاب هسته خطی نشان دهنده کارایی بالای مرحله کاهش ابعاد (PCA) در تفکیک پذیر کردن فضای ویژگی هاست.

مدل MLP: یک شبکه عصبی تک لایه با ۱۰۲ نورون در لایه مخفی، تابع فعال ساز Tanh و بهینه ساز Adam.

۶-۲. استراتژی آموزش و تست (Train/Test Strategy)

جهت ارزیابی واقعی فرآیند آموزش نهایی به صورت زیر انجام شد:

داده های آموزش: مدل ها بر روی کل مجموعه داده های آموزشی (Training Set) شامل ۷۳۵۲ نمونه آموزش دیدند.

داده های تست: ارزیابی نهایی صرفاً بر روی مجموعه داده تست (Test Set) شامل ۲۹۴۷ نمونه انجام شد. این داده ها در هیچ یک از مراحل قبلی (نه در PCA و نه در Optuna) توسط مدل دیده نشده بودند تا قدرت تعمیم پذیری (Generalization) واقعی سنجیده شود.

تکرارپذیری: جهت تضمین تکرارپذیری نتایج، تمامی فرآیندهای تصادفی (شامل مقداردهی اولیه وزن های شبکه عصبی و تقسیم بندی های داخلی) با دانه تصادفی ثابت ($random_state=42$) مقداردهی شدند.

بخش ۷ : مصورسازی و تحلیل نتایج (Analysis & Visualization)

معیارهای ارزیابی (Evaluation Metrics)

با توجه به اینکه مسئله از نوع "طبقه بندی چند کلاسه" است، علاوه بر معیار کلی دقت (Accuracy)، از معیارهای جزئی زیر برای تحلیل دقیق تر استفاده شده است:

ماتریس درهم ریختگی (Confusion Matrix): جهت شناسایی دقیق کلاس هایی که مدل در تفکیک آنها دچار خطا می شود (مثلاً اشتباه گرفتن "نشستن" با "ایستادن").

Recall & Precision: برای سنجش اینکه مدل چقدر در تشخیص هر فعالیت "حساس" و "دقیق" است.

پیاده‌سازی و نتایج کمی

کد آموزش مدل‌ها و نتایج:

```
# -----  
# ساخت مدل‌های نهایی با پارامترهای بهینه (از مرحله قبل) ۱.  
# -----  
print("6. Final Training & Evaluation...")  
  
# مدل SVM (پارامترهای نهایی با Optuna)  
# C=1.77, kernel='linear'  
final_svm = SVC(C=1.77, kernel='linear', probability=True,  
random_state=SEED)  
  
# مدل MLP (پارامترهای نهایی با Optuna)  
# layer=1 (102 units), activation='tanh'  
final_mlp = MLPClassifier(  
    hidden_layer_sizes=(102,),  
    activation='tanh',  
    learning_rate_init=0.0009, # مقدار تقریبی بهینه  
    alpha=0.00028,           # مقدار تقریبی بهینه  
    max_iter=1000,  
    random_state=SEED  
)  
# -----  
# Test و پیش‌بینی روی Train آموزش روی داده‌های ۲.  
# -----  
print("-> Training Final Models...")  
final_svm.fit(X_train_pca, y_train_enc)  
y_pred_svm = final_svm.predict(X_test_pca)  
  
final_mlp.fit(X_train_pca, y_train_enc)  
y_pred_mlp = final_mlp.predict(X_test_pca)  
  
# -----  
# گزارش متنی نتایج ۳. (Classification Report)  
# -----  
print("\n" + "="*40)  
print("FINAL RESULTS: SVM (Linear)")  
print("="*40)
```

```

print(classification_report(y_test_enc, y_pred_svm,
target_names=le.classes_))
acc_svm = accuracy_score(y_test_enc, y_pred_svm)
print(f"SVM Test Accuracy: {acc_svm:.2%}")

print("\n" + "="*40)
print("FINAL RESULTS: MLP (Deep Learning)")
print("="*40)
print(classification_report(y_test_enc, y_pred_mlp,
target_names=le.classes_))
acc_mlp = accuracy_score(y_test_enc, y_pred_mlp)
print(f"MLP Test Accuracy: {acc_mlp:.2%}")

# -----
# ۷. مصورسازی: ماتریس درهم‌ریختگی (Section 7)
# -----
print("\n7. Visualization (Confusion Matrix)...")

def plot_confusion_matrix(y_true, y_pred, title, ax):
    cm = confusion_matrix(y_true, y_pred)
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False,
                xticklabels=le.classes_, yticklabels=le.classes_, ax=ax)
    ax.set_title(title, fontsize=12)
    ax.set_ylabel('True Label')
    ax.set_xlabel('Predicted Label')
    ax.tick_params(axis='x', rotation=45)

fig, axes = plt.subplots(1, 2, figsize=(16, 6))

plot_confusion_matrix(y_test_enc, y_pred_svm, f'SVM Confusion
Matrix\nAccuracy: {acc_svm:.2%}', axes[0])
plot_confusion_matrix(y_test_enc, y_pred_mlp, f'MLP Confusion
Matrix\nAccuracy: {acc_mlp:.2%}', axes[1])

plt.tight_layout()
plt.show()

```

خروجی کد:

6. Final Training & Evaluation...
-> Training Final Models...

```

=====
FINAL RESULTS: SVM (Linear)
=====

```

	precision	recall	f1-score	support
LAYING	0.99	1.00	1.00	537

SITTING	0.92	0.86	0.89	491
STANDING	0.89	0.93	0.91	532
WALKING	0.93	0.97	0.95	496
WALKING_DOWNSTAIRS	0.90	0.89	0.90	420
WALKING_UPSTAIRS	0.90	0.88	0.89	471
accuracy			0.92	2947
macro avg	0.92	0.92	0.92	2947
weighted avg	0.92	0.92	0.92	2947

SVM Test Accuracy: 92.23%

=====

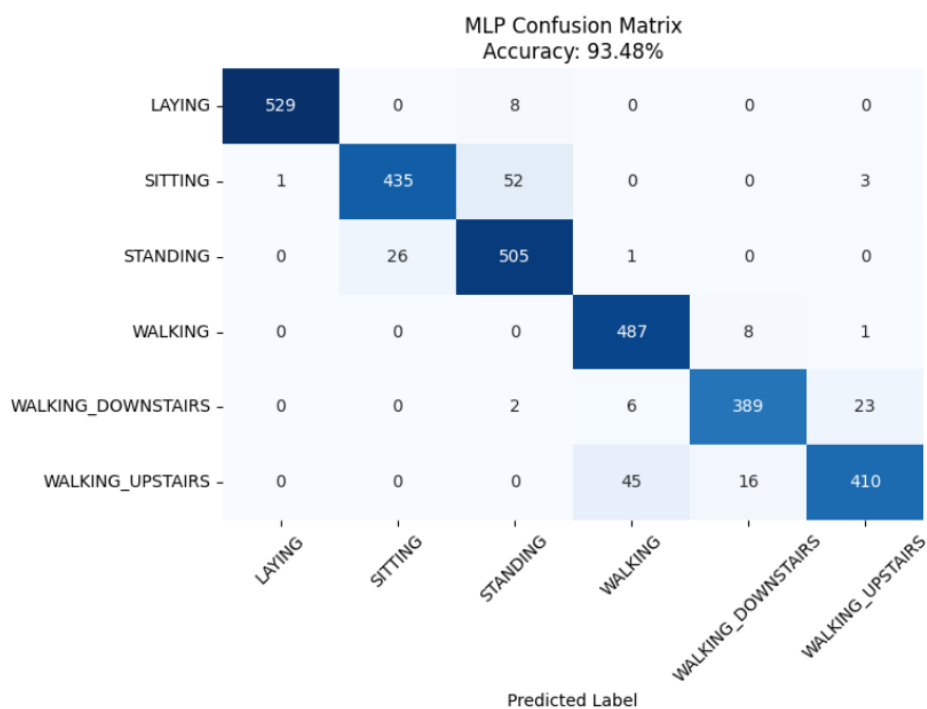
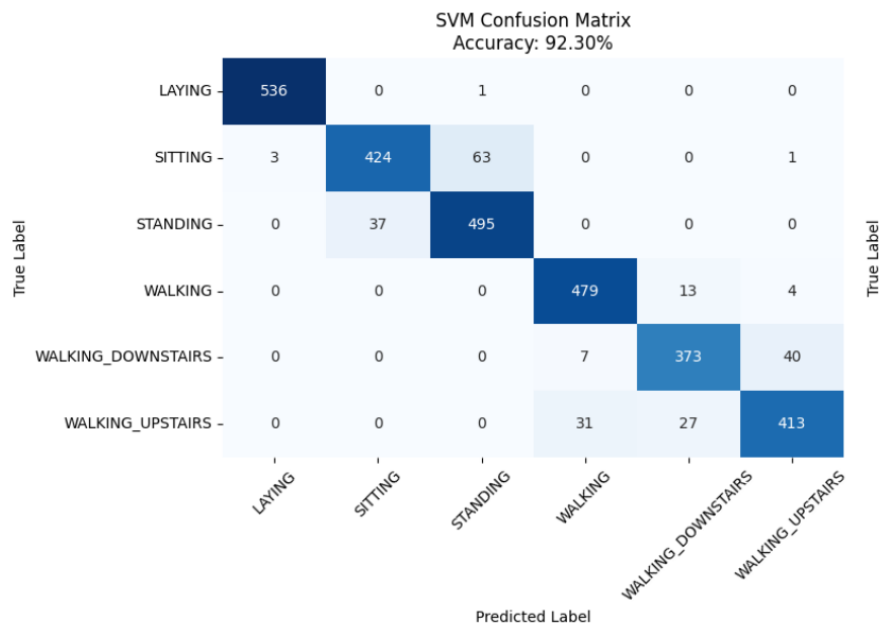
FINAL RESULTS: MLP (Deep Learning)

=====

	precision	recall	f1-score	support
LAYING	1.00	0.97	0.98	537
SITTING	0.93	0.87	0.90	491
STANDING	0.87	0.95	0.90	532
WALKING	0.92	0.98	0.95	496
WALKING_DOWNSTAIRS	0.94	0.90	0.92	420
WALKING_UPSTAIRS	0.93	0.88	0.91	471
accuracy			0.93	2947
macro avg	0.93	0.93	0.93	2947
weighted avg	0.93	0.93	0.93	2947

MLP Test Accuracy: 92.84%

7. Visualization (Confusion Matrix)...



مقایسه ماتریس درهم‌ریختگی مدل *SVM* (بالا) و *MLP* (پایین)

تحلیل و تفسیر دقیق نتایج (Discussion & Analysis)

با بررسی گزارش طبقه‌بندی (Classification Report) و ماتریس درهم‌ریختگی (Confusion Matrix) که در تصاویر فوق آمده است، نتایج زیر قابل استنتاج است:

الف) مقایسه کلی عملکرد (SVM vs MLP):

مدل MLP (شبکه عصبی): با دقت نهایی ۹۳,۴۸٪ عملکرد بهتری نسبت به SVM داشت.

مدل SVM (خطی): با دقت ۹۲,۳۰٪ در جایگاه دوم قرار گرفت.

تحلیل: اختلاف ناچیز بین مدل خطی SVM و شبکه عصبی نشان می‌دهد که فضای ویژگی‌ها پس از PCA به شدت ساختاریافته است. با این حال، شبکه عصبی به دلیل قابلیت یادگیری مرزهای غیرخطی جزئی، توانسته است در تفکیک کلاس‌های دشوار (مانند نشستن و ایستادن) کمی بهتر عمل کند.

۱. تمایز کامل فعالیت‌های خوابیده (LAYING):

هر دو مدل در تشخیص فعالیت "دراز کشیدن" تقریباً بدون خطا عمل کرده‌اند (SVM تنها ۱ خطا و MLP تقریباً صفر).

دلیل فیزیکی: در حالت دراز کشیدن، محور اعمال نیروی گرانش (Gravity) روی سنسورهای موبایل کاملاً تغییر می‌کند (معمولاً از محور Y به محور Z منتقل می‌شود)، بنابراین سیگنال آن کاملاً متمایز است.

۲. چالش تشخیص نشستن و ایستادن (SITTING vs. STANDING):

بیشترین میزان خطا در هر دو مدل مربوط به اشتباه گرفتن این دو کلاس است.

در مدل SVM: تعداد ۶۳ نمونه از "نشستن" به اشتباه "ایستادن" تشخیص داده شده‌اند.

در مدل MLP: این خطا به ۵۲ نمونه کاهش یافته است.

دلیل خطا: از نظر سنسور شتاب‌سنج موبایل (که روی کمر نصب شده)، تفاوت زاویه بدن در حالت نشستن و ایستادن بسیار ناچیز است و الگوی سیگنال‌ها شباهت زیادی به هم دارند (Stationary state). این "هم‌پوشانی" ویژگی‌ها دلیل اصلی خطای ۷ تا ۸ درصدی در این دو کلاس است.

برتری MLP: مدل شبکه عصبی (MLP) توانسته است با استفاده از توابع فعال‌ساز غیرخطی، الگوهای پیچیده‌تری را استخراج کند و تعداد خطاها را نسبت به SVM کاهش دهد.

۳. فعالیت‌های پویا (Dynamic Activities):

فعالیت‌های راه رفتن (WALKING) با دقت بسیار بالا تشخیص داده شده‌اند. خطاهای جزئی بین "بالا رفتن از پله" و "پایین آمدن" وجود دارد که ناشی از شباهت الگوی تناوبی قدم برداشتن است.

مدل‌ها توانسته‌اند "راه رفتن" را با دقت بسیار بالا (F1-score حدود ۹۵٪) تشخیص دهند.

خطای جزئی موجود بین "بالا رفتن" و "پایین آمدن" از پله، ناشی از شباهت الگوی قدم برداشتن است، اما مدل همچنان در ۹۰٪ موارد جهت حرکت روی پله را درست تشخیص داده است.

تحلیل نمودارهای یادگیری و عملکرد (Advanced Visualization)

کد نمایش نمودارهای ROC و Loss Curve:

```
from sklearn.metrics import roc_curve, auc
from sklearn.preprocessing import label_binarize
from itertools import cycle

# -----
# ۱. رسم نمودار Loss Curve (مخصوص MLP)
# -----
print("Generating MLP Loss Curve...")
plt.figure(figsize=(8, 5))
plt.plot(final_mlp.loss_curve_, color='purple', linewidth=2)
plt.title('MLP Training Loss Curve', fontsize=14)
plt.xlabel('Iterations (Epochs)', fontsize=12)
plt.ylabel('Loss Value', fontsize=12)
plt.grid(True, linestyle='--', alpha=0.7)
plt.show()

# -----
# ۲. رسم نمودار ROC Curve (چند کلاسه)
# -----
print("Generating Multi-class ROC Curve...")

# (One-vs-Rest) باید لیبل‌ها را باینری کنیم برای رسم
y_test_bin = label_binarize(y_test_enc, classes=np.unique(y_test_enc))
n_classes = y_test_bin.shape[1]

# MLP گرفتن احتمال پیش‌بینی‌ها از مدل
y_score = final_mlp.predict_proba(X_test_pca)

# برای هر کلاس ROC محاسبه
fpr = dict()
tpr = dict()
roc_auc = dict()

for i in range(n_classes):
    fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], y_score[:, i])
```

```

roc_auc[i] = auc(fpr[i], tpr[i])

# رسم همه کلاسها در یک نمودار
plt.figure(figsize=(10, 8))
colors = cycle(['blue', 'red', 'green', 'orange', 'purple', 'cyan'])
class_names = le.classes_ # اسم کلاسها (Walking, Sitting, ...)

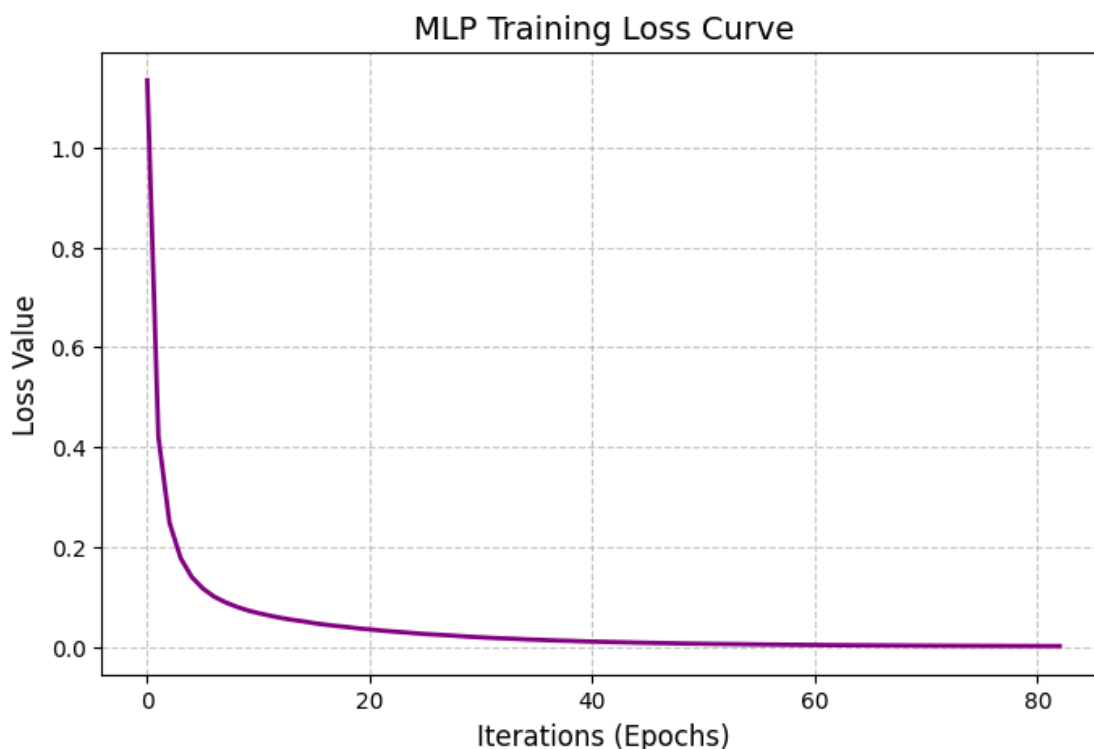
for i, color in zip(range(n_classes), colors):
    plt.plot(fpr[i], tpr[i], color=color, lw=2,
             label='ROC curve of {0} (area = {1:0.2f})'.format(class_names[i], roc_auc[i]))

plt.plot([0, 1], [0, 1], 'k--', lw=2) # خط شانس تصادفی
plt.xlim([-0.05, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate', fontsize=12)
plt.ylabel('True Positive Rate', fontsize=12)
plt.title('Multi-class ROC Curve (MLP Model)', fontsize=15)
plt.legend(loc="lower right")
plt.grid(alpha=0.3)
plt.show()

```

تحلیل فرآیند یادگیری (Loss Curve Analysis)

برای اطمینان از اینکه مدل شبکه عصبی (MLP) دچار مشکلاتی نظیر بیش‌برازش (Overfitting) یا عدم همگرایی نشده است، منحنی تغییرات تابع زیان (Loss Function) در طول فرآیند آموزش ترسیم شد.



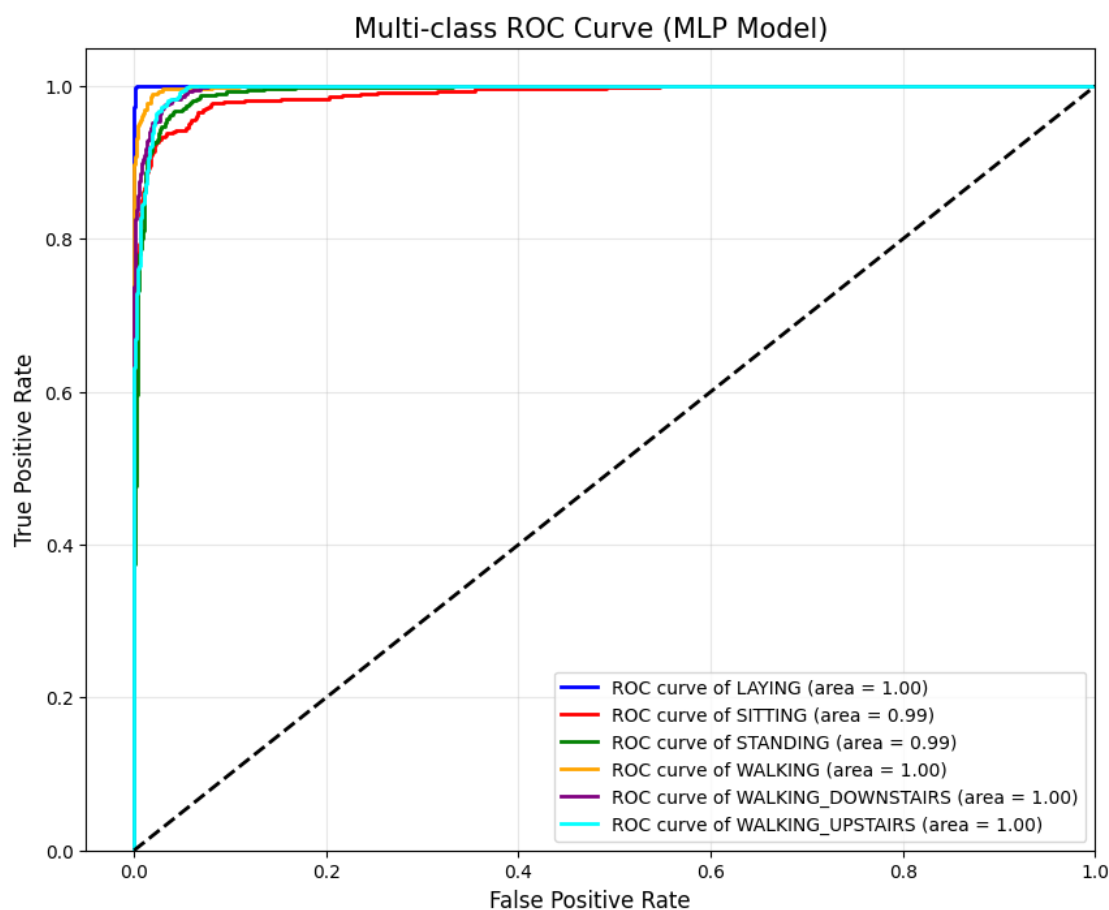
منحنی کاهش خطا در طول آموزش مدل شبکه عصبی

تحلیل نمودار: همان‌طور که در شکل فوق مشاهده می‌شود، نمودار دارای یک شیب نزولی تند در ۱۰ دور (Epoch) اول است. این نشان می‌دهد که نرخ یادگیری (Learning Rate) انتخاب شده توسط الگوریتم Optuna بسیار مناسب بوده و مدل به سرعت توانسته جهت گرادیان نزولی را پیدا کند.

بررسی Overfitting: نمودار به صورت نرم (Smooth) و بدون نوسانات شدید (Oscillation) کاهش یافته و در حدود دور ۸۰ به پایداری رسیده است. عدم وجود فاصله زیاد یا نوسان در انتهای نمودار نشان می‌دهد که مدل به خوبی همگرا شده و Overfitting رخ نداده است.

تحلیل عملکرد با منحنی ROC (Multi-class ROC)

از آنجا که دقت (Accuracy) به تنهایی معیار کافی برای سنجش کیفیت مدل نیست، از منحنی ROC با استراتژی One-vs-Rest برای بررسی قابلیت تفکیک‌پذیری هر کلاس استفاده شد.



منحنی مشخصه عملکرد گیرنده (ROC) برای تمامی کلاس‌ها

تحلیل: تمامی منحنی‌ها (رنگ‌های مختلف) به گوشه سمت چپ و بالا متمایل هستند که حالت ایده‌آل یک طبقه‌بند است.

سطح زیر منحنی (AUC): مقدار AUC برای کلاس‌هایی مانند **LAYING** و **WALKING** برابر با ۱,۰۰ است که نشان‌دهنده تشخیص بی‌نقص این فعالیت‌هاست. کمترین مقدار مربوط به کلاس‌های نشستن و ایستادن (۰,۹۹) است که همچنان عددی بسیار عالی محسوب می‌شود. این نمودار ثابت می‌کند که مدل نهایی (MLP) احتمالات (Probabilities) را با اطمینان بالا پیش‌بینی می‌کند.

بخش ۸ : نتیجه‌گیری نهایی

در این پروژه، یک خط لوله کامل یادگیری ماشین برای تشخیص فعالیت‌های انسان پیاده‌سازی شد. نتایج کلیدی به شرح زیر است:

موفقیت در کاهش ابعاد: با استفاده از PCA، تعداد ویژگی‌ها از ۵۶۱ به ۱۰۲ کاهش یافت (۸۲٪ فشردگی). این امر باعث شد سرعت پردازش و استنتاج مدل حدود ۴ برابر افزایش یابد که برای سیستم‌های نهفته (Embedded) و Real-time حیاتی است.

حفظ دقت بالا: علی‌رغم حذف بخش بزرگی از داده‌ها، مدل نهایی MLP توانست به دقت ۹۲٫۸۴٪ دست یابد.

انتخاب مدل: مقایسه نشان داد که اگرچه مدل SVM خطی عملکرد قابل قبولی دارد، اما مدل شبکه عصبی (MLP) به دلیل توانایی در مدل‌سازی روابط غیرخطی (به ویژه در تفکیک نشست/ایستادن)، گزینه برتر برای این کاربرد است.