

**EEG Spectrogram Classification using Distributed Machine Learning**  
**Automated EEG Pattern Classification for Neurological Disorder Detection using PySpark MLlib**

---

Course: DS/CMPSC 410

Semester: Fall 2025

Team Members: Krish Chavan, Ryan Hussey, Xu Wang, Hamid Shah, Kuria Mbatia, Nathan Mannings

Team Name: EEG Analysis

---

Introduction and Motivation

Application Domain:

Medical/Healthcare Analytics - Neurological Pattern Recognition

**Background:**

Electroencephalogram (EEG) signals capture brain electrical activity and are crucial for diagnosing neurological disorders like epilepsy. EEG spectrograms represent frequency domain characteristics of brain signals over time, containing complex patterns that indicate different neurological states.

**Key Problem:**

Automated classification of EEG spectrograms into six critical neurological patterns:

1. Seizure - Epileptic seizures requiring immediate medical attention
2. LPD - Lateralized Periodic Discharges
3. GPD - Generalized Periodic Discharges
4. LRDA - Lateralized Rhythmic Delta Activity
5. GRDA - Generalized Rhythmic Delta Activity
6. Other - Normal/baseline brain activity

Why Big Data Approaches Are Essential:

1. High Dimensionality: 400+ spectral frequency features per sample exceed memory capacity of single machines
2. Large Dataset Scale: Clinical EEG data involves thousands of patients with continuous monitoring
3. Complex Feature Interactions: Non-linear relationships between spectral frequencies require distributed ensemble methods
4. Real-time Requirements: Hospital settings need scalable classification for continuous patient monitoring
5. Cross-validation Complexity: Multiple model training with hyperparameter tuning is computationally intensive

## Project Objectives

Primary Goals:

1. Scalable Data Processing: Implement PySpark workflows for handling high-dimensional EEG spectrogram data (400+ features)
  2. Distributed Classification: Develop machine learning pipelines using PySpark MLlib for multi-class neurological pattern recognition
  3. Dimensionality Reduction: Apply PCA within distributed computing framework to manage feature complexity
  4. Model Comparison: Evaluate Random Forest vs. Gradient Boosted Trees performance with/without PCA
  5. Hyperparameter Optimization: Implement distributed cross-validation using TrainValidationSplit
  6. Clinical Visualization: Generate interpretable confusion matrices and feature importance plots for medical professionals
- 

## Data Description

Data Source:

EEG spectrogram dataset stored in Parquet format (course-provided medical dataset)

1. Data Size:
    - a. Records: ~6,700 EEG samples across 6 neurological classes
    - b. Features: 400+ spectral frequency components (f0-f399)
    - c. Storage: Multiple Parquet files requiring distributed loading
    - d. Estimated Size: Several GB of compressed spectral data
  2. Data Characteristics:
    - a. Structured: Tabular format with numerical spectral features
    - b. High-Dimensional: 400+ features per sample
    - c. Multi-class: 6 distinct neurological pattern categories
    - d. Temporal: Spectral features represent frequency content over time windows
  3. Data Challenges:
    - a. High Dimensionality: 400+ features create curse of dimensionality
    - b. Class Imbalance: Uneven distribution across neurological patterns
    - c. Feature Correlation: Spectral frequencies exhibit complex interdependencies
    - d. Memory Constraints: Dataset size exceeds single-machine RAM capacity
    - e. Clinical Interpretability: Need for explainable features for medical diagnosis
- 

## Technical Approach

1. Tools & Frameworks:
    - a. PySpark MLlib: RandomForestClassifier, GBTClassifier, PCA, StringIndexer
    - b. PySpark SQL: DataFrame operations for data preprocessing
    - c. Evaluation: MulticlassClassificationEvaluator, TrainValidationSplit
    - d. Visualization: Matplotlib, Seaborn for result interpretation
    - e. Storage: Parquet format for efficient columnar data access
  2. Cluster Usage Plan:
    - a. Node Configuration: Multi-node Spark cluster on ICDS Roar-Collab
    - b. Job Types:
      - i. Batch processing for data ingestion and preprocessing
      - ii. Iterative ML training with cross-validation
      - iii. Hyperparameter tuning with parallel model evaluation
    - c. Storage: Distributed Parquet files with HDFS/cluster storage
  3. Pipeline Overview:
    - a. Data Ingestion: Load Parquet files into Spark DataFrames with distributed partitioning
    - b. Preprocessing: StringIndexer for label encoding, VectorAssembler for feature preparation
    - c. Feature Engineering: PCA implementation for dimensionality reduction (400+ → reduced space)
    - d. Model Training: Parallel training of Random Forest and GBT classifiers
    - e. Hyperparameter Tuning: TrainValidationSplit with distributed cross-validation
    - f. Evaluation: Multi-class metrics computation and visualization generation
- 

## Results and Discussion

### Data Visualization Reveals:

1. Class Distribution: Imbalanced dataset with "Other" class being most prevalent (23.3%), seizures representing 18.6%
2. Feature Importance: Lower frequency spectral components (f0-f50) show highest importance for neurological pattern discrimination
3. PCA Analysis: First 50 components capture 95% of variance, enabling effective dimensionality reduction
4. Confusion Matrix: GBT model shows strong diagonal performance with minimal cross-class misclassification

### Model Performance Insights:

Best Performing Model: Gradient Boosted Trees + PCA

1. Accuracy: 72.6% (significantly outperforming alternatives)
2. F1-Score: 0.720 indicating balanced precision/recall
3. Key Finding: PCA enhances GBT performance but degrades Random Forest accuracy
4. Clinical Significance: 72.6% accuracy provides reliable automated screening, reducing manual EEG interpretation workload

### Model Comparison Results:

Model	Accuracy	F1-Score	Clinical Interpretation
GBT + PCA	72.6%	0.720	Reliable automated screening
RF (no PCA)	53.0%	0.530	Baseline performance
RF + PCA	45.6%	0.456	PCA hurts Random Forest

#### Scalability on Roar-Collab:

1. Memory Scalability: Distributed processing eliminated 400+ feature memory bottlenecks
  2. Training Acceleration: Parallel hyperparameter tuning reduced experiment time from days to hours
  3. Cross-validation Efficiency: TrainValidationSplit leveraged multiple nodes for simultaneous fold evaluation
  4. Storage Handling: Parquet format enabled efficient columnar access across cluster nodes
- 

#### Code/Model Repository:

```
GitHub Path: DS-410-EEG-Course-Project-main/
├── notebooks/Final.ipynb (Complete ML pipeline)
├── src/train_model.py (Standalone training script)
├── notebooks/cluster_eval.py (Distributed evaluation)
├── plots.py (Visualization generation)
└── data/ (Parquet files - not committed due to size)
```

---

#### Unsolved Problems & Future Work

##### Remaining Challenges:

1. Class Imbalance: Advanced sampling techniques (SMOTE) could improve minority class performance
2. Feature Interpretability: Domain-specific spectral frequency analysis needed for clinical validation
3. Real-time Processing: Streaming EEG classification for live patient monitoring
4. Multi-modal Integration: Combining EEG with other neurological signals (ECG, EMG)

##### Future Improvements:

1. Deep Learning: CNN/RNN architectures for temporal pattern recognition
  2. Ensemble Methods: Stacking multiple distributed models for improved accuracy
  3. Clinical Validation: Collaboration with neurologists for model validation
  4. Edge Deployment: Model compression for real-time bedside monitoring systems
- 

#### CREDIT Statement

[To be filled based on actual team member contributions using Elsevier guidance]

#### Suggested Roles:

1. Conceptualization: Problem formulation and medical domain research
  2. Data Curation: EEG dataset preparation and preprocessing
  3. Formal Analysis: Statistical analysis and model evaluation
  4. Investigation: Literature review and method selection
  5. Methodology: PySpark pipeline design and implementation
  6. Software: Code development and distributed computing setup
  7. Validation: Model testing and performance verification
  8. Visualization: Plot generation and result interpretation
  9. Writing: Report preparation and presentation development
- 

#### References

1. EEG Classification Literature:
  - a. Acharya, U. R., et al. "Automated EEG analysis of epilepsy: A review." *Knowledge-Based Systems* (2013)
  - b. Subasi, A. "EEG signal classification using wavelet feature extraction and a mixture of expert model." *Expert Systems with Applications* (2007)
2. Distributed Machine Learning:
  - a. Zaharia, M., et al. "Apache Spark: A unified engine for big data processing." *Communications of the ACM* (2016)
  - b. Meng, X., et al. "MLlib: Machine learning in Apache Spark." *Journal of Machine Learning Research* (2016)
3. Medical Signal Processing:
  - a. Sanei, S., & Chambers, J. A. "EEG signal processing." John Wiley & Sons (2013)
  - b. Teplan, M. "Fundamentals of EEG measurement." *Measurement Science Review* (2002)
4. Technical Resources:
  - a. PySpark MLlib Documentation: <https://spark.apache.org/docs/latest/ml-guide.html>
  - b. ICDS Roar-Collab User Guide: <https://www.icds.psu.edu/computing-services/roar-collab/>