# Report on UCI HAR Dataset Analysis(Task #1)

By : Siddharth Brijesh Tripathi
Intern ID: 21800761

## 1. Overview

This notebook focuses on building models to predict human activities using the UCI HAR dataset. The task involves two key approaches:

- **Deep Learning Models:** Implemented using LSTM and 1D CNN trained on raw accelerometer data.
- **Machine Learning Models:** Implemented using Random Forest, SVM, and Logistic Regression with features generated using the TSFEL library.
- A comparison is performed between models trained on the extracted features and those trained on the existing features provided by the dataset authors.

---

## 2. Implementation Summary

### 2.1 Dataset Handling

- The dataset is extracted and preprocessed.
- The raw accelerometer data is used for Deep Learning models.
- Feature extraction is performed using the TSFEL library for Machine Learning models.
- Existing features from the dataset are also used for performance comparison.

### 2.2 Deep Learning Models

- Implemented LSTM and 1D CNN architectures using `tensorflow.keras`.
- Evaluated models using accuracy and confusion matrices.

### 2.3 Machine Learning Models

- Features extracted using TSFEL (time-domain, frequency-domain, and statistical features).
- Implemented models: **Random Forest, SVM, Logistic Regression**.
- Evaluated using precision, recall, F1-score, and accuracy.
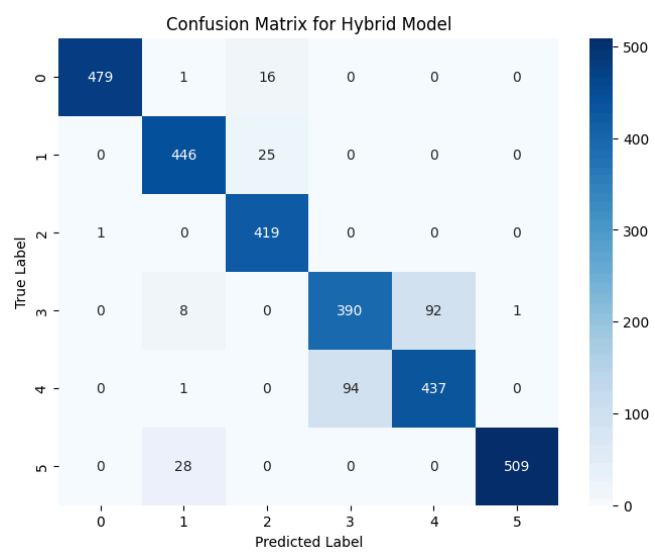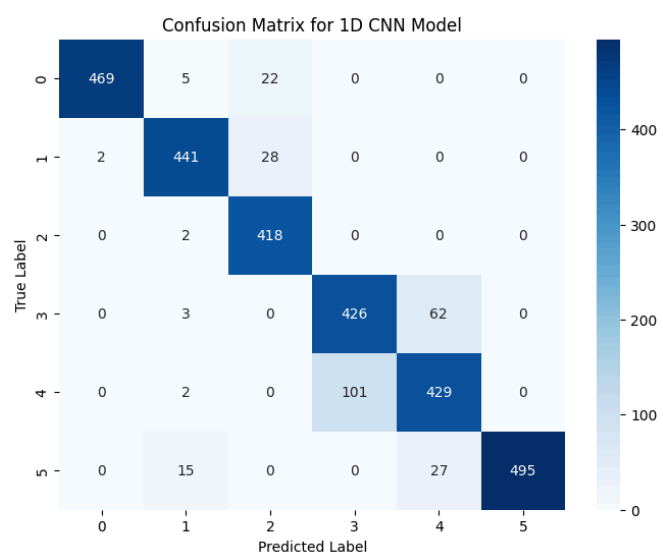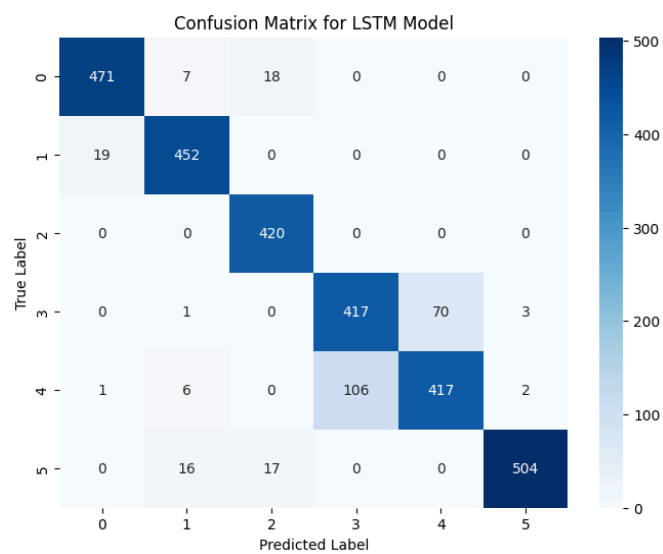- Models are trained on both extracted and existing features.

# 3. Results & Observations

## 3.1 Model Performance Comparison

- **Deep Learning Models** (LSTM, CNN):
  - Performed well with raw sensor data.
  - The 3 layered LSTM achieved accuracy of 90.97% trained and tested on the raw features.
  - The 3 layered CNN achieved accuracy of 90.87% trained and tested on the raw features.
  - The combined model (LSTM+CNN) achieved accuracy of 90.93%.


- **Machine Learning Models** (Random Forest, SVM, Logistic Regression):
  - The ML models performed somewhat worse and achieved less accuracy than the deep learning models.
  - Among the ML models, the SVM was the best achieving accuracy of ~90%.


- **Comparison of Pre-Existing and TSFEL Extracted features:**
  - The ML models achieved better accuracy when trained and tested on the features provided by the Authors(accuracy ~ 92%) performed better than on raw features(accuracy ~ 90%) and the extracted features(accuracy ~70%).
  - I also tried the selection of most important features by removing some of the features but that also did not improve the accuracy much.


## 3.2 Visualizations & Evaluation

- Confusion matrices were generated for model evaluation.
- Performance metrics like accuracy, precision, recall, and F1-score were analyzed.
- A summary table comparing Deep Learning and Machine Learning model performance was included.

## Confusion Matrix for LSTM Model

| True \ Predicted | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 471 | 7 | 18 | 0 | 0 | 0 |
| 1 | 19 | 452 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 420 | 0 | 0 | 0 |
| 3 | 0 | 1 | 0 | 417 | 70 | 3 |
| 4 | 1 | 6 | 0 | 106 | 417 | 2 |
| 5 | 0 | 16 | 17 | 0 | 0 | 504 |

## Confusion Matrix for 1D CNN Model

| True \ Predicted | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 469 | 5 | 22 | 0 | 0 | 0 |
| 1 | 2 | 441 | 28 | 0 | 0 | 0 |
| 2 | 0 | 2 | 418 | 0 | 0 | 0 |
| 3 | 0 | 3 | 0 | 426 | 62 | 0 |
| 4 | 0 | 2 | 0 | 101 | 429 | 0 |
| 5 | 0 | 15 | 0 | 0 | 27 | 495 |

## Confusion Matrix for Hybrid Model

| True \ Predicted | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 479 | 1 | 16 | 0 | 0 | 0 |
| 1 | 0 | 446 | 25 | 0 | 0 | 0 |
| 2 | 1 | 0 | 419 | 0 | 0 | 0 |
| 3 | 0 | 8 | 0 | 390 | 92 | 1 |
| 4 | 0 | 1 | 0 | 94 | 437 | 0 |
| 5 | 0 | 28 | 0 | 0 | 0 | 509 |

# 4. Areas for Improvement

### 4.1 Feature Engineering

- A deeper analysis of which TSFEL features contribute most to classification performance.
- Consider feature selection techniques to remove redundant features.

### 4.2 Hyperparameter Tuning

- Fine-tune hyperparameters for SVM, Random Forest, and Logistic Regression using GridSearchCV.
- Experiment with different architectures for LSTM and CNN models.

### 4.3 Additional Metrics

- Evaluate latency and computational efficiency of models.
- Include ROC curves for better class separation analysis.

# 5. Conclusion

- Deep Learning models (LSTM, CNN) effectively learn activity patterns from raw sensor data.
- Training on extracted features **did not provide** as good performance as using pre-existing dataset features.
- Further improvements can be achieved through feature selection, hyperparameter tuning, and additional evaluation metrics.

# 6. Future Work

- Implement hybrid models combining Deep Learning with feature-engineered Machine Learning.
- Explore Transformer-based models for time-series classification.
- Conduct real-time testing on embedded devices or mobile platforms.

**End of Report**