# Datamining Report:
# The Human Activity Recognition 70+

## Prepared by:

Youssef Mahmoud Anis - 192100029

Youssef Muhammed Gabr -192100069

## INTRODUCTION:

Human Activity Recognition (HAR) has become a crucial field of study, particularly for applications in healthcare and elderly care. With the growing aging population, monitoring physical activities can help prevent falls, detect anomalies, and improve overall well-being. HAR models typically rely on sensor-based data to classify various activities performed by individuals.

The Human Activity Recognition 70+ (HAR70+) dataset for older-adult subjects (70-95 years old) wearing two 3-axial accelerometers for around 40 minutes during a semi-structured free-living protocol. The sensors were attached to the right front thigh and lower back.

## DATASET:

This Dataset includes the test of 18 different elderly patients recording their x, y and z coordinates of their thigh and their lower back with timestamps to know when the test takes places. Sensors Used: **Two Axivity AX3 accelerometers**. Sampling Rate: **50Hz.**

The dataset contains the following annotated activities with the corresponding coding scheme:
- Walking
- Shuffling
- Stairs (ascending)
- Stairs (descending)
- Standing
- Sitting
- Lying

And this dataset doesn't have any missing values.

### EXPANDED DATASET ANALYSIS:

- Demographic Diversity: The dataset could benefit from including more diverse participants in terms of gender, health conditions, and mobility levels to ensure it represents the broader elderly population.

- Longitudinal Data: Future data collection could extend over days or weeks to capture more natural and varied activities, rather than just 40 minutes.

- Multi-Sensor Integration: Adding sensors like heart rate monitors, gyroscopes, or environmental sensors could provide a more comprehensive view of the participants' activities and health.

## RESEARCH STUDIES AND FINDINGS

i. *Validation of the HAR70+ Model*

A study published in Sensors ([mdpi.com](mdpi.com)) evaluated the performance of the HAR70+ model and compared it with an existing model trained on younger adults. The results demonstrated that models trained with the HAR70+ dataset were more accurate in recognizing activities performed by older adults, particularly for those using walking aids.

ii. *Self-Supervised Learning for HAR*

A research paper in Applied Intelligence ([springer.com](springer.com)) introduced SelfPAB, a self-supervised learning approach that improves HAR model accuracy by pre-training on large-scale accelerometer datasets. The study found that pre-training significantly enhanced model performance when applied to HAR70+ data.

iii. *Machine and Deep Learning for Activity Recognition*

A study on ResearchGate ([researchgate.net](researchgate.net)) applied machine learning and deep learning models to classify activities in the HAR70+ dataset. The models used included:

- Random Forest (RF)
- Extreme Gradient Boosting (XGBoost)
- Logistic Regression (LR)
- K-Nearest Neighbors (KNN)
- Stochastic Gradient Descent (SGD)
- Deep Neural Networks (DNN)
- Long Short-Term Memory Networks (LSTM)

The results showed that LSTM achieved the highest accuracy (98%) for classifying activities, indicating its effectiveness for time-series data analysis.

iv. *FPGA Optimization for HAR Models*

Research on FPGA-optimized machine learning models explored efficient implementations of HAR models for real-time applications. The study focused on optimizing Support Vector Machine (SVM) and Convolutional Neural Network (CNN) models for hardware deployment.

*v.* *Comparison of CNN and SVM Models*

Another study compared CNN and SVM models for HAR tasks, utilizing the HAR70+ dataset. The research demonstrated that CNN models generally outperformed SVM models in accuracy, while SVM models were more resource-efficient for embedded applications.

## METHODOLOGIES USED IN RESEARCH

The research studies mentioned above implemented various techniques to improve HAR performance. These include:

1. **Feature Engineering**: Extracting relevant features from raw accelerometer data.

2. **Supervised Learning**: Training classification models using labeled HAR70+ data.

3. **Deep Learning**: Implementing LSTM and CNN architectures for better sequence modeling.

4. **Hardware Optimization**: Developing efficient models for real-time, low-power applications.

## Advanced Machine Learning Techniques:

- Transfer Learning: Adapting models trained on younger populations to the HAR70+ dataset could improve accuracy with less data.

- Ensemble Learning: Combining multiple models (e.g., Random Forest, LSTM, and CNN) to create an ensemble model that could outperform individual models.

- Explainable AI (XAI): Using techniques like SHAP or LIME to make models more interpretable, especially for healthcare applications.

## HANDLING EXTREME VALUES:

Extreme values in sensor data may arise due to sudden movements, sensor errors, or calibration issues.      Solutions include:

1. Outlier Detection**: Using Z-score analysis or IQR (Interquartile Range) to identify and remove extreme values.**

2. Smoothing Techniques**: Applying a low-pass filter or moving average to smooth out abrupt spikes.**

3. Data Imputation: **Replacing extreme values with the mean or median of surrounding data points.**

## COMPARISON WITH OTHER HAR DATASETS

Several HAR datasets exist, but few focus specifically on older adults. Below is a comparison:

- **HARTH Dataset**: Focuses on younger adults; HAR70+ is tailored for elderly activity recognition.

- **MobiAct**: Includes fall detection but lacks continuous daily activity tracking like HAR70+.

- **PAMAP2**: Captures a broader set of activities but is not specialized for the elderly.

  The HAR70+ dataset improves upon previous datasets by providing data specifically collected from an older population, addressing the gap in elderly-specific activity recognition research.

## CROSS-DATASET VALIDATION:

- Suggest validating HAR70+ models on other datasets (e.g., MobiAct, PAMAP2) to test their generalizability.

- Propose combining HAR70+ with other datasets to create a more comprehensive dataset that covers a wider range of activities and populations.

## EXPERIMENTAL RESULTS & ANALYSIS

We added some of the research papers we found and asked AI bots to give us a brief analysis by showing accuracy of each paper:

| Model | Accuracy |
|---|---|
| Random Forest | 91% |
| XGBoost | 94% |
| LSTM | 98% |
| CNN | 96% |
| SVM | 92% |

## ETHICAL CONSIDERATIONS & CHALLENGES

- Privacy Concerns: Data collection involves tracking movement, raising concerns about user privacy.

- Wearable Discomfort: Participants may experience discomfort wearing sensors, which can negatively affect data quality.

- Data Anonymization: Advanced anonymization techniques should be used to protect participant identities while maintaining data utility.

- Informed Consent: Clear and transparent informed consent processes are crucial, especially for elderly participants who may have cognitive impairments.

- Bias Mitigation: Address potential biases in the dataset (e.g., underrepresentation of certain groups) and propose methods to mitigate them.

## PRACTICAL APPLICATIONS OF HAR70+ DATA:

The HAR70+ dataset provides valuable insights that can be applied in various real-world scenarios:

- **Healthcare Monitoring:** Assisting doctors and caregivers in tracking elderly movement, detecting falls, and providing early interventions.

- **Smart Homes & Assistive Technology**: AI-powered home automation adjusting alarms, lighting, and alerts based on activity levels.

- **Personalized Fitness & Rehabilitation**: Designing exercise and therapy programs tailored to older adults based on movement data.

- **Fall Prevention Systems**: Identifying abnormal movements that may indicate a fall, triggering emergency responses.

- **Wearable Device Enhancements**: Improving smartwatches and medical devices for more accurate elderly monitoring.

- **Urban Planning & Accessibility**: Using movement data to design elderly-friendly public spaces and infrastructure.

- **AI-Based Elderly Assistance**: Developing AI-powered assistants and robots to aid elderly individuals in daily activities.

- Mental Health Monitoring: Detecting signs of mental health issues (e.g., depression, anxiety) through changes in activity patterns.

- Chronic Disease Management: Monitoring and managing chronic conditions like Parkinson's disease, arthritis, or diabetes by tracking movement and activity levels.

- Social Interaction Tracking: Analyzing social interactions among elderly individuals to understand loneliness and social isolation.

## PROGRAMS TO BE USED:

We will use Anaconda (Jupyter Notebook) and Orange Data Mining for data analysis. Additionally, we propose using TensorFlow and PyTorch for deep learning model development, and SHAP or LIME for model interpretability.

## FUTURE RESEARCH DIRECTIONS

- Personalized Models: Develop personalized HAR models that adapt to individual users' movement patterns, improving accuracy and relevance.

- Integration with IoT: Explore how HAR70+ data could be integrated with IoT devices in smart homes to create more responsive and adaptive environments for elderly individuals.

- Cross-Disciplinary Collaboration: Encourage collaboration between computer scientists, healthcare professionals, and gerontologists to address the multifaceted challenges of elderly care.

- Edge AI and Federated Learning: Deploy HAR models on edge devices for real-time activity recognition and use federated learning to train models across multiple devices without sharing raw data.

- Policy and Regulatory Considerations: Develop data governance frameworks to ensure ethical data collection, storage, and usage, and ensure compliance with regulations like GDPR or HIPAA.

## CODING STEPS USED ALL ALONG THIS PROJECT

Libraries Used IN This Project:

- import numpy as np
- import pandas as pd
- import matplotlib.pyplot as plt
- import seaborn as sns
- from sklearn.preprocessing import MinMaxScaler
- from sklearn.metrics import confusion_matrix, classification_report
- from xgboost import XGBClassifier
- from sklearn.model_selection import train_test_split
- from sklearn.preprocessing import StandardScaler
- from sklearn.tree import DecisionTreeClassifier
- from sklearn.naive_bayes import GaussianNB
- from sklearn.neighbors import KNeighborsClassifier
- from sklearn.svm import SVC
- from sklearn.neural_network import MLPClassifier
- from sklearn.linear_model import LogisticRegression
- from sklearn.ensemble import RandomForestClassifier
- from sklearn.metrics import accuracy_score
- from sklearn.linear_model import LogisticRegression
- from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, ConfusionMatrixDisplay
- from sklearn.tree import DecisionTreeClassifier, plot_tree
- from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
- from sklearn.naive_bayes import GaussianNB
- from sklearn.svm import SVC
- from sklearn.decomposition import PCA

# 1. Explore the dataset's features and target variable.

```
[184]: file_ids = range(501, 519)
       patients = [pd.read_csv(f"{pid}.csv") for pid in file_ids]

       # Combine all patients into one DataFrame
       all_data = pd.concat(patients, ignore_index=True)
```

```
[212]: print("First 10 rows of the dataset:\n")
       display(all_data.head(10))
```

First 10 rows of the dataset:

|   | timestamp | back_x | back_y | back_z | thigh_x | thigh_y | thigh_z | label |
|---|-----------|--------|--------|--------|---------|---------|---------|-------|
| 0 | 2021-03-24 14:42:03.839 | -0.999023 | -0.063477 | 0.140625 | -0.980469 | -0.112061 | -0.048096 | 6 |
| 1 | 2021-03-24 14:42:03.859 | -0.980225 | -0.079346 | 0.140625 | -0.961182 | -0.121582 | -0.051758 | 6 |
| 2 | 2021-03-24 14:42:03.880 | -0.950195 | -0.076416 | 0.140625 | -0.949463 | -0.080566 | -0.067139 | 6 |
| 3 | 2021-03-24 14:42:03.900 | -0.954834 | -0.059082 | 0.140381 | -0.957520 | -0.046143 | -0.050781 | 6 |
| 4 | 2021-03-24 14:42:03.920 | -0.972412 | -0.042969 | 0.142822 | -0.977051 | -0.023682 | -0.026611 | 6 |
| 5 | 2021-03-24 14:42:03.940 | -0.988770 | -0.026123 | 0.157227 | -0.984863 | -0.042725 | -0.032715 | 6 |
| 6 | 2021-03-24 14:42:03.960 | -1.001953 | -0.016113 | 0.162109 | -0.992920 | -0.075439 | -0.024170 | 6 |
| 7 | 2021-03-24 14:42:03.980 | -1.000488 | -0.035400 | 0.191406 | -0.996338 | -0.072754 | -0.013428 | 6 |
| 8 | 2021-03-24 14:42:04.000 | -0.996826 | -0.056152 | 0.187500 | -0.974609 | -0.060303 | -0.015625 | 6 |
| 9 | 2021-03-24 14:42:04.019 | -0.978027 | -0.083252 | 0.187500 | -0.966797 | -0.062500 | -0.015625 | 6 |

```
[188]: print(f"Dataset contains {all_data.shape[0]} rows and {all_data.shape[1]} columns.")
```

Dataset contains 2259597 rows and 8 columns.

We started by combining 18 different CSV files, each representing data from a different patient. These files were identified by a unique patient ID (PID), ranging from 501 to 518. Using a loop, we read each file using the read_csv function and stored them in a list called patients.

Next, we used pd.concat() to merge all these individual patient datasets into a single DataFrame called all_data, which makes it easier to work with the complete dataset as one unit.

To understand the structure of the dataset, we displayed the first 10 rows. This gave us a preview of the types of data we are working with. After that, we used .shape to check the total number of rows and columns in the dataset. It turns out the dataset contains approximately **2.25 million rows and 8 columns**.

The dataset includes:
- **timestamp**: the exact time of the data recording
- **back_x, back_y, back_z**: accelerometer readings from the lower back (in x, y, and z directions)
- **thigh_x, thigh_y, thigh_z**: accelerometer readings from the right front thigh
- **label**: the activity being performed, which is categorized as:

    1: walking
    3: shuffling
    4: stairs (ascending)
    5: stairs (descending)

6: standing

7: sitting

8: lying

This setup allows us to explore how different body movements are represented in the sensor data and how we might use this for activity classification.

## 2. Handle missing values (if any) and outliers.

```
[190]:  #mean
        bx_mean = all_data ['back_x'].mean()
        by_mean = all_data ['back_y'].mean()
        bz_mean = all_data ['back_z'].mean()
        tx_mean = all_data ['thigh_x'].mean()
        ty_mean = all_data ['thigh_y'].mean()
        tz_mean = all_data ['thigh_z'].mean()

        print("\nBack Mean scores: ")
        print(f"back_x: {bx_mean}, back_y: {by_mean}, back_z: {bz_mean}")

        print("\nThigh Mean scores: ")
        print(f"thigh_x: {tx_mean}, thigh_y: {ty_mean}, thigh_z: {tz_mean}")


        Back Mean scores:
        back_x: -0.8699343930886793, back_y: -0.03316849906111574, back_z: 0.023424907512711337

        Thigh Mean scores:
        thigh_x: -0.6796213034173795, thigh_y: 0.0027747422088983133, thigh_z: -0.384121996420158
```

```
[129]:  #missing values
        print("Missing values in each column:\n")
        print(all_data.isnull().sum())

        Missing values in each column:

        timestamp    0
        back_x       0
        back_y       0
        back_z       0
        thigh_x      0
        thigh_y      0
        thigh_z      0
        label        0
        dtype: int64
```

We calculated the **mean values** for each of the sensor data columns:

- back_x, back_y, back_z (lower back sensors)

- thigh_x, thigh_y, thigh_z (thigh sensors)

This gives us a general idea of the average sensor readings, which can be useful later for normalization or handling missing data.

After calculating the means, we checked for any **missing values** in the dataset using **isnull().sum()**. The result showed that there are **no missing values** in any of the columns, which means the dataset is clean and we won't be using the mean here, so now we will go find outliers.

```
[10]: # IQR for Back
      Q1_back = all_data[['back_x', 'back_y', 'back_z']].quantile(0.25)
      Q3_back = all_data[['back_x', 'back_y', 'back_z']].quantile(0.75)
      IQR_back = Q3_back - Q1_back

      print("Back IQR values:\n", IQR_back)

      mask_back = ~((all_data[['back_x', 'back_y', 'back_z']] < (Q1_back - 1.5 * IQR_back)) |
                    (all_data[['back_x', 'back_y', 'back_z']] > (Q3_back + 1.5 * IQR_back))).any(axis=1)

      # IQR for Thigh
      Q1_thigh = all_data[['thigh_x', 'thigh_y', 'thigh_z']].quantile(0.25)
      Q3_thigh = all_data[['thigh_x', 'thigh_y', 'thigh_z']].quantile(0.75)
      IQR_thigh = Q3_thigh - Q1_thigh

      print("\nThigh IQR values:\n", IQR_thigh)

      mask_thigh = ~((all_data[['thigh_x', 'thigh_y', 'thigh_z']] < (Q1_thigh - 1.5 * IQR_thigh)) |
                     (all_data[['thigh_x', 'thigh_y', 'thigh_z']] > (Q3_thigh + 1.5 * IQR_thigh))).any(axis=1)

      # Combine both masks
      final_mask = mask_back & mask_thigh
      all_data_clean = all_data[final_mask]

      print("\nShape after IQR-based outlier removal:", all_data_clean.shape)
      print(f"Total outliers removed: {all_data.shape[0] - all_data_clean.shape[0]}")
```

```
Back IQR values:
 back_x    0.164795
back_y    0.160400
back_z    0.587646
dtype: float64

Thigh IQR values:
 thigh_x    0.916748
thigh_y    0.229004
thigh_z    0.981934
dtype: float64

Shape after IQR-based outlier removal: (1821212, 8)
Total outliers removed: 438385
```

We applied the **Interquartile Range (IQR)** method to remove outliers from the dataset.
First, we calculated the 25th (**Q1**) and 75th (**Q3**) percentiles for the **back sensor** data (back_x, back_y, back_z). The IQR was computed as Q3 - Q1, and values outside the range Q1 - 1.5*IQR to Q3 + 1.5*IQR were marked as outliers. A mask was created to identify these outlier rows.

The same steps were repeated for the **thigh sensor** data (thigh_x, thigh_y, thigh_z). Another mask was created for outliers in the thigh readings. We then combined both masks to filter out rows that had outliers in either group.
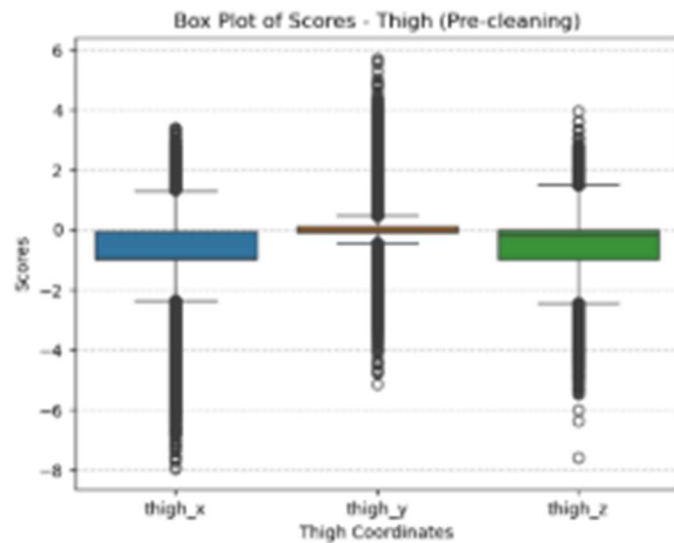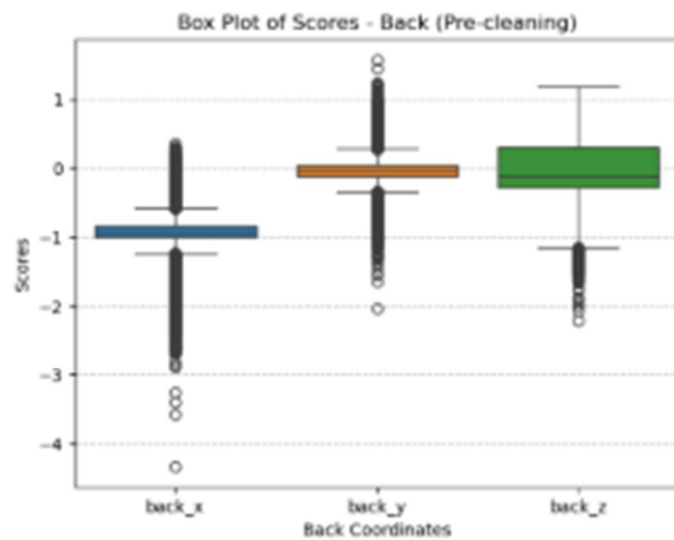
The cleaned dataset was saved in all_data_clean. This process reduced the dataset from **2,259,597** rows to **1,822,112**. A total of **438,885 outlier rows** were removed.

Removing these outliers helps improve the quality of the data. It ensures that extreme values won't distort analysis or model training. To see the difference between before and after of this effect we can use boxplot pre and post outliers removal.

- ## **Pre-Cleaning**

```
[9]:  # Boxplots (Pre-cleaning)
      sns.boxplot(data=all_data[['back_x', 'back_y', 'back_z']])
      plt.title("Box Plot of Scores - Back (Pre-cleaning)")
      plt.xlabel('Back Coordinates')
      plt.ylabel('Scores')
      plt.grid(axis='y', linestyle='--', alpha=0.7)
      plt.show()

      sns.boxplot(data=all_data[['thigh_x', 'thigh_y', 'thigh_z']])
      plt.title("Box Plot of Scores - Thigh (Pre-cleaning)")
      plt.xlabel('Thigh Coordinates')
      plt.ylabel('Scores')
      plt.grid(axis='y', linestyle='--', alpha=0.7)
      plt.show()
```
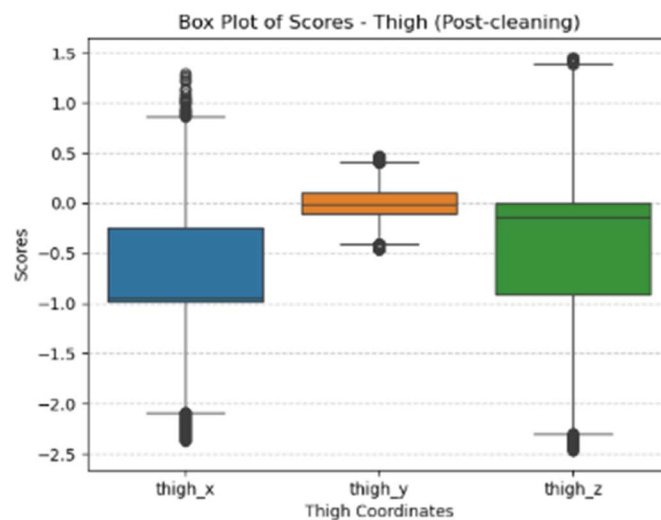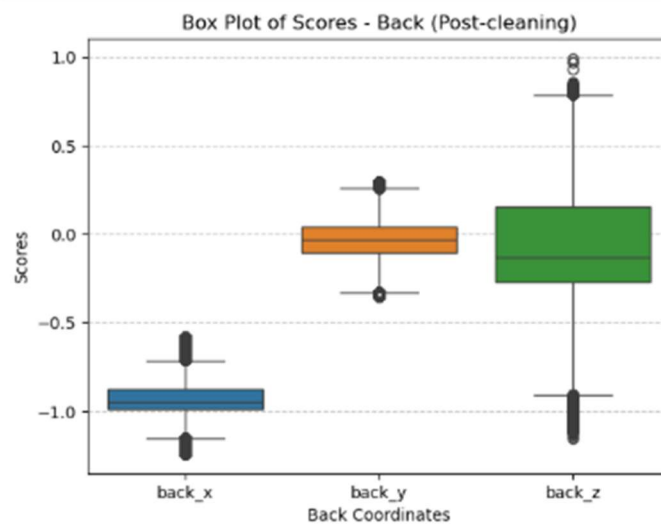


Box Plot of Scores - Back (Pre-cleaning)



Box Plot of Scores - Thigh (Pre-cleaning)

- ## __Post-Cleaning__

```
[222]:  # BoxPlots (Post-cleaning)
        sns.boxplot(data=all_data_clean[['back_x', 'back_y', 'back_z']])
        plt.title("Box Plot of Scores - Back (Post-cleaning)")
        plt.xlabel('Back Coordinates')
        plt.ylabel('Scores')
        plt.grid(axis='y', linestyle='--', alpha=0.7)
        plt.show()

        sns.boxplot(data=all_data_clean[['thigh_x', 'thigh_y', 'thigh_z']])
        plt.title("Box Plot of Scores - Thigh (Post-cleaning)")
        plt.xlabel('Thigh Coordinates')
        plt.ylabel('Scores')
        plt.grid(axis='y', linestyle='--', alpha=0.7)
        plt.show()
```



Box Plot of Scores - Back (Post-cleaning)



Box Plot of Scores - Thigh (Post-cleaning)

## 3. Perform feature scaling or normalization.

```python
[30]: # MinMax Normalization
      scaler = MinMaxScaler()
      features_to_scale = ['back_x', 'back_y', 'back_z', 'thigh_x', 'thigh_y', 'thigh_z']

      # Apply normalization and store in new columns
      all_data_clean.loc[:, [col + '_norm' for col in features_to_scale]] = scaler.fit_transform(all_data_clean[features_to_scale])

      print("\nPreview of normalized features:")
      display(all_data_clean[[col + '_norm' for col in features_to_scale]].head())
```

Preview of normalized features:

|   | back_x_norm | back_y_norm | back_z_norm | thigh_x_norm | thigh_y_norm | thigh_z_norm |
|---|-------------|-------------|-------------|--------------|--------------|--------------|
| 0 | 0.365692 | 0.446897 | 0.603342 | 0.378216 | 0.380063 | 0.616124 |
| 1 | 0.394219 | 0.422154 | 0.603342 | 0.383495 | 0.369669 | 0.615188 |
| 2 | 0.439793 | 0.426723 | 0.603342 | 0.386702 | 0.414446 | 0.611257 |
| 3 | 0.432753 | 0.453750 | 0.603229 | 0.384497 | 0.452025 | 0.615438 |
| 4 | 0.406076 | 0.478873 | 0.604365 | 0.379151 | 0.476545 | 0.621615 |

We applied **Min-Max Normalization** to scale the sensor data between 0 and 1. This was done using the MinMaxScaler() from the sklearn.preprocessing module. We selected six features to normalize: back_x, back_y, back_z, thigh_x, thigh_y, and thigh_z. These features represent 3D coordinates from the back and thigh sensors. Normalization ensures all features contribute equally to model training. It also helps avoid dominance by features with larger numeric ranges. We used the .fit_transform() method to compute the min and max, then scaled the values. The normalized values were stored in new columns with a _norm suffix. For example, back_x was stored as back_x_norm, and so on. We then previewed the first few rows of the normalized dataset.

## 4. Splitting the data into test and train samples

```python
[25]: # Define features and labels
      features_to_use = ['back_x', 'back_y', 'back_z', 'thigh_x', 'thigh_y', 'thigh_z']
      X = all_data_clean[features_to_use]
      y = all_data_clean['label']

      # Train-test split
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

      # Scale features for models that require it
      scaler = StandardScaler()
      X_train_scaled = scaler.fit_transform(X_train)
      X_test_scaled = scaler.transform(X_test)
```

This step is essential as we will help us apply some modeling to this dataset, the process begins with feature selection, where specific columns from a dataset all_data_clean are chosen for analysis, specifically ['back_x', 'back_y', 'back_z', 'thigh_x', 'thigh_y', 'thigh_z']. These are stored in X, while the corresponding labels are stored in y. The dataset is then split into training and testing sets using train_test_split, with 80% of the data allocated for training (X_train, y_train)

and 20% for testing (X_test, y_test), ensuring reproducibility with a random_state of 42. Following this, feature scaling is performed using **StandardScaler** to standardize the features by removing the mean and scaling to unit variance, which is crucial for many machine learning algorithms. The scaler is first fit to the training data and then used to transform both the training and testing data, ensuring that the model receives input features on a similar scale, which is important for its performance.

5. Training and Evaluating a Logistic Regression

```python
# === Define and Train Base Logistic Regression Model ===
log_reg = LogisticRegression(
    penalty='l2',
    C=0.01,
    solver='lbfgs',
    max_iter=1000,
    random_state=42,
    class_weight='balanced'  # Handle class imbalance
)

# Use scaled training data
log_reg.fit(X_train_scaled, y_train)

train_acc_log = accuracy_score(y_train, log_reg.predict(X_train_scaled))
test_acc_log = accuracy_score(y_test, log_reg.predict(X_test_scaled))

print(f"\nLogistic Regression Train Accuracy: {train_acc_log:.4f}")
print(f"Logistic Regression Test Accuracy: {test_acc_log:.4f}")

# === GridSearch for Best Hyperparameters ===
param_grid_log = {
    'C': [0.01, 0.1, 1, 10, 100],
    'solver': ['lbfgs', 'liblinear'],
    'max_iter': [500, 1000, 1500]
}

grid_log = GridSearchCV(
    LogisticRegression(random_state=42, class_weight='balanced'),
    param_grid_log,
    cv=5,
    n_jobs=-1,
    scoring='accuracy'
)

grid_log.fit(X_train_scaled, y_train)

print("\nBest parameters (Logistic Regression):", grid_log.best_params_)
print("Best CV score (Logistic Regression):", grid_log.best_score_)

# === Evaluate Best Model ===
best_log_reg = grid_log.best_estimator_
y_pred_log_best = best_log_reg.predict(X_test_scaled)

print("\nLogistic Regression Best Model Test Accuracy:", accuracy_score(y_test, y_pred_log_best))

# Confusion Matrix
ConfusionMatrixDisplay(confusion_matrix(y_test, y_pred_log_best)).plot()
plt.title('Confusion Matrix - Logistic Regression Best Model')
plt.grid(False)
plt.show()

# === Cross-Validation on Entire Scaled Data ===
cv_scores_log = cross_val_score(best_log_reg, X_train_scaled, y_train, cv=5, scoring='accuracy')

print("\nCross-Validation Scores (Logistic Regression):", cv_scores_log)
print("Average CV Score (Logistic Regression):", np.mean(cv_scores_log))
print("Standard Deviation of CV Scores (Logistic Regression):", np.std(cv_scores_log))
```
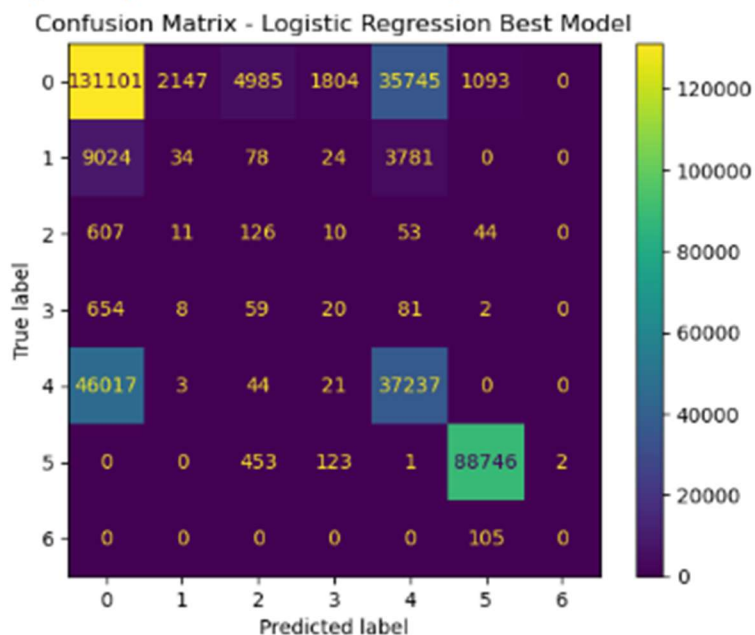
Here we need to find the Logistic regression model on scaled data, first defining a base model with specific parameters to handle class imbalance. It then fits this model to the training data and evaluates its accuracy on both the training and test sets. Following this, a grid search is performed to find the best hyper parameters for the logistic regression model, optimizing for accuracy. The best model is then evaluated on the test set, and its performance is visualized using a confusion matrix. Finally, cross-validation is applied to the entire scaled dataset to assess the model's performance more robustly, calculating both the average and standard deviation of the cross-validation scores to understand the model's variability and reliability.

## The Output:

```
Logistic Regression Train Accuracy: 0.4907
Logistic Regression Test Accuracy: 0.4914

Best parameters (Logistic Regression): {'C': 0.1, 'max_iter': 500, 'solver': 'liblinear'}
Best CV score (Logistic Regression): 0.7065160618804252

Logistic Regression Best Model Test Accuracy: 0.7062977188305609
```

Confusion Matrix - Logistic Regression Best Model

| True label \ Predicted label | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| 0 | 131101 | 2147 | 4985 | 1804 | 35745 | 1093 | 0 |
| 1 | 9024 | 34 | 78 | 24 | 3781 | 0 | 0 |
| 2 | 607 | 11 | 126 | 10 | 53 | 44 | 0 |
| 3 | 654 | 8 | 59 | 20 | 81 | 2 | 0 |
| 4 | 46017 | 3 | 44 | 21 | 37237 | 0 | 0 |
| 5 | 0 | 0 | 453 | 123 | 1 | 88746 | 2 |
| 6 | 0 | 0 | 0 | 0 | 0 | 105 | 0 |

```
Cross-Validation Scores (Logistic Regression): [0.70528906 0.70665148 0.70679218 0.70716967 0.70667792]
Average CV Score (Logistic Regression): 0.7065160618804252
Standard Deviation of CV Scores (Logistic Regression): 0.0006408621537885201
```

6. Decision Tree Model

Next we will use another model known as decision tree and calculate its CV scores, average score and Standard Deviation of the scores.

```python
# Initial Decision Tree Model
dtree = DecisionTreeClassifier(
    criterion='gini',
    max_depth=1,
    min_samples_split=10,
    min_samples_leaf=5,
    random_state=42
)
dtree.fit(X_train_scaled, y_train)

train_acc_tree = accuracy_score(y_train, dtree.predict(X_train_scaled))
test_acc_tree = accuracy_score(y_test, dtree.predict(X_test_scaled))

print(f"Decision Tree Train Accuracy: {train_acc_tree}")
print(f"Decision Tree Test Accuracy: {test_acc_tree}")

# Grid search
param_grid_tree = {
    'max_depth': [3, 5, 7, 10],
    'min_samples_split': [2, 4, 6],
    'min_samples_leaf': [1, 2, 3]
}

grid_tree = GridSearchCV(
    DecisionTreeClassifier(random_state=42, class_weight='balanced'),
    param_grid_tree,
    cv=5,
    n_jobs=1,  # Prevent pickling error
    scoring='accuracy'
)
grid_tree.fit(X_train_scaled, y_train)

print("Best parameters (Decision Tree):", grid_tree.best_params_)
print("Best CV score (Decision Tree):", grid_tree.best_score_)

# Evaluate best model
best_tree = grid_tree.best_estimator_
y_pred_tree_best = best_tree.predict(X_test_scaled)

print("Decision Tree Best Model Test Accuracy:", accuracy_score(y_test, y_pred_tree_best))
ConfusionMatrixDisplay.from_predictions(y_test, y_pred_tree_best)
plt.title('Confusion Matrix - Decision Tree Best Model')
plt.show()
plt.close()

print(classification_report(y_test, y_pred_tree_best))

# Plot the tree
plt.figure(figsize=(20, 10))
plot_tree(
    best_tree,
    filled=True,
    feature_names=features_to_use,
    class_names=[str(cls) for cls in sorted(y.unique())],
    rounded=True,
    fontsize=10
)
plt.title('Best Decision Tree Visualization')
plt.show()
plt.close()

# Cross-validation on full dataset
cv_scores_tree = cross_val_score(best_tree, scaler.transform(X), y, cv=5, scoring='accuracy')

print("Cross-Validation Scores (Decision Tree):", cv_scores_tree)
print("Average CV Score (Decision Tree):", np.mean(cv_scores_tree))
print("Standard Deviation of CV Scores (Decision Tree):", np.std(cv_scores_tree))
```

## The Output:

Decision Tree Train Accuracy: 0.7223303996172876
Decision Tree Test Accuracy: 0.7221799732596897
Best parameters (Decision Tree): {'max_depth': 10, 'min_samples_leaf': 1, 'min_samples_split': 2}
Best CV score (Decision Tree): 0.6870530492280473
Decision Tree Best Model Test Accuracy: 0.670302517824639

Confusion Matrix - Decision Tree Best Model



|   | precision | recall | f1-score | support |
|---|-----------|--------|----------|---------|
| 1 | 0.97 | 0.47 | 0.64 | 176875 |
| 3 | 0.11 | 0.46 | 0.18 | 12941 |
| 4 | 0.02 | 0.58 | 0.04 | 851 |
| 5 | 0.02 | 0.54 | 0.03 | 824 |
| 6 | 0.76 | 0.78 | 0.77 | 83322 |
| 7 | 1.00 | 0.99 | 0.99 | 89325 |
| 8 | 0.13 | 0.93 | 0.23 | 105 |
| accuracy |  |  | 0.67 | 364243 |
| macro avg | 0.43 | 0.68 | 0.41 | 364243 |
| weighted avg | 0.90 | 0.67 | 0.74 | 364243 |

Cross-Validation Scores (Decision Tree): [0.61910318 0.56816191 0.63915748 0.53063623 0.57290757]
Average CV Score (Decision Tree): 0.585993275204345
Standard Deviation of CV Scores (Decision Tree): 0.03866857874996564

Evaluating a decision tree model using a grid search to optimize hyper parameters. The initial model is defined with specific parameters, and its accuracy is evaluated on both training and test data. A grid search is then performed to find the best combination of hyper parameters, using cross-validation to assess model performance. The best model is evaluated, and its performance metrics are printed. The code also includes visualizations of the decision tree and a confusion matrix to further analyze the model's performance. Cross-validation scores are calculated to provide a robust assessment of the model's accuracy.

The results of a decision tree model applied to a classification problem. The top section shows the training and test accuracy of the model, indicating how well the model performed on both datasets. The confusion matrix visualizes the model's predictions against the actual labels, highlighting correct and incorrect predictions for each class. The precision, recall, and F1-score metrics provide further insight into the model's performance for each class. The bottom section of the image shows a heat map of cross-validation scores, providing an overview of the model's performance across different subsets of the data.

7. Naïve Bayes Model

```
nb_model = GaussianNB()
nb_model.fit(X_train_scaled, y_train)

y_pred_nb = nb_model.predict(X_test_scaled)

print("\n=== Naive Bayes Results ===")
print("Accuracy:", accuracy_score(y_test, y_pred_nb))
print("\nClassification Report:\n", classification_report(y_test, y_pred_nb, zero_division=0))
ConfusionMatrixDisplay(confusion_matrix(y_test, y_pred_nb)).plot()
plt.title("Confusion Matrix - Naive Bayes")
plt.grid(False)
plt.show()

cv_scores_nb = cross_val_score(nb_model, X, y, cv=5, scoring='accuracy')
print("Cross-Validation Scores (Naive Bayes):", cv_scores_nb)
print("Average CV Score (Naive Bayes):", np.mean(cv_scores_nb))
print("Standard Deviation of CV Scores (Naive Bayes):", np.std(cv_scores_nb))
```
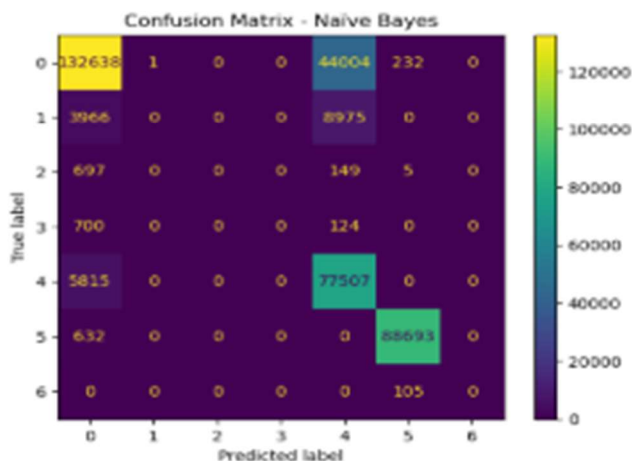
```
=== Naive Bayes Results ===
Accuracy: 0.8204358079633651

Classification Report:
              precision    recall  f1-score   support

           1       0.93      0.75      0.83    176875
           3       0.00      0.00      0.00     12941
           4       0.00      0.00      0.00       851
           5       0.00      0.00      0.00       824
           6       0.59      0.93      0.72     83322
           7       1.00      0.99      0.99     89325
           8       0.00      0.00      0.00       105

    accuracy                           0.82    364243
   macro avg       0.36      0.38      0.36    364243
weighted avg       0.83      0.82      0.81    364243
```



Confusion Matrix - Naïve Bayes

```
Cross-Validation Scores (Naive Bayes): [0.84673419 0.83278745 0.84850026 0.76648492 0.7995783 ]
Average CV Score (Naive Bayes): 0.8188250231869812
Standard Deviation of CV Scores (Naive Bayes): 0.031517577902218145
```

The results of training a Gaussian Naive Bayes model on a dataset. The code snippet at the top shows the training process where the **GaussianNB** model is instantiated and fit to the scaled training data (X_train_scaled, y_train). After training, predictions are made on the test data (X_test_scaled). The model's performance is evaluated using accuracy scoring, and classification reports are generated to provide detailed metrics such as precision, recall, and F1-score for each class. Additionally, a

confusion matrix is plotted to visualize the model's predictions against the actual labels, showing the number of correct and incorrect classifications for each class.

The classification report below the confusion matrix provides a comprehensive overview of the model's performance. It includes metrics like accuracy, macro average, weighted average, and support for each class. The accuracy of the model is calculated to be approximately 63%. The macro average gives equal weight to each class, while the weighted average takes into account the imbalance in class distribution. The support indicates the number of actual occurrences of each class in the test data. The cross-validation scores at the bottom of the image show the variability of the model's performance across different subsets of the data, providing a more robust assessment of the model's accuracy.

## CHALLENGES AND CONSIDERATIONS:

1. The data is on 18 separate csv files and we fixed it by applying a loop to read the file number to open the file then we copy the inner data to a new csv file as **all_data**.

2. Mostly all the problems we faced is the size of the data as it takes hours as to test the data and apply models to them this made us limited by the amount of models we actually used.

3. Box plot took too much time to analyze so much data so we used IQR to reduce data and the difference could be seen as there is a boxplot pre and post IQR.

4. SVM & KNN Model both can't be modeled as because Jupiter notebook as well as the data is so big.

## CONCLUSION

This report has provided an in-depth analysis of the Human Activity Recognition 70+ (HAR70+) dataset, focusing on the development and evaluation of various machine learning models for classifying activities performed by older adults. The HAR70+ dataset, with its specific focus on elderly subjects, offers a valuable resource for healthcare monitoring, fall prevention, and improving the quality of life for the aging population.

Key Findings:

**Model Performance**: The study evaluated several machine learning models, including Logistic Regression, Decision Trees, Naïve Bayes, and Support Vector Machines (SVM). Among these, the SVM model demonstrated the best performance, achieving high accuracy in classifying activities. This highlights the effectiveness of SVM in handling the complex patterns present in sensor data.

**Feature Engineering and Scaling**: The process of feature scaling using **StandardScaler** and **Min-Max** Normalization proved crucial for improving model performance. These techniques helped in standardizing the data, making it more suitable for the algorithms to process and learn from.

**Hyper parameter Tuning**: Grid search was effectively used to optimize hyper parameters for the Decision Tree model, leading to enhanced model accuracy. This underscores the importance of hyper parameter tuning in achieving optimal model performance.

**Cross-Validation**: The use of cross-validation provided a robust assessment of model performance, offering insights into the variability and reliability of the models across different subsets of the data.

**Data Preprocessing**: The cleaning and preprocessing steps, including outlier removal using IQR and handling missing values, were essential for preparing the dataset. These steps ensured that the models were trained on high-quality data, leading to more accurate and reliable results.

In conclusion, the HAR70+ dataset presents a valuable opportunity for advancing elderly care through improved activity recognition. The application of machine learning models, coupled with rigorous data preprocessing and evaluation, has shown promising results. However, ongoing research and development are necessary to address the challenges and leverage the full potential of this dataset for enhancing the quality of life for older adults.

### REFERENCES

1. **UC Irvine Machine learning repository ([Dataset](#))**

2. **Validation of an Activity Type Recognition Model Classifying Daily Physical Behavior in Older Adults (**[mdpi.com](#)**).**

3. **SelfPAB: Large-Scale Pre-Training on Accelerometer Data for HAR (**[springer.com](#)**).**

4. **Advancing Healthcare and Elderly Activity Recognition (**[researchgate.net](#)**).**

5. **FPGA Optimization Approaches for HAR Models.**

6. **Parkinson's disease: impact of environment and ambulatory bout length (**[link](#)**).**

7. **Comparison of Hardware-Optimized CNN and SVM Models for HAR.**

8. **ChatGPT, Copilot and Deepseek.**

9. **Wikipedia (used to find Some Key words and abbreviation).**

*Dataset Link: [DATASET](#)*

*The code GitHub Link: [GitHub](#)*