
ST443 Group 1 Project

Abstract

This report investigates the performance of machine learning models across two high-dimensional classification tasks. Task 1 addresses multi-class land-type classification using pixel-level image features, while Task 2 focuses on feature selection for binary decision classification based on neural activity recordings. Both tasks begin with exploratory data analysis, covering data structure, distributions, and feature characteristics. A range of models is then trained and evaluated using multiple performance metrics. Task 1 also includes a brief real-world application, whereas Task 2 emphasises identifying informative features in a highly sparse, high-dimensional dataset.

1 Task 1: multi-class classification

1.1 Methodology

The high-dimensional pixel dataset was preprocessed by detecting invalid entries and extracting all spectral-band features. We then split the data into an 80% training set (with a small subsample used for particular models) and a 20% test set using stratified sampling to preserve class proportions. All models were implemented using `scikit-learn` pipelines with standardisation. We additionally evaluated PCA versions of all models, using 10 principal components fitted on the training set.

Each classifier—*Discriminant Analysis*, *Logistic Regression*, *k-NN*, *Tree-based Models*, and *SVM*—was first evaluated using 5-fold stratified cross-validation, computing *Accuracy*, *Balanced Accuracy*, *Macro-F1*, *Macro-AUC*, and confusion matrices. The models were then retrained on the complete 80% training set and assessed once on the untouched 20% test set to obtain unbiased performance estimates and identify candidates for the final glacier–ice binary classification task.

1.2 Visualisation and summary statistics

We identified the dataset as having dimensions $215,604 \times 223$, with 218 spectral band features and no missing values or duplicates. Since reflectance values must lie within $[0,1]$, we checked for invalid entries and found about 0.84% outside this range; these were clipped to maintain consistency. Approximately 5.17% of values were flagged as outliers under a 3-sigma rule. But given the small proportion and the fact that extreme reflectance (e.g., from bright snow) can occur naturally, we retained them.

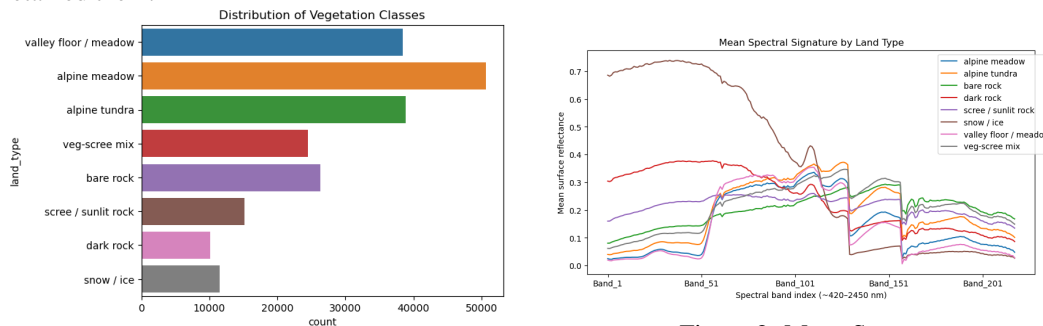


Figure 1: Class-wise distribution

Figure 2: Mean Spectra

For visual exploration, we first examined class distributions and inter-class relationships. From Figure 1, *alpine meadow* is the most frequent class, accounting for about a quarter of all observations,

while *valley floor* and *alpine tundra* are also common, together making up around 58% of the dataset. This imbalance highlights the need for careful stratification in cross-validation, as uniform splits could cause minority classes to vanish in some folds. Figure 2 shows that classes with distinct spectral curves should be easier to separate, whereas overlapping curves indicate spectrally similar classes where linear methods may struggle, and nonlinear models may perform better. The smooth variation across neighbouring bands further suggests strong correlation among features, supporting the use of PCA or regularisation in later modelling.

To examine these correlations more closely and assess whether PCA is appropriate for this feature space, we plot the correlation heatmaps for adjacent bands and for all spectral bands:

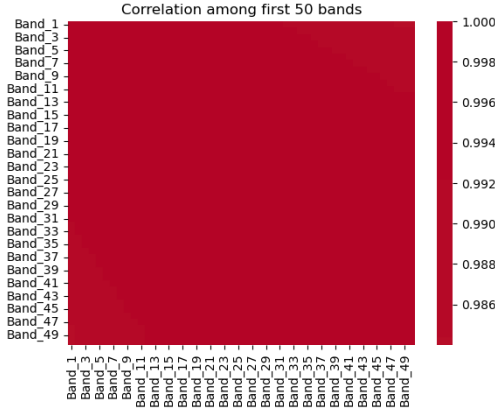


Figure 3: Adjacent correlation heatmap

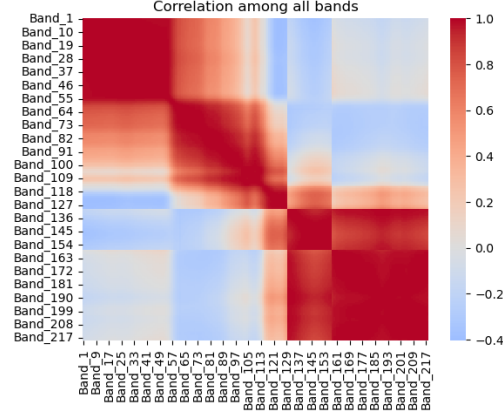


Figure 4: All bands correlation heatmap

Figure 3 shows strong correlations among adjacent spectral bands, meaning many channels capture almost the same information. In contrast, the blue regions in Figure 4 represent bands from different spectral ranges, showing lower or even negative correlation. This block-like pattern highlights substantial multicollinearity among the features and motivates using PCA to compress redundant bands into a smaller set of orthogonal components while retaining most spectral variation.

As the last step before modelling, we applied PCA and clustering as exploratory methods. PCA reduces the highly correlated, high-dimensional band data to a smaller set of informative components, helping reveal broad structure in the data. A test PCA with four components explains 99.87% of the variance, but K-means on the first two components offers limited interpretive value, with clusters appearing arbitrarily separated (Figure 5). Therefore, we use PCA solely for dimensionality reduction before model fitting.

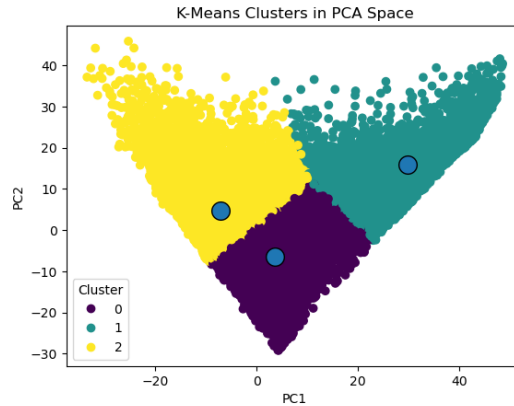


Figure 5: K-Means clustering for the 2 testing principal components

1.3 Model development

In this part of the investigation, we aim to predict vegetation type from pixel reflectance values. We consider several types of classifiers for this multi-class classification problem. Categorical labels are

encoded numerically using `scikit-learn`’s `LabelEncoder`. Our methodology is as follows: we evaluate each model using cross-validation on the training set, select the best-performing models, and then retrain them on the full training portion before evaluating performance on the test set. The test set is used only for final reporting.

1.3.1 Inputs and evaluation framework

We use all available spectral-band measurements as the feature matrix X and the land-type labels as the target variable y . Before fitting models, we standardise all features to ensure comparability across bands, and we evaluate each classifier on both the raw features and a reduced-dimensional representation obtained via PCA with 10 components.

To reliably assess model performance, we use 5-fold stratified cross-validation on the training set. Stratification preserves class proportions in each fold, which is especially important given the imbalanced land-type distribution observed earlier. For each model, we compute several performance metrics: *Accuracy*, *Balanced Accuracy*, *Macro-F1*, and *Macro-AUC*. The *Macro* versions average performance across classes, making them more sensitive to imbalance. After cross-validation, models are retrained on the full 80% training split and evaluated once on the untouched 20% test set to obtain unbiased performance estimates.

1.3.2 Linear and Quadratic Discriminant Analysis (LDA & QDA)

LDA assumes that each class follows a multivariate Gaussian distribution with its own mean vector but a shared covariance matrix, leading to linear decision boundaries. The discriminant function is:

$$\delta_k(x) = x^\top \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \log \pi_k,$$

QDA relaxes the *LDA* assumption by allowing each class to have its own covariance matrix. This yields more flexible quadratic decision boundaries. The discriminant function is:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^\top \Sigma_k^{-1} (x - \mu_k) + \log \pi_k,$$

Both classifiers predict based on the rule of:

$$\psi(x) = \arg \max_k \delta_k(x).$$

1.3.3 Logistic Classifier

As one of the most fundamental models in machine learning, the *Logistic Classifier* estimates the probability of each class given the input features. In the binary case, it models the log-odds of class membership as a linear function of the predictors. For multi-class problems, logistic regression uses the softmax function to model probabilities across all classes:

$$P(y = k | x) = \frac{\exp(\beta_k^\top x)}{\sum_{j=1}^K \exp(\beta_j^\top x)}.$$

In addition, to ensure parameter convergence, we trained the classifier for 2000 iterations, providing faster, more stable optimisation.

1.3.4 Tree-type models

Gradient Boosting Decision Trees (GBDT) build an additive ensemble of shallow trees, where each new tree is trained to correct the residual errors of the previous ones. By iteratively minimising a differentiable loss via gradient descent in function space, GBDT captures nonlinear relationships and class-specific patterns. Although it is more sensitive to hyperparameters, it often performs well on high-dimensional, structured data. We chose the GBDT with 100 trees, a learning rate of 0.15 to ensure stable operation and regularisation, shallow trees with `max_depth=2` to prevent overfitting, and used only 70% of samples per tree to reduce variance:

$$F_m(x) = F_{m-1}(x) + \nu h_m(x), \quad m = 1, \dots, M,$$

$$F_0(x) = \arg \min_c \sum_i \ell(y_i, c),$$

An alternative was the *Random Forest* model, which leveraged multiple simple decision trees to reduce variance when averaging the trees. In this case, we set the forest to have 100 trees and a minimum sample split of 2.

1.3.5 Nearest Neighbour Classifier (k-NN)

K-Nearest Neighbours is a simple, non-parametric classification algorithm that predicts the label of a new sample by examining the labels of its k closest training samples under a chosen distance metric (usually Euclidean distance). In high-dimensional settings like task 1, distances between points become less informative, thereby reducing k -NN’s discriminative power:

$$\psi_{\text{kNN}}(x) := \arg \max_k \sum_{i=1}^K \mathbf{1}\{Y^{(i)} = k\},$$

However, due to the high dimensionality of the data, k -NN ran out of memory. To address this while preserving fair comparability across models, we trained k -NN on a randomly selected subset of the training data. Crucially, the final evaluation was still performed on the same 20% untouched test set used by all other models, ensuring consistency in evaluation and avoiding data leakage. Additionally, we chose $k=7$ to balance model stability and flexibility, and used distance weights to reduce the influence of irrelevant, far-away neighbours.

1.3.6 Support Vector Machine (SVM)

Support Vector Machines classify data by finding the optimal separating hyperplane that maximises the margin between classes. Only a subset of the training data—the support vectors—determines this boundary, making SVM robust and effective in high-dimensional feature spaces:

$$\psi_{\text{SVM}}(x) = \text{sign} \left(\beta_0^* + \sum_{i \in \mathcal{S}} \alpha_i^* y_i k(x, x_i) \right),$$

To ensure comparability across models, the SVM used the same training–testing configuration as k -NN. We set the SVM misclassification penalty to $C=20$; such a large value forces the classifier to fit the training data more tightly by penalising errors more heavily, which is helpful in our setting because the classes become more separable after PCA. We also set $k(x, x_i)$ as the Gaussian (RBF) kernel, which allows the SVM to model nonlinear relationships in high-dimensional spectral data by measuring similarity via distance.

1.4 Model performance

As mentioned in the methodology, we split the data into two parts: 80% for training (with a further subset for SVM and k -NN) and 20% for testing, and used 5-fold cross-validation to prevent overfitting. Now that we have fitted all the desired models, it is time to examine their performances using various evaluation metrics, all computed on the same untouched 20% test set.

Model	Acc	BalAcc	AUC	F1
Logistic	0.9922	0.9911	0.99995	0.99116
LDA	0.8613	0.8598	0.99080	0.86179
QDA	0.9217	0.9183	0.99558	0.91886
RF	0.9849	0.9832	0.99987	0.98353
GBDT	0.9475	0.9521	0.99842	0.95305
KNN	0.9613	0.9595	0.99869	0.96081
SVM	0.9904	0.9895	0.99995	0.98986

(a) Raw Models

(b) PCA-10 Models

Table 1: Performance comparison of models with and without PCA.

As shown in Table 1a, most raw models performed very well across all four metrics, though with some noticeable differences. Logistic Regression, SVM, and RF consistently achieved high Accuracy, Balanced Accuracy, AUC, and F1 score, suggesting that they handled the data’s high dimensionality effectively. In contrast, LDA and QDA performed weaker across most metrics, likely due to their stronger assumptions on class covariance that do not hold for this dataset. GBDT and k -NN fell somewhere in between, performing reasonably but not reaching the top-performing models.

After applying PCA with 10 components (Table 1b), the overall ranking of models remained similar. Logistic Regression and SVM stayed stable and strong, while RF showed a slight drop—as expected, since tree-based models do not always benefit from linear dimension reduction. LDA and QDA continued to struggle after PCA, suggesting that the generative approach to classification may not be suitable for this dataset, especially if the features do not follow class-conditional Multivariate Normal assumptions. In contrast, k-NN performance remained broadly consistent.

In the end, we implemented a `mypredict()` function to train our final chosen model on the full labelled dataset and generate predictions for the external test set provided. This function automates the complete pipeline for Task 1, including preprocessing, model fitting, prediction, and output.

1.5 Application: glacier-ice detection

For this task, we relabelled *glacier* as the positive class and merged all other land types into the negative class, reducing the original multi-class problem to a binary one. The *F1 score*—the harmonic mean of precision and recall—provides a balanced measure of performance by penalising both false positives and false negatives, making it more suitable than accuracy or *AUC* when identifying a comparatively rare class. The three best-performing models were PCA10 Logistic Regression (0.99248), raw SVM (0.98904), and raw Random Forest (0.98353).

These models use the same configurations as before: 2000 iterations for Logistic Regression, an RBF kernel with $C=20$ for SVM, and 100 trees (minimum split 2) for Random Forest. The key difference is the binary target, which allows use of the standard (non-macro) F1 score and enables SVM to train on the full dataset rather than a subsample, since the two-class formulation is far less computationally demanding than the multi-class case.

Model	F1
Logistic_PCA10	0.997177
RF_Raw	0.990224
SVM_Raw	0.996745

Table 2: *F1 scores*

Table 2 summarises the F1 scores of the three evaluated models. Logistic_PCA10 achieved the highest score (0.9972), indicating that PCA preserved nearly all discriminative information. Raw SVM performed similarly well ($F1 = 0.9967$), while Raw RF obtained a slightly lower score (0.9902) but still showed strong predictive ability. Overall, the slight performance differences suggest that the underlying signal is highly separable and that multiple model classes can achieve excellent Accuracy. We may gain more intuitive insight from their corresponding confusion matrices:

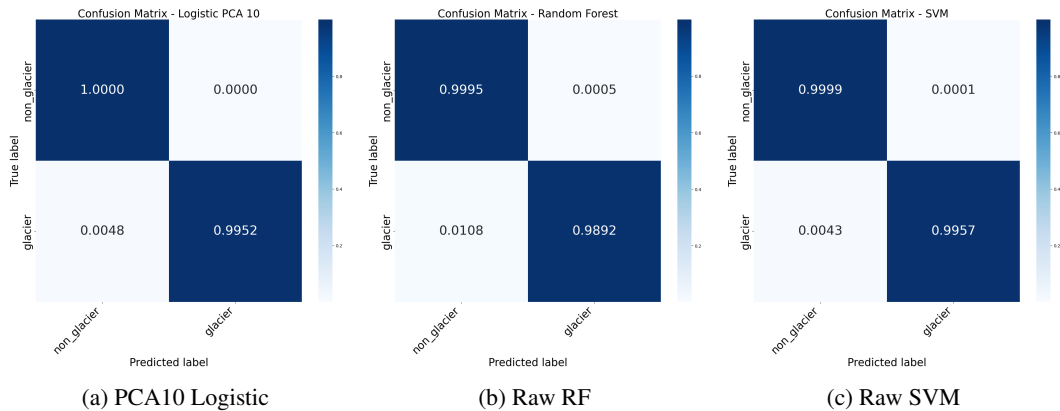


Figure 6: Comparison of confusion matrices for the three models.

All three confusion matrices (Figure 6) show that the models correctly classify almost all samples in both classes, consistent with their very high F1 scores. PCA10 Logistic Regression and the raw SVM achieve nearly perfect separation, with misclassification rates below 0.01% for both classes. RF performs similarly across the non-glacier class but shows a slightly higher false-negative rate for glacier pixels, consistent with its relatively lower F1 score. Overall, the matrices confirm that all models are highly accurate, with only minor differences in how they handle glacier misclassifications.

2 Task 2: feature selection

This task, with its primary focus on feature selection, differs significantly from the previous classification task. Tools used for feature selection are called *Selector*. The three main types of selectors are: *Embedded Selector*, *Filter Selector*, and *Wrapping Selector*. For this large dataset, the Wrapping method isn't suitable due to its high computational cost, so we only investigate the first two types.

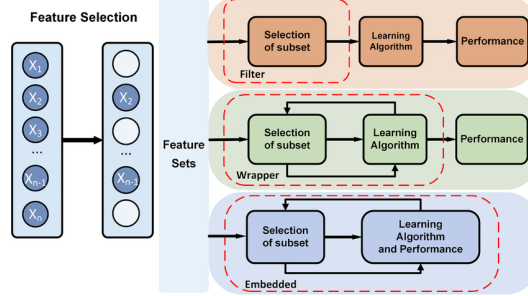


Figure 7: Different types of selectors (ResearchGate, 2020)

2.1 Methodology

Our analysis followed a structured modelling pipeline consisting of data preparation, feature selection, model training, and evaluation. After creating a stratified 8:2 train–test split, all models were implemented using `scikit-learn` pipelines to ensure that preprocessing, feature selection, and classification were applied consistently across different methods. This unified structure allowed us to compare feature-selection techniques fairly while keeping the classifier fixed. We keep our core logic for this task as:

$$Selector \rightarrow Classifier$$

We evaluated several feature-selection methods—L1, Random Forest, XGBoost, and Mutual Information—covering both embedded and filter-based approaches, each tested across multiple feature counts. All pipelines were assessed using stratified 5-fold cross-validation on the training data, followed by evaluation on the held-out test set. Balanced accuracy and confusion matrices served as the primary performance metrics, and a unified evaluation framework enabled consistent comparison across models and feature numbers $k \in \{20, 50, 100, 200\}$. This workflow provided a clear, reproducible procedure for identifying the most effective feature selection strategies.

2.2 Exploratory Data Analysis (EDA)

The high-dimensional ($n < p$) neural spike dataset 11,191×683 contained no missing values, duplicates, negative entries, or non-integer counts. We observed 536 neuron features that never fired (all zeros); these were retained, as silent neurons may simply be inactive during the recorded time window. Because the observations are discrete spike counts, the data are non-Normal. Outlier detection using the IQR rule—flagging values outside $Q_1 - 3IQR$ to $Q_3 + 3IQR$ —identified no outliers. Now, let's have an idea of the data distribution:

Figure 8 shows histogram examples from twelve feature channels, which are generally right-skewed and contain many zeros—approximately 74.03% of all entries by computation. The 200-sample correlation heatmap (Figure 9) also shows weak correlation structure, with mostly grey tones and only a few isolated coloured blocks; strong correlations appear only along the diagonal, reflecting individual feature variances. These observations suggest that a substantial proportion of features are likely uninformative. This motivates and justifies the application of feature selection and dimensionality-reduction techniques in the following analysis.

2.3 Model development

During model development, each feature selector is combined with the fixed classifier to form a pipeline, which is evaluated across multiple feature counts.

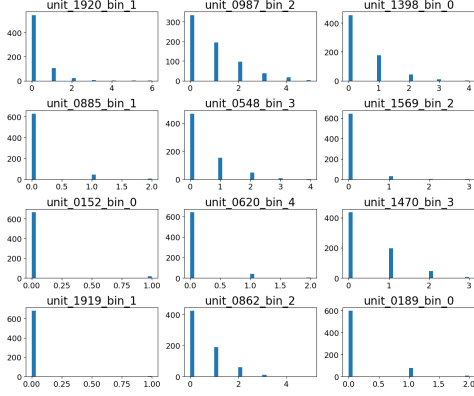


Figure 8: Sample distribution of features

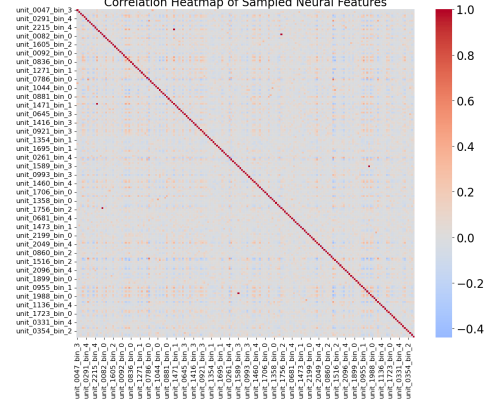


Figure 9: Sample heatmap of features

2.3.1 Inputs and evaluation framework

We use the full neural feature matrix as the input X and the binary class labels as the target variable y . To ensure comparability across feature-selection methods, each selector is paired with the same downstream classifier—Logistic Regression with L2 regularisation. This model is chosen because it is fast, stable, and well-suited to high-dimensional data, allowing differences in performance to be attributed directly to the feature selector rather than to the choice of classifier. For every method, we evaluate multiple feature-set sizes by varying the number of selected features k .

Model performance is assessed using 5-fold stratified cross-validation on the training set, ensuring equal class balance across folds. For each selector–classifier pipeline, we compute Balanced Accuracy as the primary metric, as it accounts for class imbalance and gives equal weight to both classes. After cross-validation, each pipeline is retrained on the full training split and evaluated once on the untouched test set to produce an unbiased estimate of predictive performance. Confusion matrices further illustrate how each model behaves under different feature-selection configurations.

2.3.2 Embedded Selectors

Embedded selectors perform feature selection as part of the model training process itself. Instead of ranking features independently (as in filter methods), these approaches learn feature importance directly from a fitted model—typically through regularisation penalties or tree-based splitting criteria.

2.3.3 Filter Selector

Here we introduce another type of selector: *Filter*. Filter Selectors evaluate each feature independently based on its statistical relationship with the target variable, without involving any classifier (only in the selection step; in the later step, a classifier is still needed to deliver a classification result). Because they operate directly on the data rather than through model training, filter methods are computationally efficient and scale well to high-dimensional settings (Figure 7). Common Filter Selectors are: *ANOVA F-test*, *Chi-square Test*, etc.

In this analysis, we employ Mutual Information as our filter-based selector. MI measures the dependency between each feature and the class label, capturing both linear and nonlinear associations. In practice, MI is estimated using a k -nearest-neighbours (kNN) density estimator, where the choice of $n_neighbors$ affects the smoothness and stability of the scores:

$$I(X; Y) = \mathbb{E} \left[\log \frac{p(X, Y)}{p(X)p(Y)} \right].$$

We apply MI through `SelectKBest`, and a pipeline that standardises the data to ensure reproducibility, ranks features by their MI scores, and then trains the common Logistic Regression classifier on the top k selected features. We also set $n_neighbors = 10$ (instead of the default 3) to stabilise k-NN density estimation and improve performance. A larger neighbourhood reduces variance in the MI scores, leading to more reliable feature rankings across runs.

2.4 Model performance

3 Summary

References

- [1] ResearchGate (2020) *Feature selection methods such as filter, wrapper and embedded method*. Available at: https://www.researchgate.net/figure/Feature-selection-methods-such-as-filter-wrapper-and-embedded-method_fig1_345579532 (Accessed: 27 November 2025).

A Appendix 1: Acknowledgement: The use of generative AI tools

We acknowledge the use of GPT-5 (OpenAI, <https://openai.com/gpt-5>) to generate, debug, and assess portions of the Python code in this investigation, particularly during feature selector construction.