# Wood classification

**Flemming Morsch**
**Lys Sanz Moreta**
**Zhi Ye**

# Data Exploration

**Overall**

**3 types**

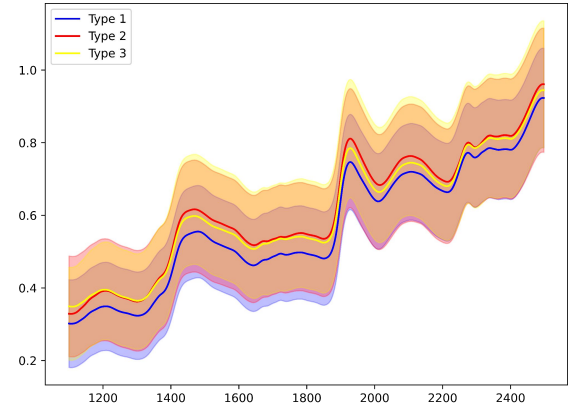**Mean Normalization**

# Data Exploration

# Data Exploration



**Spectrums' trends are different for wood types in some regions**

# Feature Engineering

- Take the 2nd derivative of spectra.
- Type 1 is easy to be distinguished.
- Type 2 and 3 are very similar.

- The 2nd derivative implies the data difference clearly instead of just trends difference.

# Feature Engineering - 2nd derivative

# Feature Engineering - 2nd derivative

# Feature Engineering - 2nd derivative

# Feature Engineering - 2nd derivative

# Unbalanced data

- Unbalanced data set:

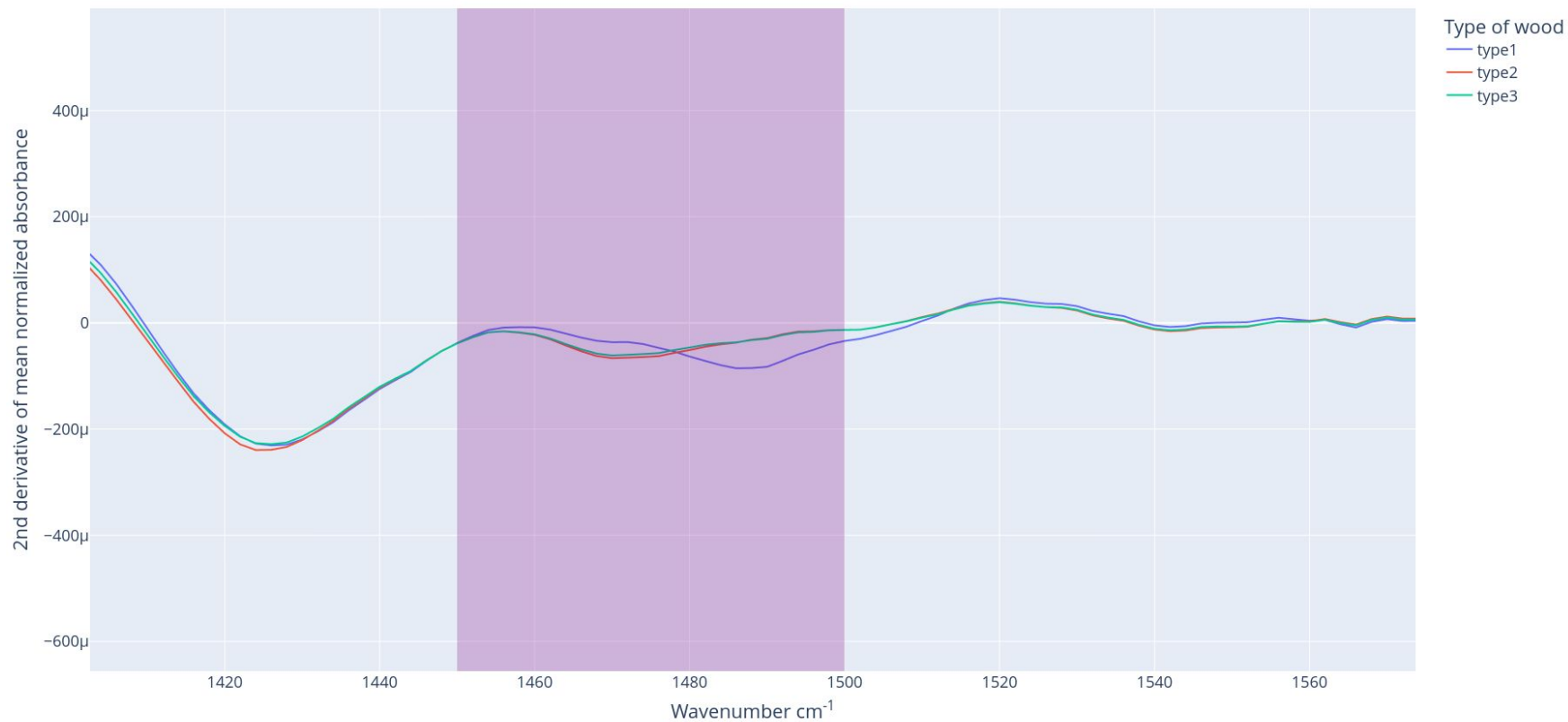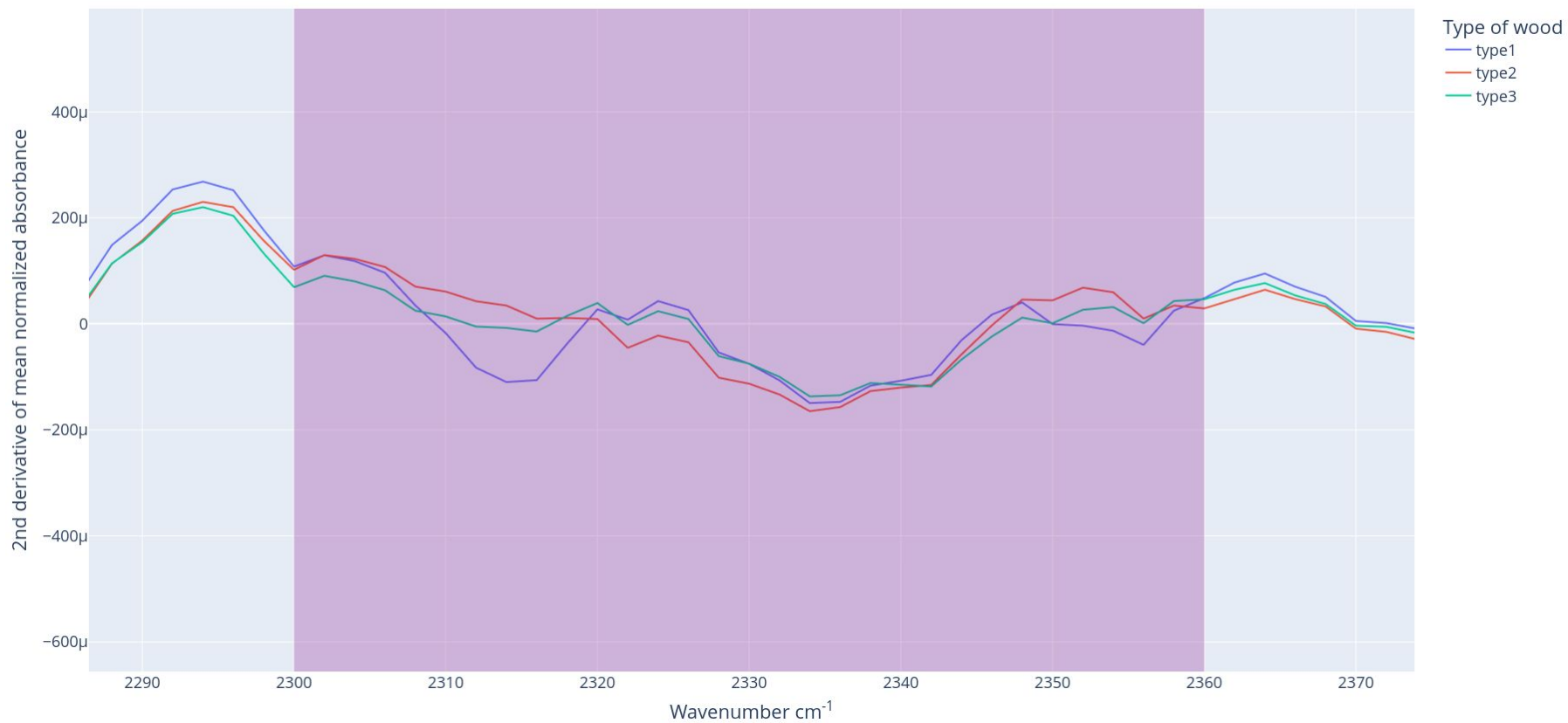| Type 1 | Type 2 | Type 3 |
|--------|--------|--------|
| 4416   | 2784   | 1344   |

- Solution - Upsampling!
  - Synthesize new samples for the minority classes to obtain a balanced data set.
  - SMOTE (Synthetic Minority Oversampling Technique).
  - Choose a random sample from the minority class and compute its 5 nearest neighbors.
  - Randomly selected a neighbor from 5-nearest neighbors and generate a synthetic sample between these two samples in feature space.
- After upsampling we obtain a balanced data set for classification task:

| Type 1 | Type 2 | Type 3 |
|--------|--------|--------|
| 4416   | 4416   | 4416   |

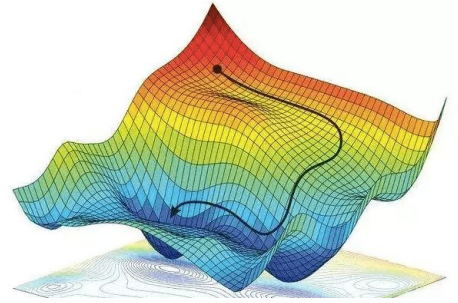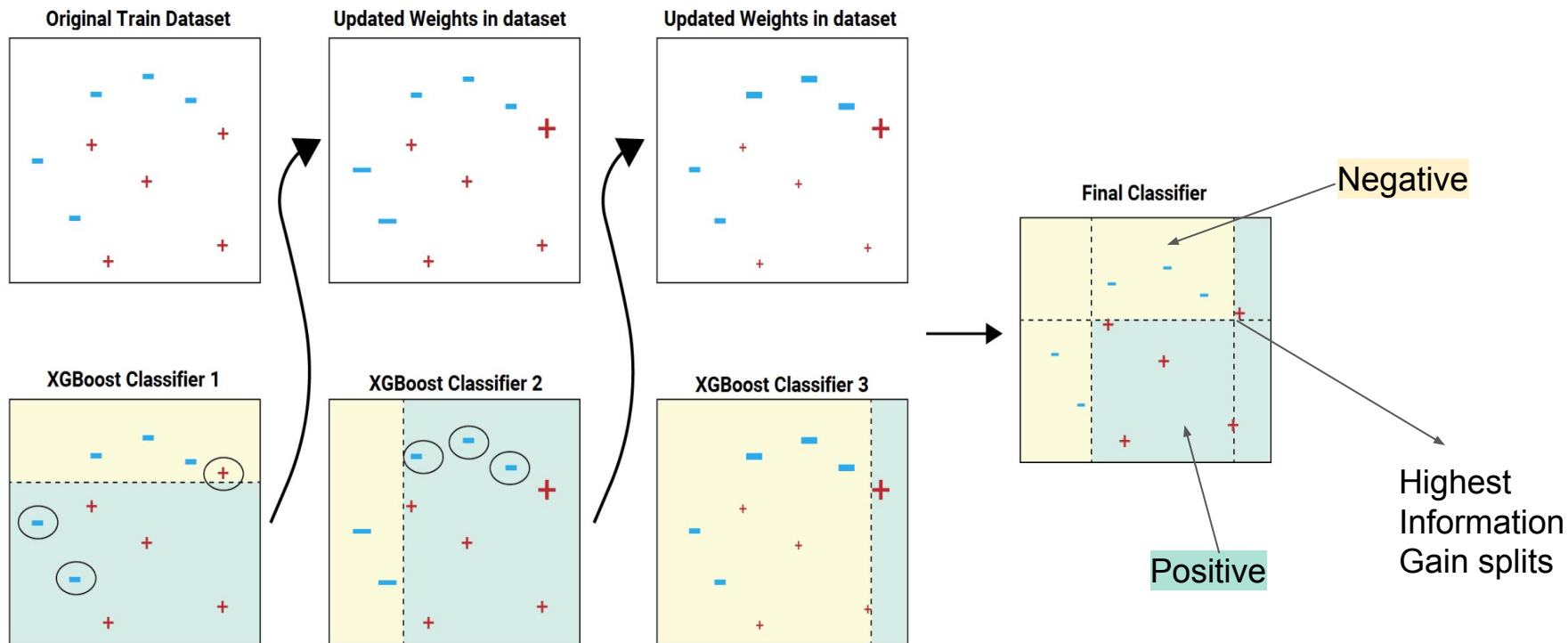# XGB Boost Classifier



- **Extreme Gradient Boosting Decision tree.**

- Minimizes error loss function using **gradient descent**

- Makes use of **gradient boosting.** Trees are trained sequentially
- **Parallelized tree building** for extreme computational performance.

- Many hyperparameters: tree depth, number of trees, regularization,...

# XGB Boost Classifier



Original Train Dataset

Updated Weights in dataset

Updated Weights in dataset

XGBoost Classifier 1

XGBoost Classifier 2

XGBoost Classifier 3

Final Classifier

Negative

Positive

Highest Information Gain splits

# Results

| Class | Precision | Recall | F1-score | Support |
|-------|-----------|--------|----------|---------|
| **1** | 0.90 | 0.78 | 0.84 | 36 |
| **2** | 0.73 | 1.00 | 0.84 | 27 |
| **3** | 0.94 | 0.71 | 0.81 | 21 |
| | | | | |
| **average** | **0.86** | **0.83** | **0.83** | |
| **accuracy** | | | | **0.83** |

# Future work

- Consistency: Model not very consistent across scans/replicates within same sample

- Try different model types

- Time limitation: Comprehensive Hyperparameter search