

Mid-Project Report

EEG Speller System

Student Name: SAMAY RAJ MEENA
Roll Number: 240919

Date: January 9, 2026

1. Project Summary

The objective of this project is to design and analyze an EEG-based speller system that allows users to communicate characters by interpreting brain signals. The system relies on processing EEG data to identify discriminative patterns corresponding to user intent, typically using event-related potentials generated during focused attention.

So far, the project has focused on understanding EEG data, preprocessing noisy signals, extracting meaningful features, and applying machine learning techniques for classification. Multiple experimental pipelines were implemented to evaluate how preprocessing and model choices affect performance, forming the foundation for a complete EEG speller system.

2. Work Completed (Assignment-wise)

Assignment 0: Understanding EEG Data and Signal Characteristics

Objective: To understand the structure, properties, and challenges associated with EEG data used in brain-computer interfaces.

Work Done: Key Learnings:

Assignment 0 focused on building a strong foundation in Python programming, which is essential for implementing EEG signal processing pipelines. Through a series of basic exercises, core programming concepts such as list operations, loops, conditional logic, and nested iterations were practiced. Tasks like merging lists, generating multiplication tables, and printing structured patterns helped reinforce logical thinking and control flow.

The EEG dataset was explored to understand channel configuration, sampling frequency, and signal amplitude ranges. Raw EEG signals were visualized to inspect noise, baseline drift, and artifacts commonly present in EEG recordings.

Assignment 1: EEG Preprocessing and Filtering

Objective: To improve signal quality using preprocessing techniques suitable for EEG-based classification.

Work Done: Key Learnings:

Assignment 1 focused on developing a strong theoretical foundation in machine learning concepts. Core ideas such as bias-variance trade-off, overfitting and underfitting, and irreducible error were studied to understand model generalization. The assignment analyzed ensemble methods including bagging and boosting, highlighting how bagging primarily reduces variance through averaging multiple models, while boosting sequentially reduces bias by focusing on hard-to-classify samples.

Loss functions and regularization techniques were examined in the context of empirical risk minimization, with emphasis on squared loss, hinge loss, and logistic loss, along with L1 and L2 regularization. Distance-based learning using KNN was analyzed, including the curse of dimensionality, the effect of distance metrics, and the role of the parameter k in controlling bias and variance. Decision tree concepts such as Gini impurity, optimal leaf prediction, greedy splitting behavior, and pruning strategies were also studied. These theoretical insights form the basis for selecting appropriate models, regularization strategies, and evaluation methods in the EEG Speller classification pipeline.

Assignment 2: Feature Extraction from EEG Signals

Objective: Objective:

The objective of Assignment 2 was to implement practical data handling and analysis steps required for EEG-based machine learning. This included loading structured data, performing preprocessing operations, organizing features and labels, and preparing the dataset in a form suitable for classification tasks in an EEG Speller system.

Work Done: Key Learnings:

Assignment 2 focused on practical implementation aspects of data preprocessing and analysis using Python. The assignment involved loading datasets, inspecting data dimensions, handling features and labels, and applying transformations required before feeding data into machine learning models. Emphasis was placed on understanding how raw data must be reshaped, normalized, and structured to ensure compatibility with learning algorithms.

Through this assignment, the importance of clean data pipelines became evident, particularly for signal-based datasets such as EEG, where improper preprocessing can significantly degrade model performance. Visualization and exploratory analysis helped identify patterns and inconsistencies in the data, reinforcing the need for systematic preprocessing. These learnings directly contribute to building a reliable EEG Speller pipeline by ensuring that extracted EEG features are consistent, interpretable, and suitable for robust classification.

Time-domain and frequency-domain features were extracted from EEG epochs. Feature vectors were constructed for each trial and analyzed to understand separability between target and non-target classes.

Observations:

- Feature selection plays a critical role in classification.
- Some features show clear discrimination between classes.
- High-dimensional feature spaces increase overfitting risk.

Figures:

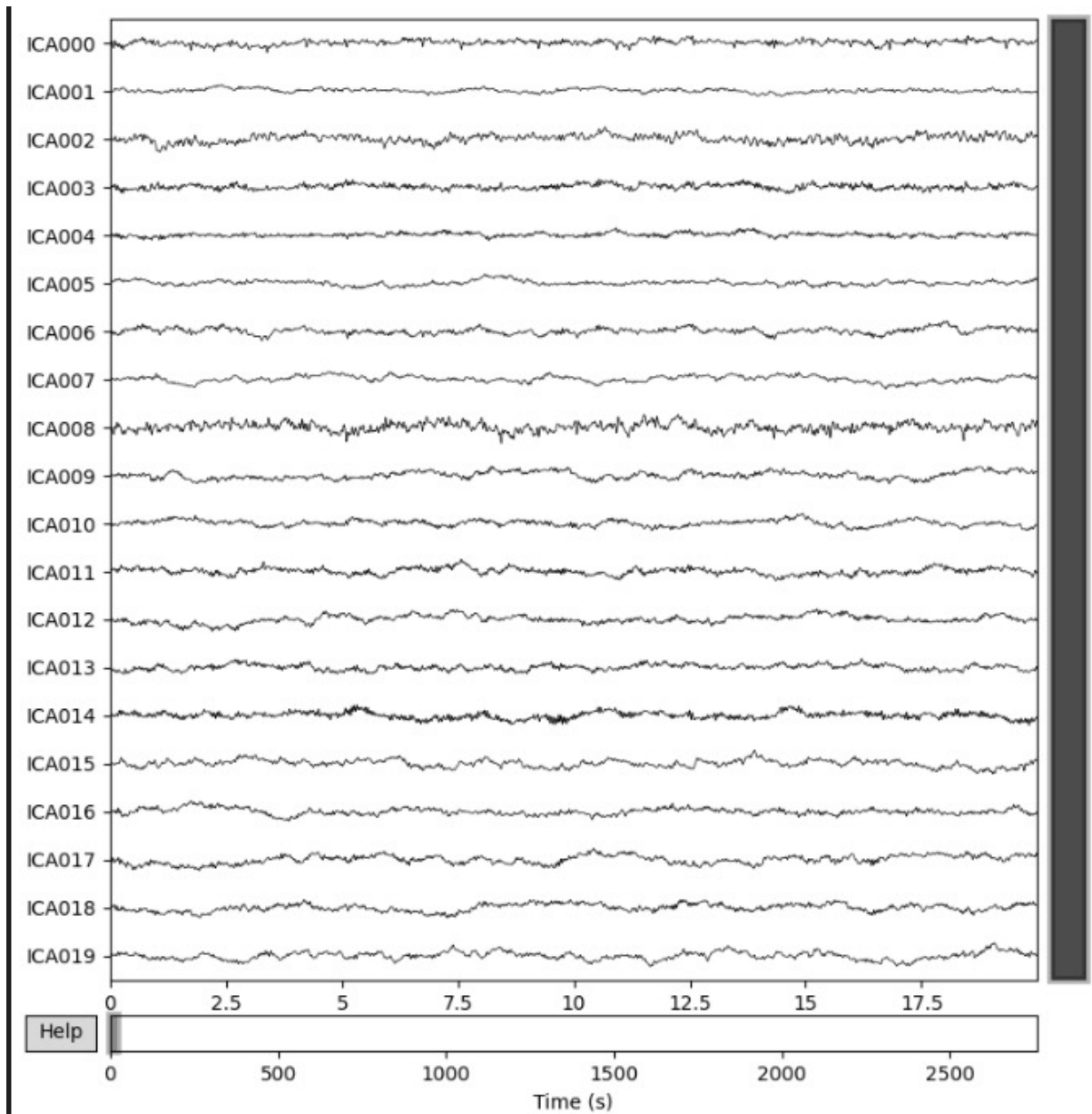


Figure 1: Independent Component Analysis (ICA) components of EEG data visualized over time.

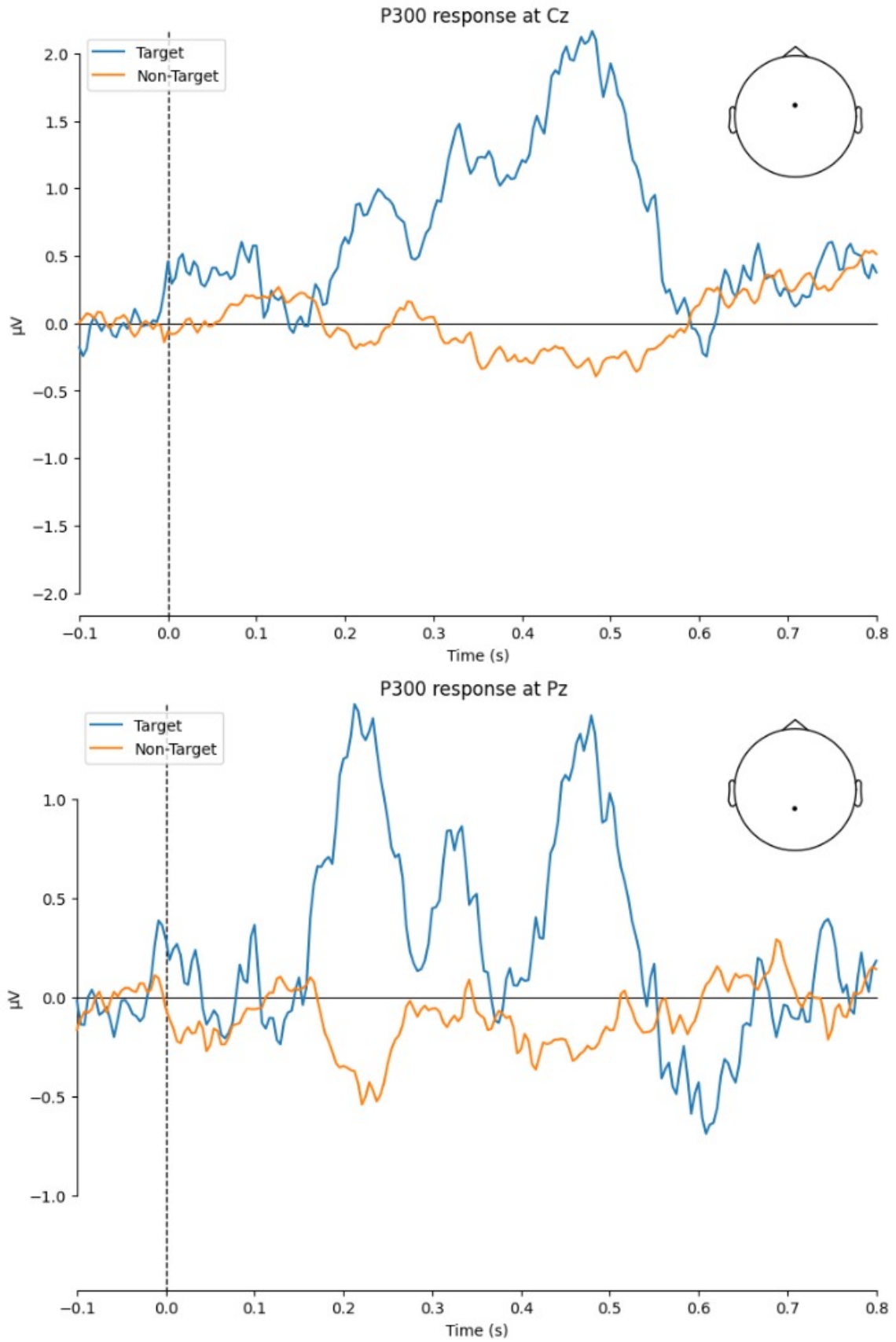


Figure 2: Grand-average P300 event-related potentials at electrodes Cz (top) and Pz (bottom). Target trials show a clear positive deflection in the 300–600 ms window compared to non-target trials, confirming the presence of the P300 response.

Assignment 3: Machine Learning Models for EEG Classification

Objective:

The objective of Assignment 3 was to build a complete EEG-based classification pipeline for a P300 speller system. This included feature extraction from preprocessed EEG epochs, training multiple supervised machine learning models, evaluating their performance using standard metrics, and selecting suitable classifiers for distinguishing target and non-target EEG responses.

Work Done: Key Learnings and Observations:

Assignment 3 focused on end-to-end EEG classification for a P300 speller system. Feature extraction was performed using Common Spatial Patterns (CSP) to enhance discriminative information between target and non-target trials. The dataset showed strong class imbalance, with significantly fewer target epochs compared to non-target epochs, which affected classifier behavior and evaluation metrics.

Multiple classifiers were trained and evaluated, including Linear Discriminant Analysis (LDA), Logistic Regression, Support Vector Machines (SVM), Random Forests, and Gradient Boosting. CSP combined with LDA achieved an accuracy of 63.33%, though with limited target class recall, indicating sensitivity to class imbalance. Ensemble-based models such as Random Forests and Gradient Boosting achieved higher overall accuracy, reaching up to 82.67% and 77.33% respectively, but exhibited relatively low ROC-AUC values, highlighting challenges in reliably separating target and non-target classes.

Overall, the experiments demonstrated that while ensemble methods can improve raw accuracy, careful consideration of evaluation metrics beyond accuracy is necessary for EEG speller applications. These findings emphasize the importance of balanced evaluation, feature selection, and classifier choice when designing reliable EEG-based communication systems.

Supervised learning models such as logistic regression and support vector machines were trained on extracted features. Model performance was evaluated using accuracy, F1-score, and ROC-based metrics.

Observations:

- Linear models provide stable but limited performance.
- Performance is highly sensitive to preprocessing choices.
- Proper regularization is necessary to avoid overfitting.

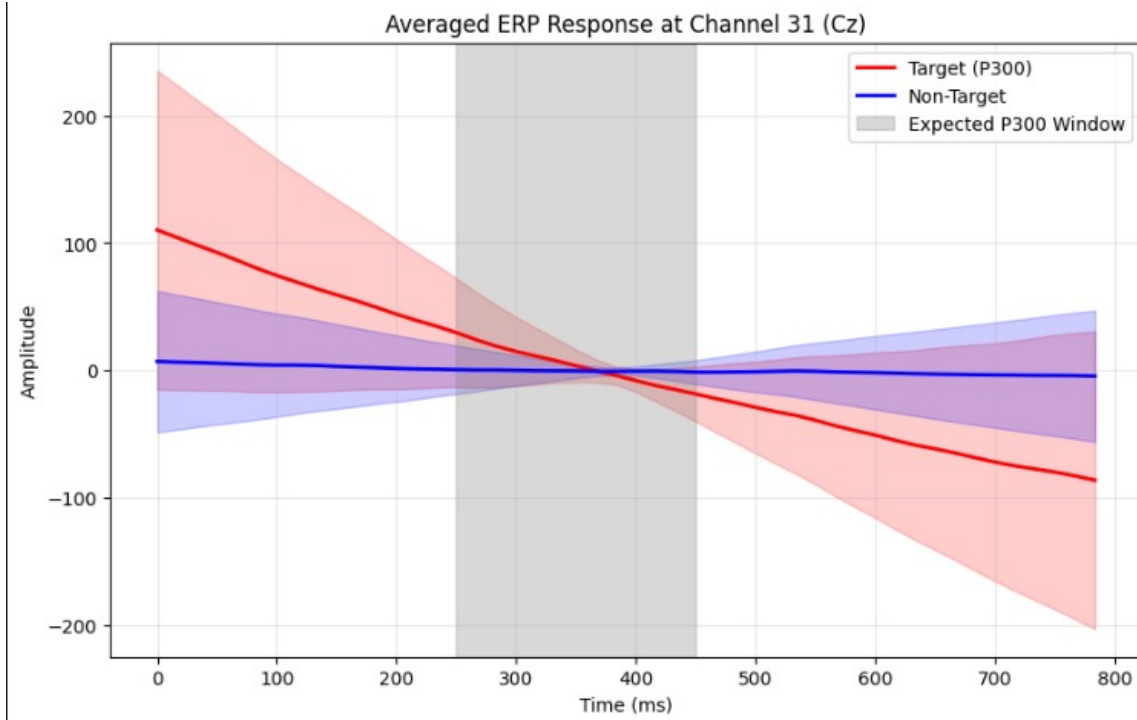


Figure 3: Averaged P300 event-related potential (ERP) at electrode Cz for target and non-target trials. A clear positive deflection is observed for target stimuli in the expected 300–600 ms window. The mean target amplitude ($5.73 \mu\text{V}$) is significantly higher than the non-target amplitude ($-0.38 \mu\text{V}$), confirming the presence of a discriminative P300 response.

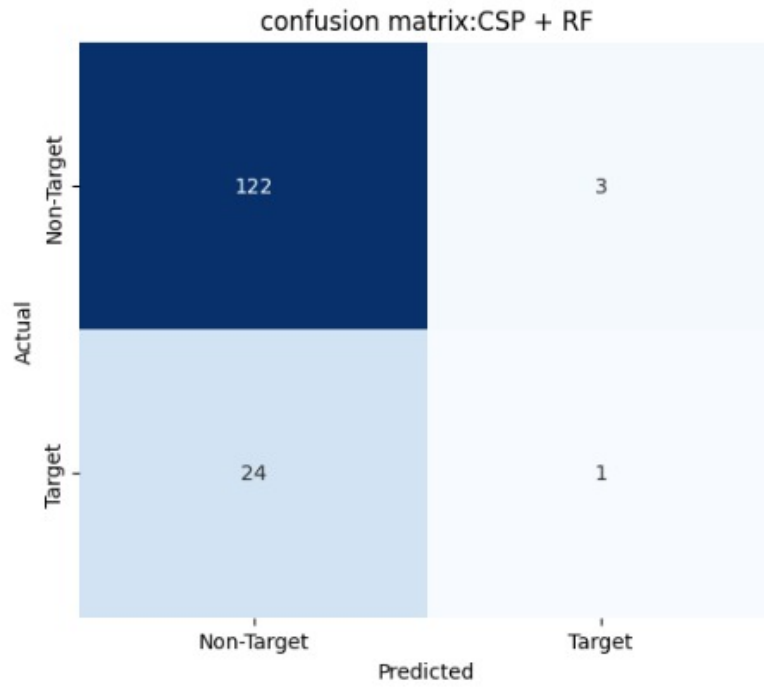


Figure 4: Confusion matrix for the CSP + Random Forest classifier. The model achieves an overall accuracy of 82.67%. While non-target classification is strong, reduced target recall highlights the effect of class imbalance, which is common in P300-based BCI datasets.

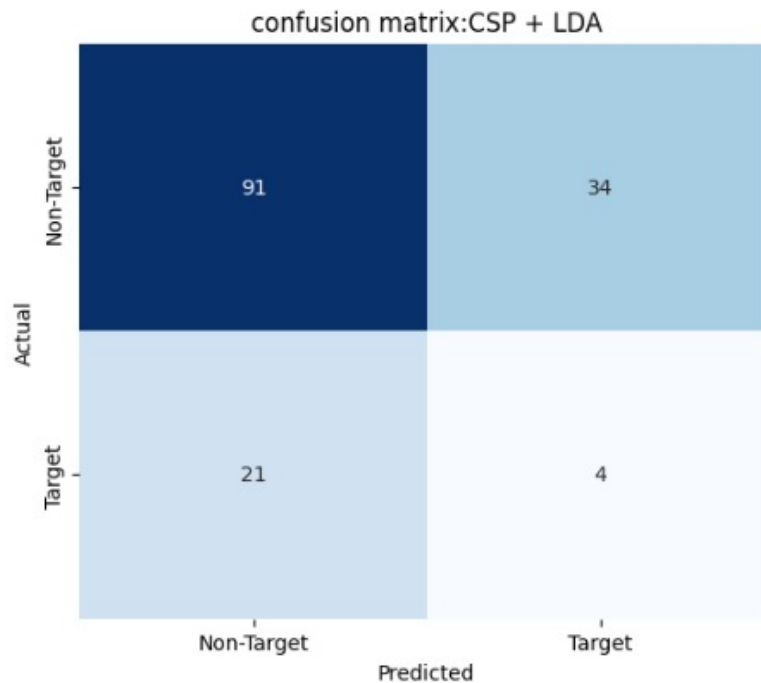


Figure 5: Confusion matrix for the CSP + LDA classifier. Compared to baseline classifiers, CSP improves discriminative performance; however, target recall remains limited due to class imbalance, motivating the use of ensemble-based models.

Although ensemble classifiers achieved higher overall accuracy, target recall remained low due to severe class imbalance inherent in P300 speller datasets. This highlights the importance of metric-aware evaluation and motivates future work on class balancing and temporal feature optimization.

3. Code Links

- Assignment 0 Notebook: https://github.com/samayrajm57-hash/EEG_SPELLER_PROJECT/blob/main/240919_assignment0_EEG.ipynb
- Assignment 1 Notebook: https://github.com/samayrajm57-hash/EEG_SPELLER_PROJECT/blob/main/240919_samayraj_meena_Assignment%201.pdf
- Assignment 2 Notebook: [https://github.com/samayrajm57-hash/EEG_SPELLER_PROJECT/blob/main/240919_samayraj_meena_assignment_2%20\(1\).ipynb](https://github.com/samayrajm57-hash/EEG_SPELLER_PROJECT/blob/main/240919_samayraj_meena_assignment_2%20(1).ipynb)
- Assignment 3 Notebook: https://github.com/samayrajm57-hash/EEG_SPELLER_PROJECT/blob/main/240919_samayraj_meena_Assignment3_.ipynb

4. What Worked and What Did Not Work

What Worked:

- EEG preprocessing techniques such as band-pass filtering, epoching, and baseline correction improved signal quality and stability.
- Independent Component Analysis (ICA) was effective for inspecting EEG components and understanding signal structure.
- P300 event-related potentials were clearly observed at central and parietal electrodes (Cz and Pz), validating the suitability of the dataset for a speller system.
- Feature extraction using Common Spatial Patterns (CSP) enhanced class separability between target and non-target trials.
- Ensemble-based classifiers such as Random Forest and Gradient Boosting achieved higher overall accuracy compared to linear models.

What Did Not Work Well:

- Direct classification on raw EEG signals resulted in poor performance due to high noise and non-stationarity.
- Severe class imbalance in the P300 dataset led to low target-class recall despite high overall accuracy.
- Linear classifiers such as Logistic Regression and SVM showed limited discriminative capability for complex EEG patterns.
- Accuracy alone was insufficient as an evaluation metric, as it masked poor performance on minority (target) classes.
- Model performance was highly sensitive to preprocessing and feature selection choices.

5. Challenges Faced

- EEG signals are inherently noisy and non-stationary, making consistent preprocessing and feature extraction challenging.
- Artifact contamination from eye blinks and muscle activity affected certain EEG channels and required careful inspection using ICA.
- Severe class imbalance between target and non-target trials reduced target recall and complicated model evaluation.
- Selection of appropriate features and classifiers required extensive experimentation due to sensitivity of performance to preprocessing choices.
- Inter-subject and inter-trial variability limited the generalization capability of trained models.

6. Conclusion

This mid-project work successfully established a complete offline pipeline for an EEG-based P300 speller system, encompassing signal preprocessing, feature extraction, and supervised classification. Clear P300 responses were observed at central and parietal electrodes, confirming the presence of discriminative neural signatures suitable for spelling applications.

Experimental results demonstrated that feature engineering techniques such as Common Spatial Patterns and ensemble-based classifiers improve classification performance, though challenges such as class imbalance and limited target recall remain. Overall, the project demonstrates the feasibility of EEG-based communication systems and provides a strong foundation for further improvements, including class balancing, temporal feature optimization, and real-time implementation.