

WESAD database

Stress Analysis Using ECG, EDA,
and EMG Signals

Team HARmonizers
Kanishk Goel- 2021325
Milind Jain- 2021165



01. Problem Statement

02. Progress

03. Challenges with Raw Data

*04. Hypothesis Tests and Experiments
for Validation*

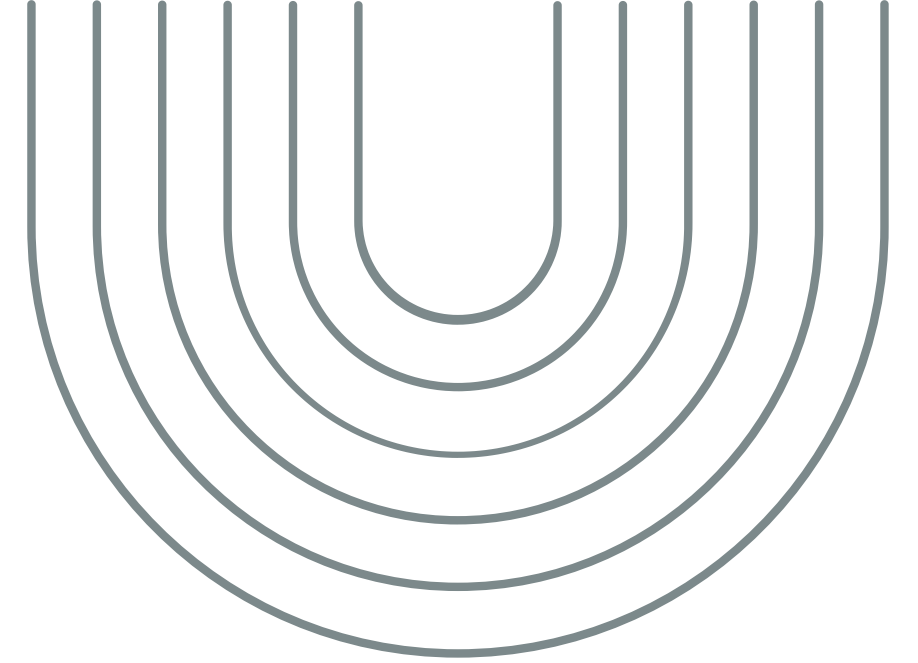
05. Conclusion & Plan of action



TABLE OF CONTENT

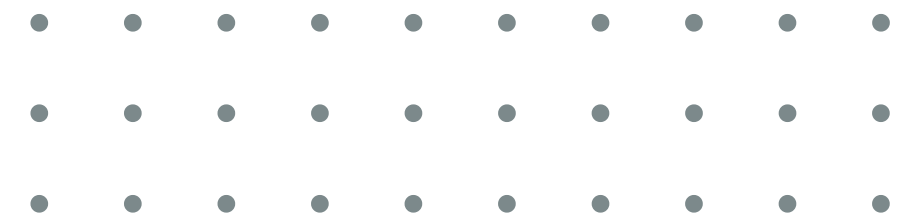


PROBLEM STATEMENT



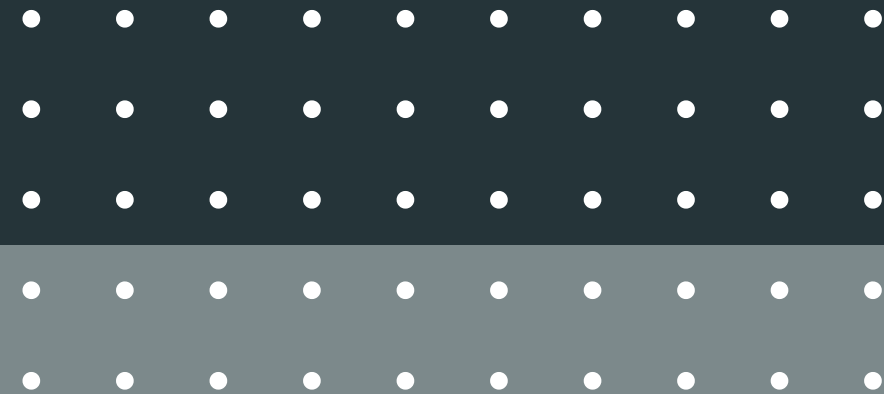
- "Understanding the physiological response to stress using non-invasive biomedical signals."
- Importance: "Timely detection and quantification of stress can lead to effective interventions, improving mental and physical well-being."

01.



02.

DATA PREPROCESSING





1. Data Cleaning

- Removal of irrelevant data points.
- Handling and imputation of missing values.
- Rectification of inconsistencies in the dataset.

2. Normalization (Handling Noisy Data)

- Applied Min-Max scaling to transform the data into the range [0, 1].
- Reduced the effect of outliers and ensured uniformity in the data distribution.

```
# Drop NaN values
ecg_raw = ecg_raw[~np.isnan(ecg_raw)]
eda_raw = eda_raw[~np.isnan(eda_raw)]
emg_raw = emg_raw[~np.isnan(emg_raw)]
```

```
def normalize_data(data):
    """Normalize data using Min-Max scaling."""
    return (data - np.min(data)) / (np.max(data) - np.min(data))
```

3. Feature Selection

- Identified and retained the most informative features: ECG, EDA, and EMG.
- Excluded other features to focus on primary physiological signals related to stress.

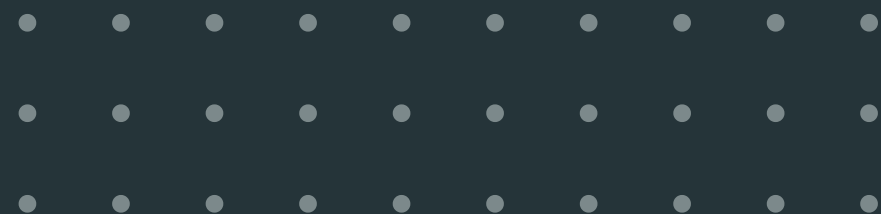
4. Data Splitting

- Divided the dataset into training and evaluation subsets.
- Ensured the model has unseen data for validation and performance assessment.

```
# Load data from .pkl file
pkl_path = os.path.join(subject, subject + '.pkl')
with open(pkl_path, 'rb') as file:
    data = pickle.load(file, encoding='latin1')

# Extracting chest data
chest_data = data['signal']['chest']

ecg_raw = chest_data['ECG']
eda_raw = chest_data['EDA']
emg_raw = chest_data['EMG']
```



03.

CHALLENGES



PROBLEM

Complex Structure of WESAD Dataset.

Noisy Data and Signal Interference.

SOLUTION

Reviewed dataset documentation and extracted relevant signals systematically.

Implemented signal processing techniques, including filtering, to reduce noise and enhance the clarity of physiological signals.

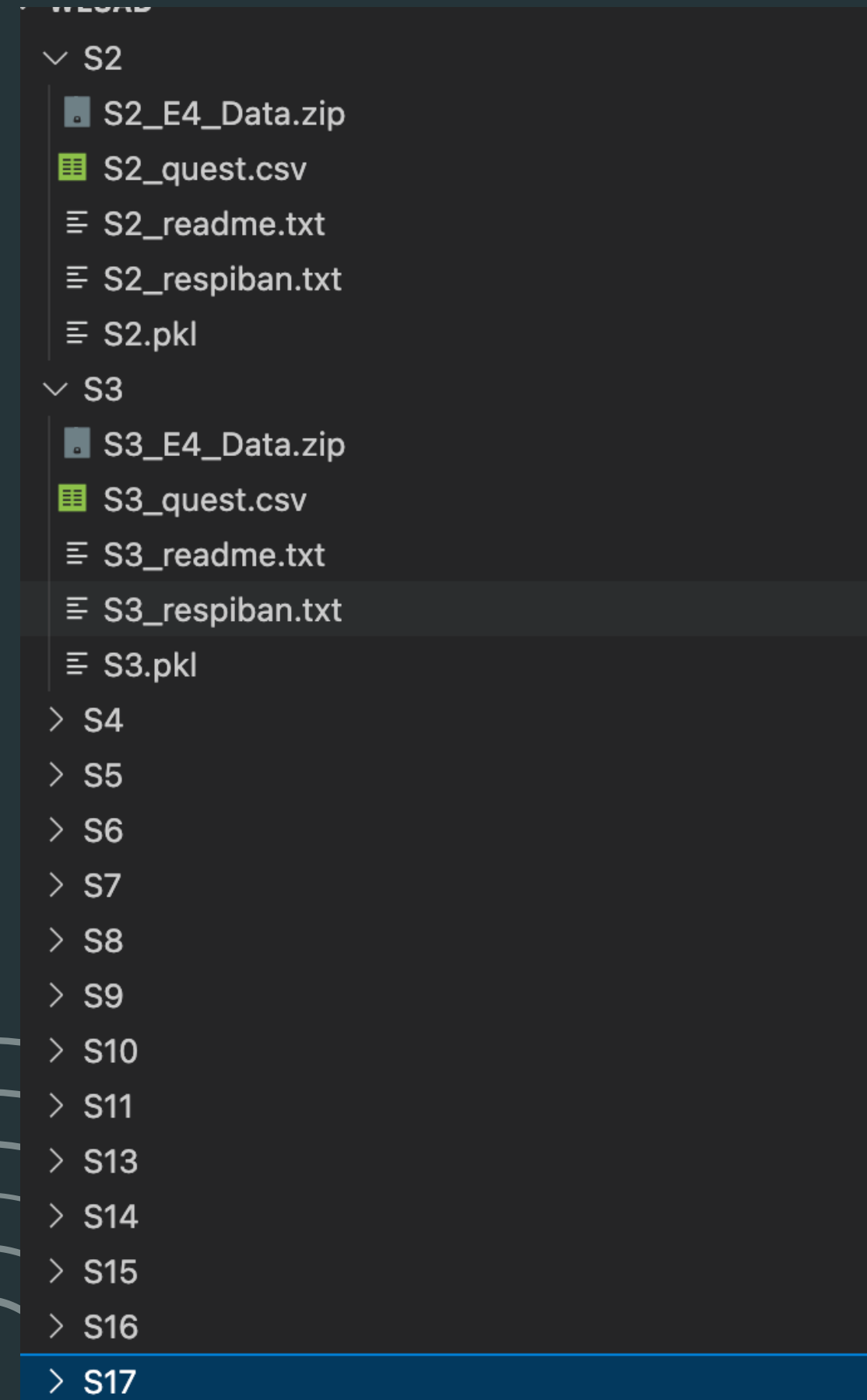


PROBLEM

Combining Multiple Data Streams

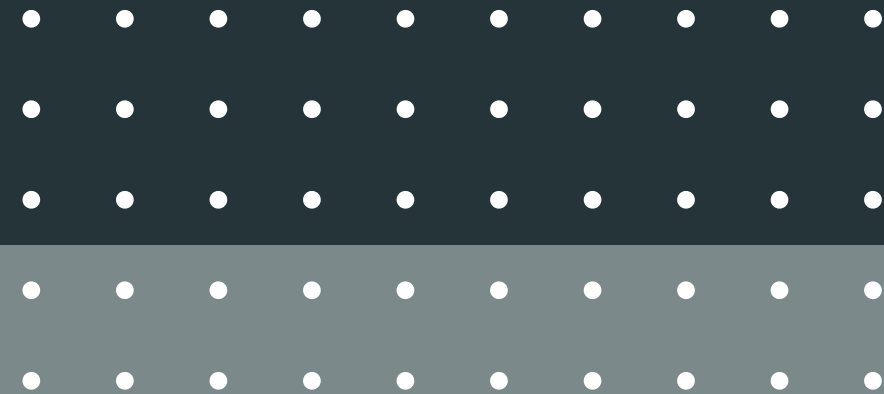
SOLUTION

Developed a systematic pipeline for data concatenation and alignment, preserving the temporal integrity of the signals.



04.

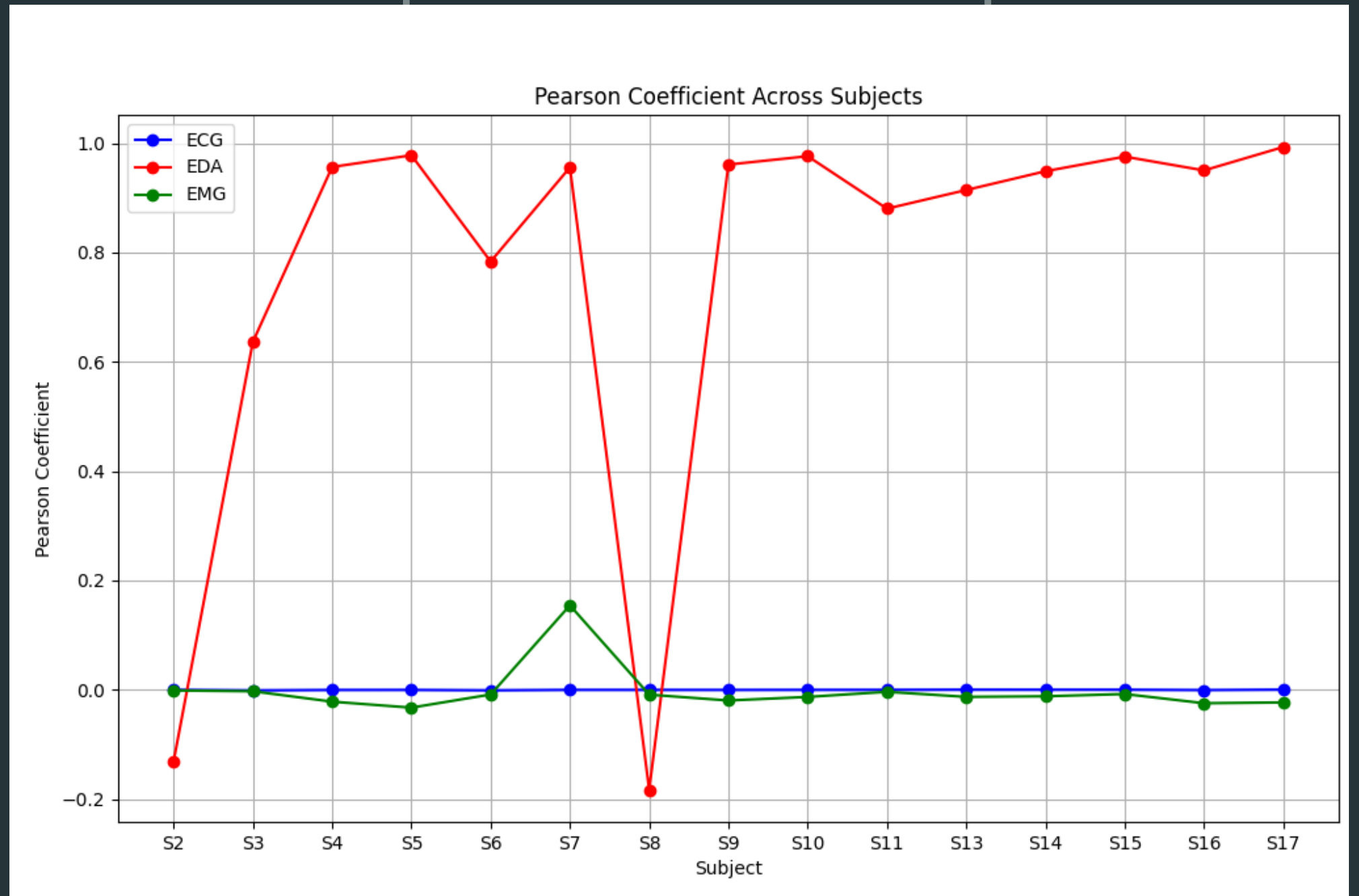
HYPOTHESIS TESTING



PEARSON COEFFICIENT

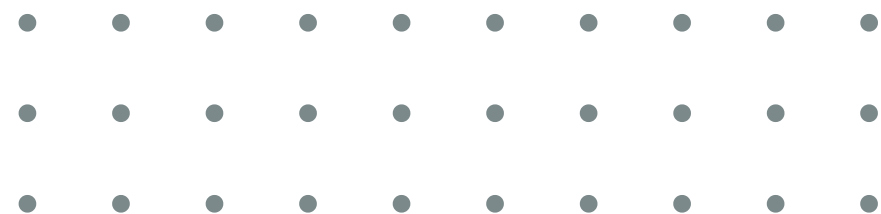
We decided to first check the pearson correlation coefficient between baseline and the stress level for ECG, EDA and EMG. Here is the plot for the same for all the subjects.

Here, we observed that the EDA and stress levels were highly correlated(positive linear relationship)



HYPOTHESIS

- "Is there a significant difference in the means of biomedical signals (ECG, EDA, EMG) between baseline and stress conditions?"
- "Ho (Null Hypothesis): There isn't a significant difference between the means."
- "H1 (Alternative Hypothesis): There is a significant difference between the means."



PROCEDURE

Data Segregation:

- Segregated the data based on labels, distinguishing between baseline and stress conditions.

Sampling:

- Randomly sampled 10% of the data for each physiological signal under both conditions.

Z-test Application:

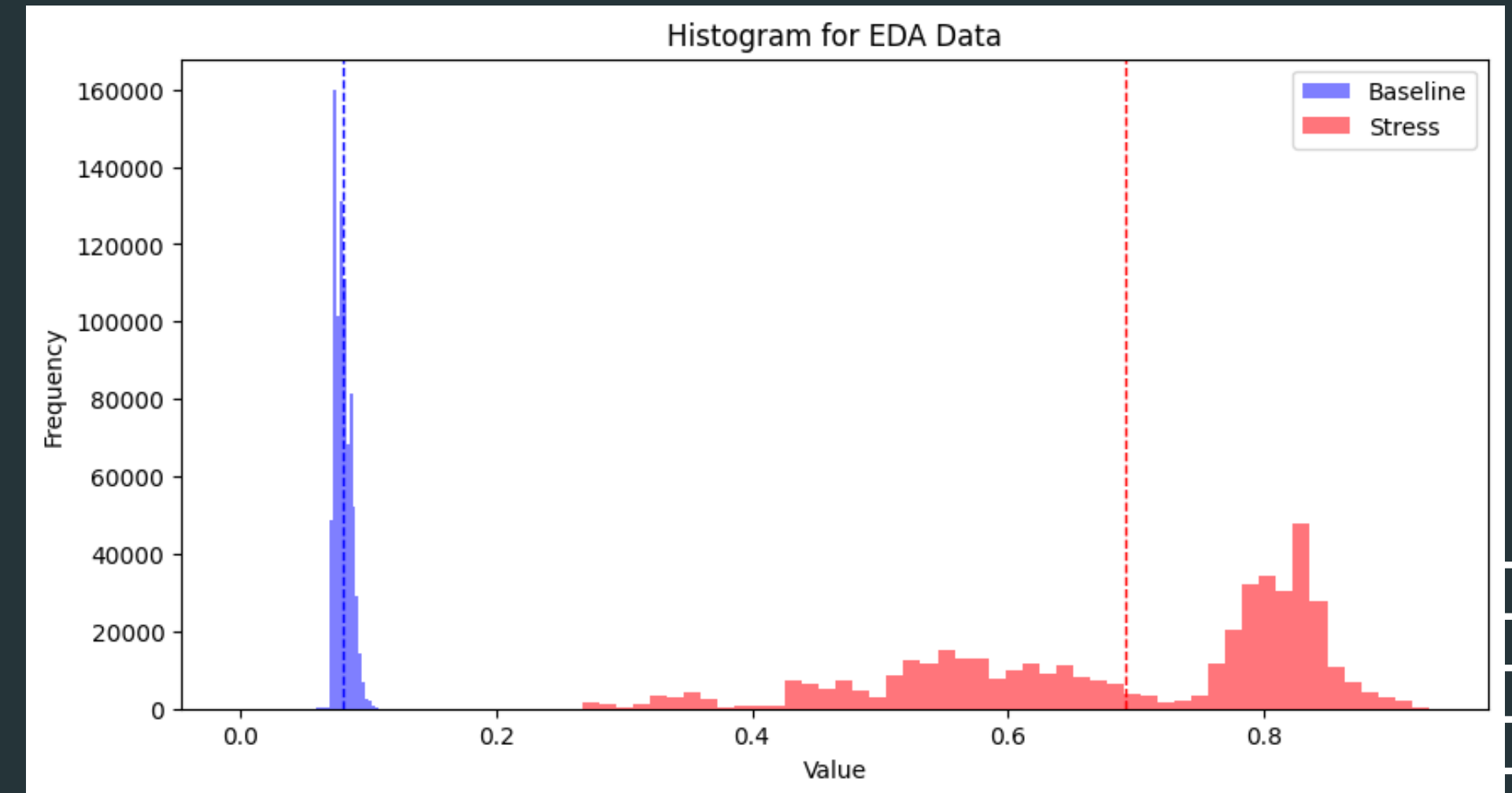
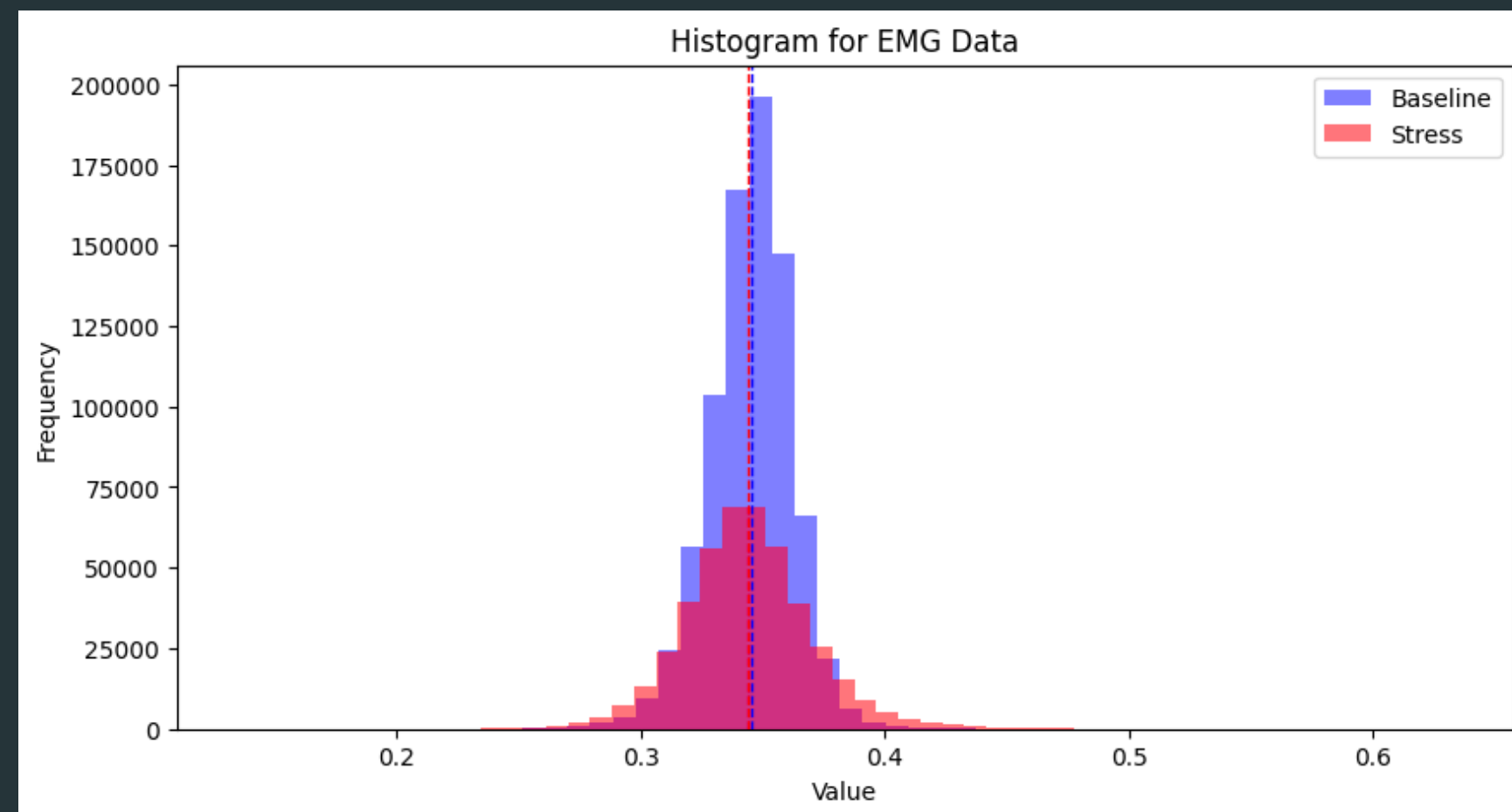
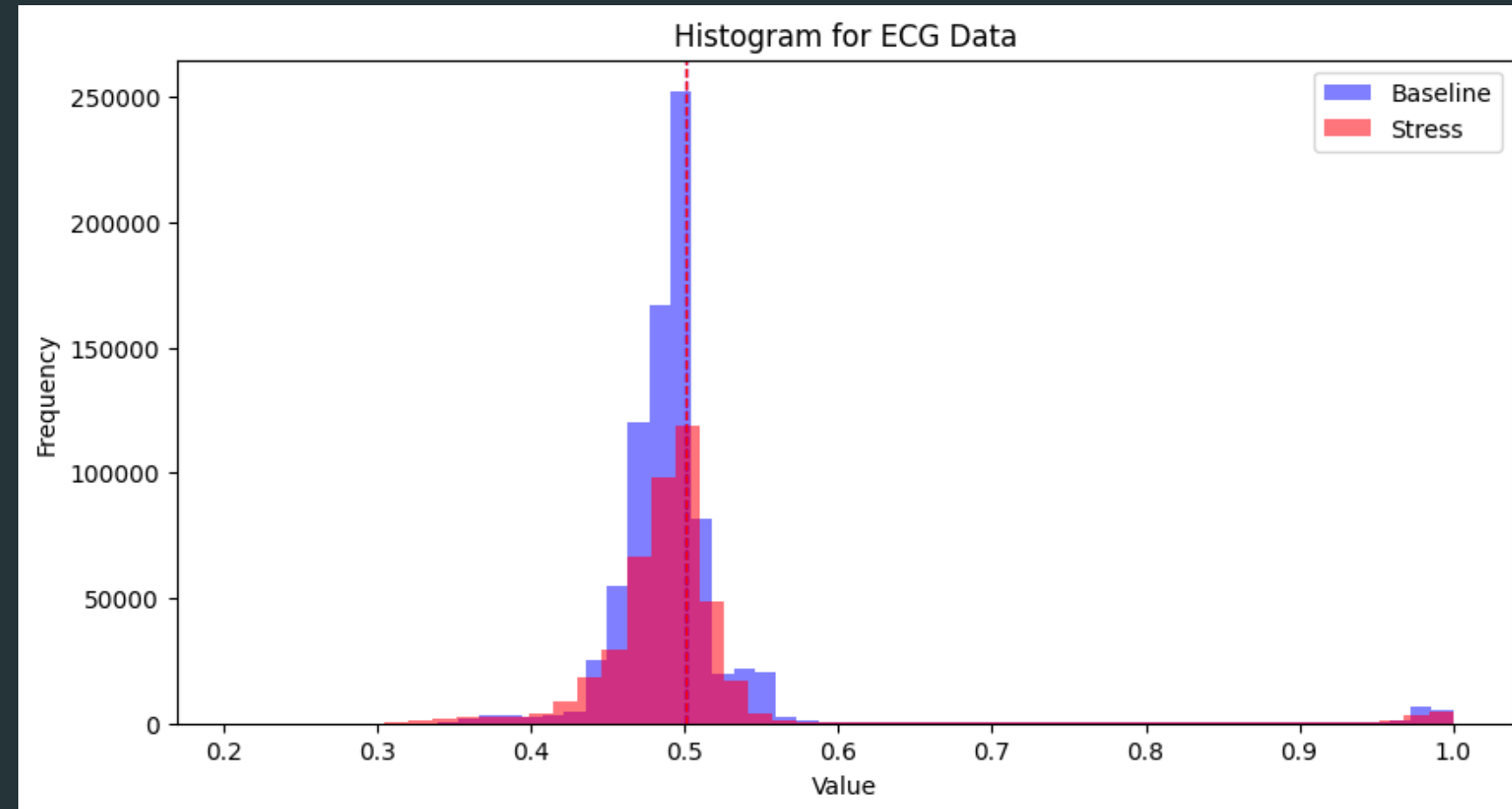
- Applied the Z-test to compare the means of the sampled data for baseline and stress conditions.

Decision Criteria:

- Used a significance level (α) of 0.05.
- If $p\text{-value} < \alpha$, reject the null hypothesis.
- If $p\text{-value} > \alpha$, fail to reject the null hypothesis.



VISUALISATION





HYPOTHESIS RESULTS AND VALIDATION

For nearly all subjects, a significant difference in means was observed between baseline and stress conditions, both in sampled and entire dataset. This aligns with the strong correlation seen between EDA and stress.

For ECG and EMG, no significant correlation was observed.

```
Subject: S17
ECG-Baseline vs Stress Correlation: 0.0002656950030589104
EDA-Baseline vs Stress Correlation: 0.9930300191900447
EMG-Baseline vs Stress Correlation: -0.02310494688339912

Validation for ECG data:

Hypothesis Test on 10.0% of ECG data:
Z-score: 0.4353304115961519, p-value: 0.663322613490565
Fail to reject the null hypothesis for ECG on the sample: There isn't a significant difference between the means.
Validation on entire dataset contradicts this, as means of ECG Baseline and Stress are different.

Validation for EDA data:

Hypothesis Test on 10.0% of EDA data:
Z-score: -3490.9355480744925, p-value: 0.0
Reject the null hypothesis for EDA on the sample: There is a significant difference between the means.
Validation on entire dataset supports this, as means of EDA Baseline and Stress are different.

Validation for EMG data:

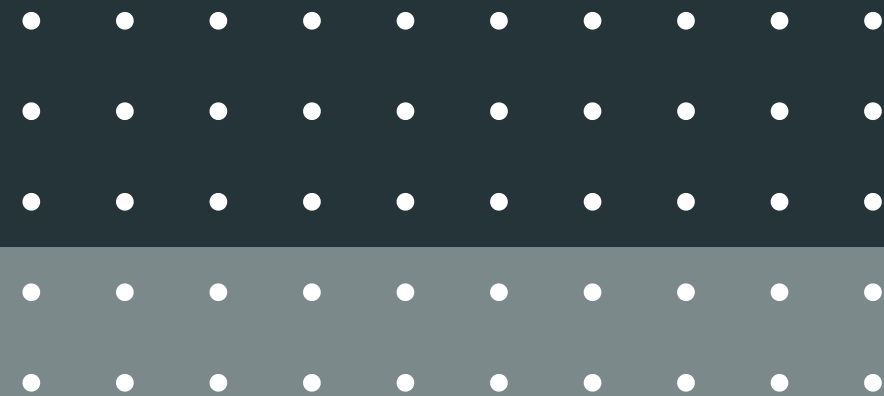
Hypothesis Test on 10.0% of EMG data:
Z-score: 11.050213504346823, p-value: 0.0
Reject the null hypothesis for EMG on the sample: There is a significant difference between the means.
Validation on entire dataset supports this, as means of EMG Baseline and Stress are different.
=====
```

For EDA, findings remained consistent for nearly all subjects, validating the results.

05.

CONCLUSION

- Successfully pre-processed and analyzed the complex WESAD dataset.
- Identified significant correlations and patterns in ECG, EDA, and EMG signals related to stress.
- We found a strong positive correlation between stress levels and EDA signals.
- Validated hypotheses revealing differences in baseline and stress states across subjects.



01. Recalling Problem Statement

02. Data preprocessing

03. ML Models

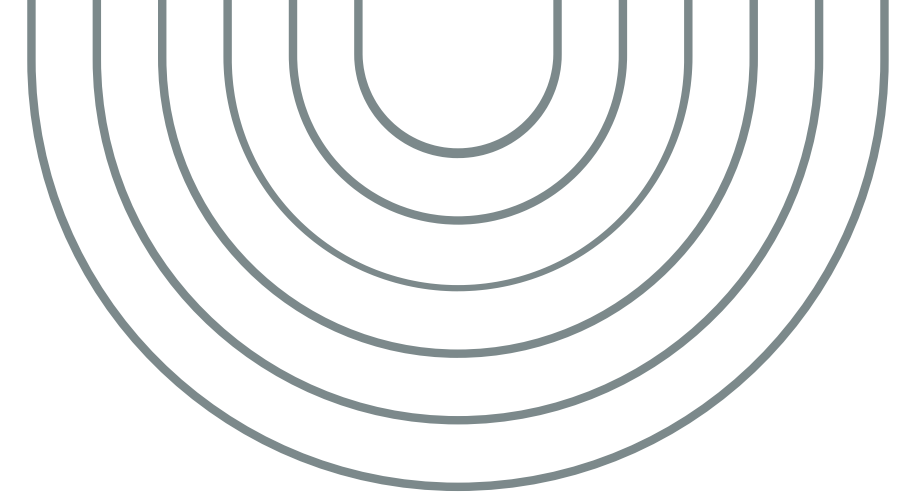
04. Random Forest, SVM and Logistic Regression

05. Comparison of the different Ml models



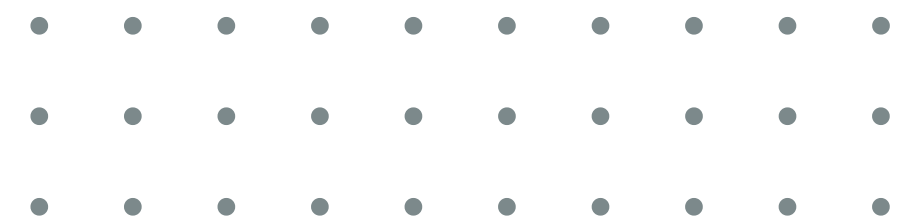
**MOVING
FORWARD**

PROBLEM STATEMENT



This project aims to leverage the WESAD (WEarable Stress and Affect Detection) dataset to develop personalized machine learning models that accurately predict stress levels for individuals. By analyzing physiological and demographic data collected through wearable sensors, the project seeks to understand the variance in stress responses among different subjects. The goal is to create a robust model for each individual that effectively distinguishes between stress and non-stress states, potentially contributing to personalized stress management solutions.

01.



FURTHER DATA PREPROCESSING

WESAD

> __pycache__

> cvxEDA

> src

> __pycache__

cvxEDA.m

cvxEDA.py

> test

.gitignore

LICENSE.txt

README.txt

> data

> subject_feats

S2_feats_4.csv

S3_feats_4.csv

S4_feats_4.csv

S5_feats_4.csv

S6_feats_4.csv

S7_feats_4.csv

S8_feats_4.csv

S9_feats_4.csv

S10_feats_4.csv

S11_feats_4.csv

S13_feats_4.csv

S14_feats_4.csv

data > subject_feats > S2_feats_4.csv

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

,net_acc_mean,net_acc_std,net_acc_min,net_acc_max,EDA_phasic_mean,EDA_phasic_std,EDA_phasic_min,EDA_phasic_max,EDA_smna_mean,EDA_smna_std,EDA_smna_min,

0,0.609521853499563,0.1414808818605863,-358.13,554.77,0.609521853499563,1.0883282078827645,-358.13,554.77,0.609521853499563,1.952822653105783,-358.13,5

0,0.5380042561332239,0.09188242918526915,-392.28,438.16,0.5380042561332239,1.2235284466452476,-392.28,438.16,0.5380042561332239,2.8525099643883873,-392

0,0.5747834522945172,0.10231487191145877,-240.61,209.89,0.5747834522945172,0.12889555567315103,-240.61,209.89,0.5747834522945172,0.2448912309527342,-240

0,0.5908793286276521,0.046391093143305494,-289.26,145.36,0.5908793286276521,0.1261666514224916,-289.26,145.36,0.5908793286276521,0.1997043206547065,-289

0,0.567452195664804,0.03453958636272787,-197.37,194.12,0.567452195664804,0.039616218775094955,-197.37,194.12,0.567452195664804,0.11473779195936254,-197

0,0.5765672803768963,0.030052817368396118,-367.11,363.29,0.5765672803768963,0.05187081817631527,-367.11,363.29,0.5765672803768963,0.11219155530599768,-3

0,0.5905417430311125,0.018591445079357678,-670.2,361.86,0.5905417430311125,0.14789909342142904,-670.2,361.86,0.5905417430311125,0.5932838764656052,-670

0,0.5349033251929577,0.017300961607875073,-182.86,182.92,0.5349033251929577,0.05975018119363328,-182.86,182.92,0.5349033251929577,0.1761955558235208,-18

0,0.544739927671242,0.017555040059839478,-120.43,139.13,0.544739927671242,0.012294712157708068,-120.43,139.13,0.544739927671242,0.004402871272127186,-12

0,0.592366136472579,0.017095621839545518,-39.38,77.19,0.592366136472579,0.0044412709461499,-39.38,77.19,0.592366136472579,0.013326580402636413,-39.38,7

0,0.5159207379131975,0.008785377737088344,-596.97,368.75,0.5159207379131975,0.02688954065275901,-596.97,368.75,0.5159207379131975,0.07936218958014021,-5

0,0.5874940949071376,0.013722298952503632,-449.14,338.57,0.5874940949071376,0.029751614611405863,-449.14,338.57,0.5874940949071376,0.08161297986138642,-

0,0.5810779110655182,0.01814573385506129,-400.42,421.48,0.5810779110655182,0.025146486869367404,-400.42,421.48,0.5810779110655182,0.1032086667105597,-40

0,0.5605698092814017,0.015354349252197905,-428.17,285.6,0.5605698092814017,0.03198285278631688,-428.17,285.6,0.5605698092814017,0.0978703816414689,-428

0,0.6175694440142145,0.009597120600196324,-727.64,595.6,0.6175694440142145,0.018005964807266524,-727.64,595.6,0.6175694440142145,0.07870259113750801,-7

0,0.4890710902922207,0.009703592911034833,-335.37,351.65,0.4890710902922207,0.011270269819860911,-335.37,351.65,0.4890710902922207,0.002802863769478963,

0,0.5656286537446896,0.00558201019101538,-34.68,37.36,0.5656286537446896,0.0016749033033470254,-34.68,37.36,0.5656286537446896,0.004327540891224482,-34

0,0.5287147596008881,0.006217793122105529,-409.52,347.37,0.5287147596008881,0.016144276226611173,-409.52,347.37,0.5287147596008881,0.06672805026721546,-

0,0.5896058142568802,0.00796747336025052,-351.18,402.66,0.5896058142568802,0.017234819310553227,-351.18,402.66,0.5896058142568802,0.016694616423042454,-

0,0.5724750553569308,0.003838960143066366,-116.17,380.05,0.5724750553569308,0.009883142057826824,-116.17,380.05,0.5724750553569308,0.036469003966793416,

0,0.5636727604337,0.002107620011637863,-38.32,35.95,0.5636727604337,0.0045730362426433455,-38.32,35.95,0.5636727604337,0.0044848253082944885,-38.32,35.

0,0.5539298554674729,0.006210066746916988,-314.16,225.36,0.5539298554674729,0.008901259094837796,-314.16,225.36,0.5539298554674729,0.025735799001093854,

0,0.5943513237177004,0.0030779970908238084,-186.68,214.5,0.5943513237177004,0.0038835093667029,-186.68,214.5,0.5943513237177004,0.014021104316460283,-

0,0.5331439040296682,0.002611415603280077,-126.98,81.53,0.5331439040296682,0.0026766498944738584,-126.98,81.53,0.5331439040296682,0.006821965333678672,-

0,0.5640153137842114,0.0026507423521210916,-182.35,204.13,0.5640153137842114,0.00268792967491335,-182.35,204.13,0.5640153137842114,0.010997863654577228,

0,0.5605969245791881,0.0017458052165547962,-47.62,36.83,0.5605969245791881,0.00019066887990243454,-47.62,36.83,0.5605969245791881,0.0006976208116269899,

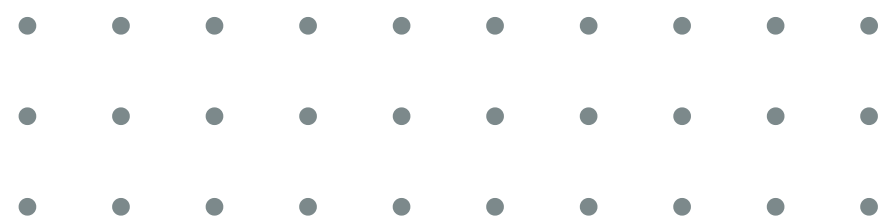
0,0.5440484040547346,0.007612303152725211,-230.28,181.44,0.5440484040547346,0.007148341670780159,-230.28,181.44,0.5440484040547346,0.022243230885227394,

0,0.6041364007937963,0.0009416247298085708,-78.59,87.06,0.6041364007937963,0.0005950599146580347,-78.59,87.06,0.6041364007937963,0.002071913133665327,-

0,0.5636846346070677,0.0009513772881053551,-60.03,36.42,0.5636846346070677,0.0007555438068194456,-60.03,36.42,0.5636846346070677,0.0025757182455460473,-

ML MODELS

- Random Forest: Achieved unexpectedly high accuracy, suggesting potential overfitting or data leakage. Our dataset was limited in the number of data points for each subject.
- Support Vector Machine (SVM): Lower accuracy compared to Random forrest, indicating a more nuanced understanding of the data, potentially reducing overfitting.
- Logistic Regression: Logistic Regression outperformed other models due to its suitability for smaller datasets, where complex models like Random Forest tend to overfit.



RANDOM FOREST

Processing Subject: S17
Accuracy for S17: 0.64
Classification Report for S17:

	precision	recall	f1-score	support
False	0.76	0.72	0.74	18
True	0.38	0.43	0.40	7
accuracy			0.64	25
macro avg	0.57	0.58	0.57	25
weighted avg	0.66	0.64	0.65	25

Final Results:
Accuracies: {'S2': 0.782608695652174, 'S3': 0.9166666666666666, 'S4': 0.7391304347826086, 'S5': 0.8333333333333334, 'S6': 0.75, 'S7': 0.9166666666666666, 'S8': 0.9166666666666666, 'S9': 0.9166666666666666, 'S10': 0.84, 'S11': 0.875, 'S13': 0.7916666666666666, 'S14': 0.8333333333333334, 'S15': 0.7916666666666666, 'S16': 0.8333333333333334, 'S17': 0.64}

```
subject_ids = [f"S{i}" for i in range(2, 18) if i != 12] # Subjects S2 to S17 excluding S12
base_path = 'data/subject_feats/'

# Initialize storage for results
accuracy_scores = {}
classification_reports = {}

for subject_id in subject_ids:
    print(f"\nProcessing Subject: {subject_id}")

    # Load the data
    file_path = f"{base_path}{subject_id}_feats_4.csv"
    data = pd.read_csv(file_path)

    # Select only EDA-related columns for features
    eda_columns = [col for col in data.columns if 'EDA' in col]
    X = data[eda_columns] # Use only EDA columns as features
    y = data['2'] # Assuming '2' is the stress label

    # Split the data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.8, random_state=42)

    # Initialize and train the Random Forest Classifier
    rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
    rf_classifier.fit(X_train, y_train)

    # Predictions
    y_pred = rf_classifier.predict(X_test)

    # Evaluate the model
    accuracy = accuracy_score(y_test, y_pred)
    report = classification_report(y_test, y_pred)
    accuracy_scores[subject_id] = accuracy
    classification_reports[subject_id] = report

    print(f"Accuracy for {subject_id}: {accuracy}")
    print(f"Classification Report for {subject_id}:\n{report}")

# Results
print("\nFinal Results:")
print("Accuracies:", accuracy_scores)
```

SVM

Processing Subject: S16

Accuracy for S16: 0.75

Classification Report for S16:

	precision	recall	f1-score	support
False	0.80	0.89	0.84	18
True	0.50	0.33	0.40	6
accuracy			0.75	24
macro avg	0.65	0.61	0.62	24
weighted avg	0.72	0.75	0.73	24

Final Results:

Accuracies: {'S2': 0.782608695652174, 'S3': 0.7916666666666666, 'S4': 0.782608695652174, 'S5': 0.75, 'S6': 0.75, 'S7': 0.7916666666666666, 'S8': 0.875, 'S9': 0.75, 'S10': 0.72, 'S11': 0.75, 'S13': 0.7916666666666666, 'S14': 0.75, 'S15': 0.75, 'S16': 0.75}

```
subject_ids = [f"S{i}" for i in range(2, 18) if i != 12] # Subjects S2 to S17 excluding S12
base_path = 'data/subject_feats/'
```

```
# Initialize storage for results
accuracy_scores = {}
classification_reports = {}
```

```
count = 0
for subject_id in subject_ids:
    count=count+1
    if count == 15:
        break
    print(f"\nProcessing Subject: {subject_id}")
```

```
# Load the data
file_path = f"{base_path}{subject_id}_feats_4.csv"
data = pd.read_csv(file_path)

# Select only EDA-related columns for features
eda_columns = [col for col in data.columns if 'EDA' in col]
X = data[eda_columns] # Use only EDA columns as features
y = data['2'] # Assuming '2' is the stress label
```

```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

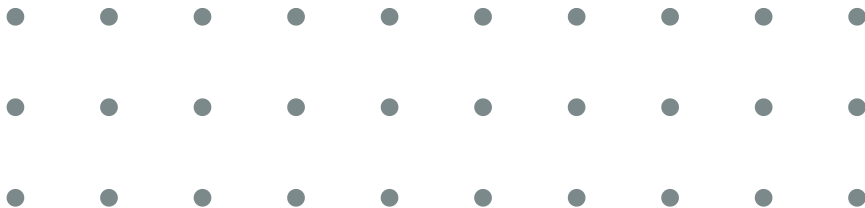
```
# Initialize and train the SVM (variable) y_train: Any
svm_classifier = SVC(kernel='linear')
svm_classifier.fit(X_train, y_train)
```

```
# Predictions
y_pred = svm_classifier.predict(X_test)
```

```
# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
report = classification_report(y_test, y_pred)
accuracy_scores[subject_id] = accuracy
classification_reports[subject_id] = report
```

```
print(f"Accuracy for {subject_id}: {accuracy}")
print(f"Classification Report for {subject_id}:\n{report}")
```

```
# Results
```



LOGISTIC REGRESSION

Processing Subject: S17
Accuracy for S17: 0.64
Classification Report for S17:

	precision	recall	f1-score	support
False	0.76	0.72	0.74	18
True	0.38	0.43	0.40	7
accuracy			0.64	25
macro avg	0.57	0.58	0.57	25
weighted avg	0.66	0.64	0.65	25

Final Results:
Accuracies: {'S2': 0.782608695652174, 'S3': 0.9166666666666666, 'S4': 0.7391304347826086, 'S5': 0.8333333333333334, 'S6': 0.75, 'S7': 0.9166666666666666, 'S8': 0.9166666666666666, 'S9': 0.9166666666666666, 'S10': 0.84, 'S11': 0.875, 'S13': 0.7916666666666666, 'S14': 0.8333333333333334, 'S15': 0.7916666666666666, 'S16': 0.8333333333333334, 'S17': 0.64}

```
subject_ids = [f"S{i}" for i in range(2, 18) if i != 12] # Subjects S2 to S17 excluding S12
base_path = 'data/subject_feats/'

# Initialize storage for results
accuracy_scores = {}
classification_reports = {}

count = 0
for subject_id in subject_ids:
    print(f"\nProcessing Subject: {subject_id}")

    # Load the data
    file_path = f"{base_path}{subject_id}_feats_4.csv"
    data = pd.read_csv(file_path)

    # Select only EDA-related columns for features
    eda_columns = [col for col in data.columns if 'EDA' in col]
    X = data[eda_columns] # Use only EDA columns as features
    y = data['2'] # Assuming '2' is the stress label

    # Split the data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

    # Initialize and train the Logistic Regression Classifier
    logistic_classifier = LogisticRegression(random_state=42)
    logistic_classifier.fit(X_train, y_train)

    # Predictions
    y_pred = logistic_classifier.predict(X_test)

    # Evaluate the model
    accuracy = accuracy_score(y_test, y_pred)
    report = classification_report(y_test, y_pred)
    accuracy_scores[subject_id] = accuracy
    classification_reports[subject_id] = report

    print(f"Accuracy for {subject_id}: {accuracy}")
    print(f"Classification Report for {subject_id}:\n{report}")

# Results
print("\nFinal Results:")
print("Accuracies:", accuracy_scores)
```

COMPARISON OF THE DIFFERENT ML MODELS

