# Signal Acquisition Team
## Spring 2025 Semester Project
## Real-Time Data Validation Pipeline

Longhorn Neurotech - Software Division

January 2025

# 1 Project Overview

Your mission is to implement a **real-time data validation layer** that catches hardware and signal quality issues *before* data enters the processing pipeline. This is the first line of defense against poor data quality.

## 1.1 Deliverables

- `data_validation.py` - New module with all validation functions

- Integration into `eeg_processor.py` at specified lines

- No modification of existing filtering code

- Documentation of validation metrics

# 2 Technical Background

## 2.1 Why Data Validation Matters

In real-time BCI systems, bad data causes:

- False classifications (user thinks left, system reads right)

- Model training on garbage data

- Inconsistent performance across sessions

Your validation layer detects:

1. **Hardware failures** - disconnected electrodes, broken channels

2. **Signal anomalies** - flatlines, extreme noise, NaN values

3. **Data integrity** - duplicate channels, missing samples

## 2.2 Validation Happens in Two Places

- **Per-sample validation** - Fast checks on every new data chunk (Lines 151-152 in eeg_processor.py)

- **Buffer validation** - Comprehensive checks on accumulated data (Line 177 in eeg_processor.py)

# 3 Component 1: NaN and Inf Detection

## 3.1 What It Is

NaN (Not a Number) and Inf (Infinity) values appear when:

- Division by zero in calculations

- Hardware read errors

- Buffer overflow

## 3.2 Mathematical Definition

For data matrix $X \in R^{C \times T}$ where $C$ = channels, $T$ = time points:

$$\text{Valid}(X) = \begin{cases} \text{True} & \text{if } \forall i, j : X_{ij} \in R \text{ and } |X_{ij}| < \infty \\ \text{False} & \text{otherwise} \end{cases} \quad (1)$$

## 3.3 Implementation

Listing 1: Function: detect_nan_inf

```python
def detect_nan_inf(data):
    """
    Detect NaN or Inf values in EEG data.

    INPUT:
        data: np.ndarray, shape (channels, samples)
                Raw EEG data from board

    OUTPUT:
        is_valid: bool
                    True if no NaN/Inf found
        nan_channels: list
                    Indices of channels containing NaN/Inf

    EXAMPLE:
        >>> data = np.array([[1.0, 2.0, np.nan], [3.0, 4.0, 5.0]])
        >>> is_valid, bad_ch = detect_nan_inf(data)
        >>> print(is_valid)  # False
        >>> print(bad_ch)    # [0]
    """
    # YOUR IMPLEMENTATION HERE
    # Hint: Use np.isnan() and np.isinf()
    pass
```

**Your Task:**

1. Check each channel for NaN values using `np.isnan()`

2. Check each channel for Inf values using `np.isinf()`

3. Return list of bad channel indices

4. Set is_valid = False if any bad values found

## 4 Component 2: Flatline Detection

### 4.1 What It Is

A **flatline** occurs when a channel outputs constant values, indicating:

- Disconnected electrode

- Hardware malfunction

- Saturated amplifier

### 4.2 Detection Method

A channel is flatline if its standard deviation is below threshold:

$$\text{Flatline}(x) = \begin{cases} \text{True} & \text{if } \sigma(x) < \theta_{\text{flat}} \\ \text{False} & \text{otherwise} \end{cases} \tag{2}$$

where $\sigma(x)$ is standard deviation and $\theta_{\text{flat}} = 0.1 \ \mu\text{V}$ (typical threshold).

### 4.3 Implementation

Listing 2: Function: detect_flatline

```
def detect_flatline(data, threshold=0.1):
    """
    Detect flatline channels (constant/near-constant signal).

    INPUT:
        data: np.ndarray, shape (channels, samples)
                EEG data segment
        threshold: float
                    Standard deviation threshold in microvolts
                    Default: 0.1 uV

    OUTPUT:
        flatline_channels: list
                            Indices of flatline channels
        channel_stds: np.ndarray
                        Standard deviation of each channel

    REASONING:
        Healthy EEG has std > 1 uV typically
        Flatline means electrode disconnected or broken
```

```
21        """
22        # YOUR IMPLEMENTATION HERE
23        # Hint: Calculate std per channel, compare to threshold
24        pass
```

# 5 Component 3: Extreme Noise Detection

## 5.1 What It Is

**Extreme noise** appears as:

- Muscle artifacts (EMG contamination)

- Movement artifacts

- Electrical interference spikes

## 5.2 Detection Method

Use $z$-score to detect outliers:

$$z_{ij} = \frac{X_{ij} - \mu_i}{\sigma_i} \tag{3}$$

Spike detected if $|z_{ij}| > \theta_z$ where $\theta_z = 5$ (typical threshold).

## 5.3 Implementation

Listing 3: Function: detect_extreme_noise

```
1  def detect_extreme_noise(data, z_threshold=5.0):
2      """
3      Detect extreme noise/spike artifacts.
4
5      INPUT:
6          data: np.ndarray, shape (channels, samples)
7          z_threshold: float
8                       Z-score threshold for spike detection
9                       Default: 5.0 (5 standard deviations)
10
11     OUTPUT:
12          has_spikes: bool
13                      True if spikes detected in any channel
14          spike_channels: list
15                      Channels containing spikes
16          spike_percentage: float
17                      Percentage of samples that are spikes
18
19     DETECTION LOGIC:
20          For each channel:
21              1. Calculate mean and std
22              2. Compute z-scores: z = (x - mean) / std
23              3. Flag samples where |z| > threshold
24              4. If > 1% of samples are spikes, mark channel bad
```

```
25      """
26      # YOUR IMPLEMENTATION HERE
27      pass
```

# 6  Component 4: Channel Duplication Check

## 6.1  What It Is

Sometimes hardware errors cause:

- Same channel duplicated in multiple positions

- Copy of previous buffer instead of new data

## 6.2  Detection Method

Use correlation between channels:

$$\rho(x_i, x_j) = \frac{\text{cov}(x_i, x_j)}{\sigma_{x_i} \sigma_{x_j}} \tag{4}$$

If $\rho > 0.99$, channels are duplicates.

## 6.3  Implementation

Listing 4: Function: detect_channel_duplication
```
1   def detect_channel_duplication(data, correlation_threshold=0.99):
2       """
3       Detect if channels are duplicates of each other.
4
5       INPUT:
6           data: np.ndarray, shape (channels, samples)
7           correlation_threshold: float
8                                  Pearson correlation threshold
9                                  Default: 0.99
10
11      OUTPUT:
12          duplicate_pairs: list of tuples
13                           [(ch_i, ch_j), ...] duplicate channel pairs
14          correlation_matrix: np.ndarray, shape (channels, channels)
15                              Full correlation matrix
16
17      USAGE:
18          If duplicates found, data is invalid - hardware issue
19      """
20      # YOUR IMPLEMENTATION HERE
21      # Hint: Use np.corrcoef(data) to get correlation matrix
22      pass
```

# 7 Component 5: Validation Package Class

## 7.1 What It Is

A **ValidationPackage** bundles data with its quality metrics.

Listing 5: Class: ValidationPackage

```python
class ValidationPackage:
    """
    Container for validated EEG data with quality metrics.

    ATTRIBUTES:
        data: np.ndarray
                The validated EEG data
        is_valid: bool
                    Overall validity flag
        timestamp: float
                    Timestamp when data was collected
        quality_metrics: dict
                        Dictionary of quality measurements
                        {
                                'has_nan': bool,
                                'has_inf': bool,
                                'flatline_channels': list,
                                'noisy_channels': list,
                                'duplicate_pairs': list,
                                'sampling_rate': float
                        }

    USAGE:
        package = ValidationPackage(data, timestamp)
        if package.is_valid:
            # Safe to process
            process(package.data)
        else:
            # Log the issue
            log_error(package.quality_metrics)
    """
    def __init__(self, data, timestamp):
        self.data = data
        self.timestamp = timestamp
        self.is_valid = True
        self.quality_metrics = {}
```

# 8 Integration Points

## 8.1 File: eeg_processor.py

**Line 1-10:** Add import statement:

```python
# Add after existing imports
from data_validation import (
    detect_nan_inf,
```

```
4        detect_flatline ,
5        detect_extreme_noise ,
6        detect_channel_duplication ,
7        ValidationPackage
8   )
```

**Line 151-152:** Add per-sample validation in `get_recent_data()`:

```
1   # Current line 151:
2   eeg_data = data[self.eeg_channels , :]
3
4   # INSERT AFTER LINE 151:
5   # Validate new data chunk
6   is_valid, nan_channels = detect_nan_inf(eeg_data)
7   if not is_valid:
8       print(f"Warning: NaN/Inf detected in channels {nan_channels}")
9       # Skip this chunk or use last known good data
10      return self.processed_data_buffer[:, -int(duration * self.
            sampling_rate):]
```

**Line 177:** Add buffer validation before returning:

```
1   # Current line 177 returns model_input
2   # INSERT BEFORE RETURN:
3
4   # Comprehensive validation on buffer
5   flatline_ch = detect_flatline(recent_data)
6   if len(flatline_ch) > 0:
7       print(f"Warning: Flatline detected in channels {flatline_ch}")
8
9   has_spikes , spike_ch , spike_pct = detect_extreme_noise(recent_data)
10  if has_spikes and spike_pct > 5.0:
11      print(f"Warning: Excessive noise ({spike_pct:.1f}%) in channels {
            spike_ch}")
12
13  # Then continue with existing return statement
```

# 9   Testing Your Implementation

## 9.1   Test Cases

Listing 6: Test Script - test_validation.py

```
1   import numpy as np
2   from data_validation import *
3
4   def test_nan_detection():
5       """Test NaN detection"""
6       # Create test data with NaN
7       data = np.random.randn(8, 100)
8       data[2, 50] = np.nan
9
10      is_valid, bad_ch = detect_nan_inf(data)
11      assert not is_valid, "Should detect NaN"
```

```
12      assert 2 in bad_ch, "Should identify channel 2"
13      print("    NaN detection works")
14
15  def test_flatline():
16      """Test flatline detection"""
17      data = np.random.randn(8, 100)
18      data[5, :] = 0.0  # Flatline channel 5
19
20      flat_ch, stds = detect_flatline(data)
21      assert 5 in flat_ch, "Should detect flatline in channel 5"
22      print("    Flatline detection works")
23
24  def test_noise():
25      """Test noise detection"""
26      data = np.random.randn(8, 100)
27      data[3, 10:20] = 50.0  # Large spike
28
29      has_spike, spike_ch, pct = detect_extreme_noise(data)
30      assert has_spike, "Should detect spike"
31      print("    Noise detection works")
32
33  if __name__ == "__main__":
34      test_nan_detection()
35      test_flatline()
36      test_noise()
37      print("\nAll tests passed!")
```

# 10    Expected Outcomes

## 10.1    Performance Metrics

- **Validation time:** $< 5$ms per chunk

- **False positive rate:** $< 1\%$

- **True positive rate:** $> 95\%$ on simulated bad data

## 10.2    Documentation Required

1. Comments explaining each validation function

2. Test results showing validation catches errors

3. Performance benchmarks (timing measurements)

4. Integration checklist confirming all insertion points completed

# 11 Grading Rubric

| Component | Points | Criteria |
|---|---|---|
| NaN/Inf Detection | 15 | Correctly identifies bad values, returns channel indices |
| Flatline Detection | 15 | Uses std threshold, identifies constant channels |
| Noise Detection | 20 | Implements z-score method, reports percentage |
| Duplication Check | 15 | Calculates correlations, identifies duplicates |
| Integration | 20 | Code properly inserted at specified lines |
| Testing | 10 | Test script runs and passes all tests |
| Documentation | 5 | Clear comments and docstrings |
| **Total** | **100** | |

# 12 Timeline

- **Week 1-2:** Implement validation functions

- **Week 3:** Integration into eeg_processor.py

- **Week 4:** Testing with synthetic bad data

- **Week 5:** Hardware testing with real OpenBCI board

- **Week 6:** Documentation and final report

# 13 Resources

- NumPy documentation: `https://numpy.org/doc/`

- EEG artifact detection paper: Delorme et al. (2007) - EEGLAB artifact rejection

- BrainFlow API: `https://brainflow.readthedocs.io/`