



PAGE BLOCKS CLASSIFICATION

**Presented by Shangzhi LOU,
Wenbo SUI, Limin TIAN**

CONTENT

01

RESEARCH PURPOSES

02

EXPLORE DATASET

03

DATA MODEL ANALYSIS

04

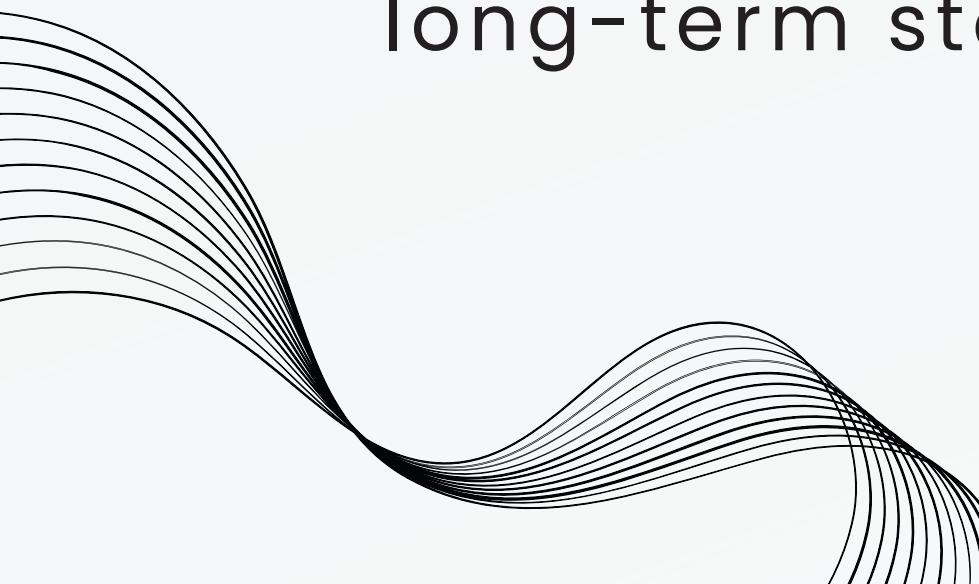
MODEL EVALUATION

05

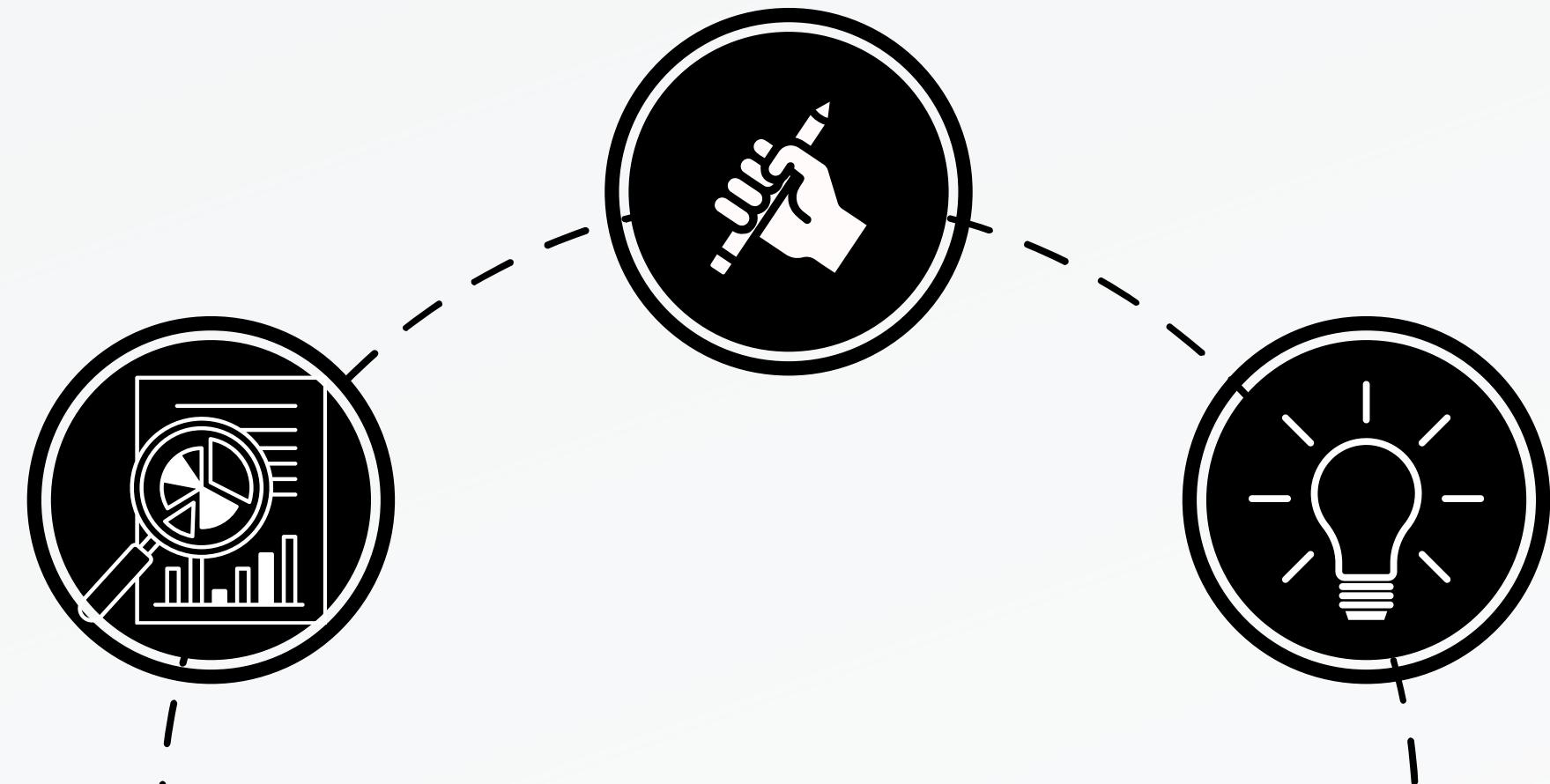
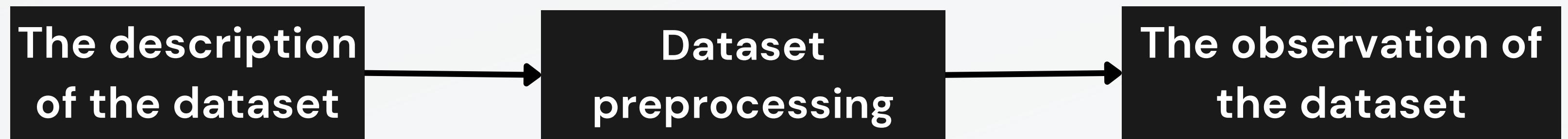
RESULTS

RESEARCH PURPOSES

- **Information extraction and indexing:** Page layout analysis can help identify and extract different types of information blocks such as text, images, and tables.
- **Learning and AI applications:** Page block classification is a fundamental task in the field of machine vision and AI, helping to improve the ability of computer vision to identify document content.
- **Digital Archiving:** When digitizing paper documents, page chunking can help create a more structured digital copy. This is critical for long-term storage, information retrieval and digital protection.



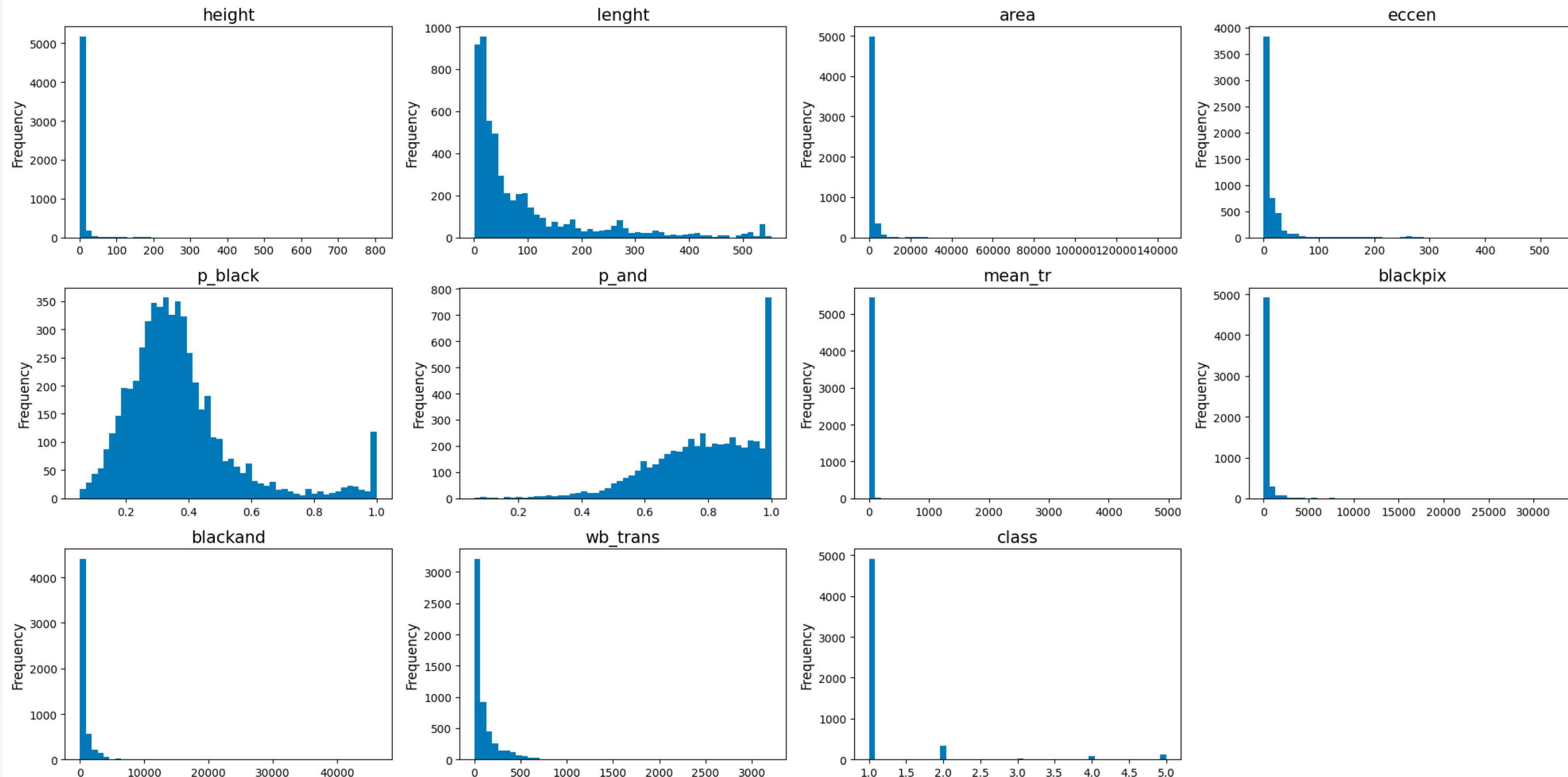
EXPLORE DATASET



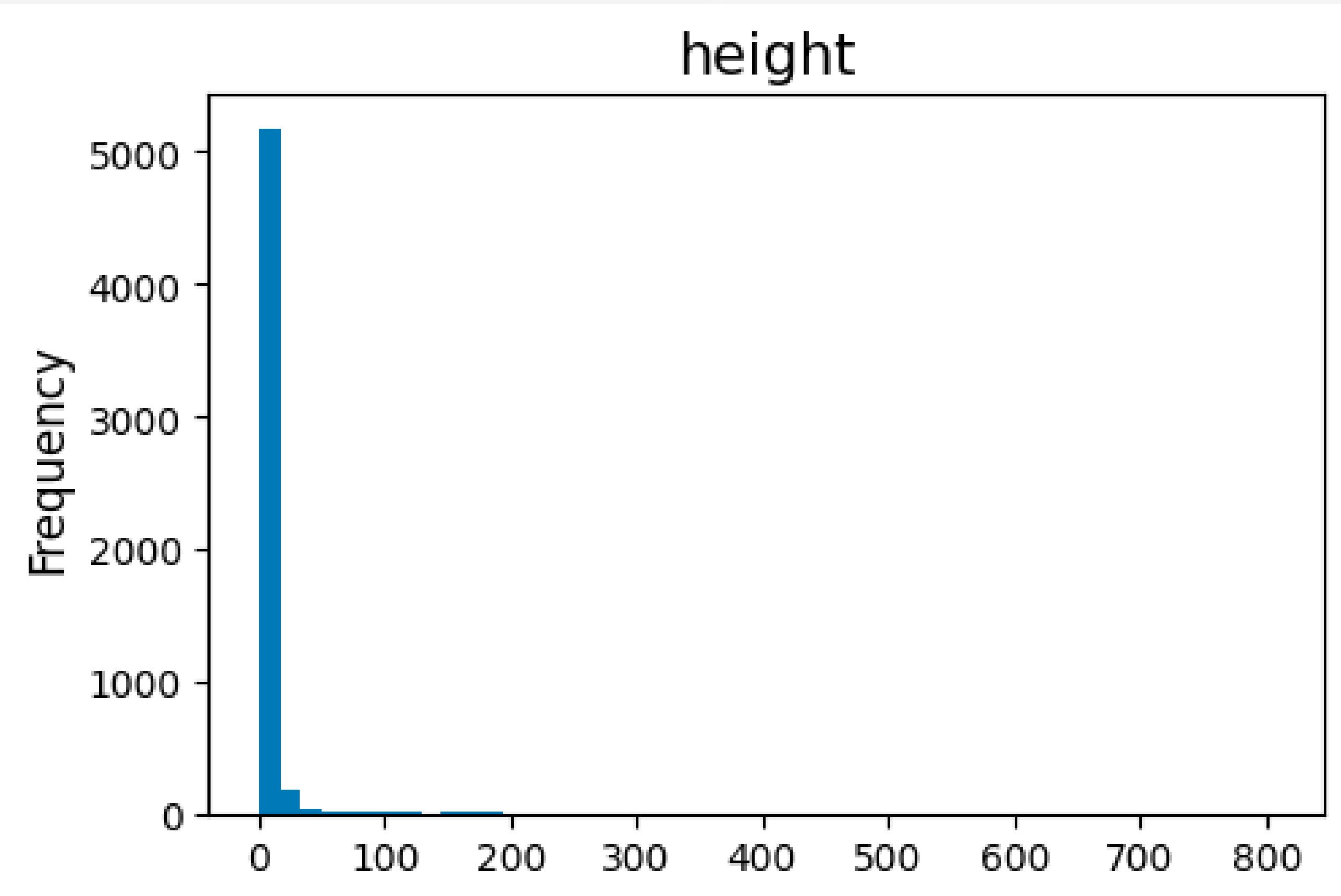
THE DESCRIPTION OF THE DATASET

- Sample size: **5473** samples.
- Number of features: There are **10** feature columns ('height', 'length', 'area', 'eccen', 'p_black', 'p_and', 'mean_tr', 'blackpix', 'blackand', 'wb_trans'), and 1 target column or category column ('class').
- Feature type: Some features are **integer types** ('height', 'length', 'area', 'blackpix', 'blackand', 'wb_trans'), and the rest are **continuous types** ('eccen', 'p_black', 'p_and', 'mean_tr').
- Target variable: "The 'class' column serves as the target variable, containing five classes: text (1), horizontal line (2), picture (3), vertical line (4), and graphic (5)."

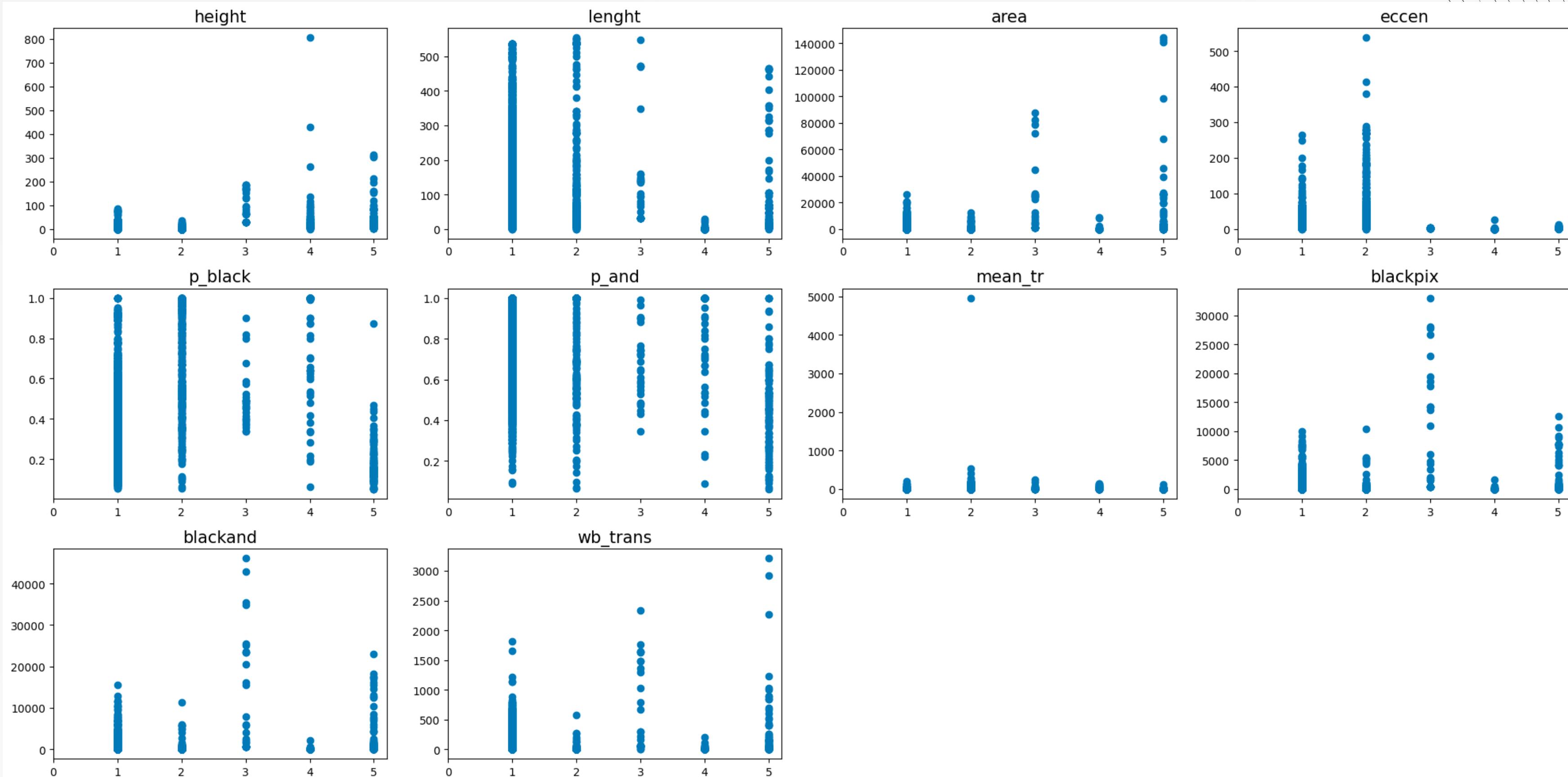
DISTRIBUTION OF ALL ATTRIBUTES



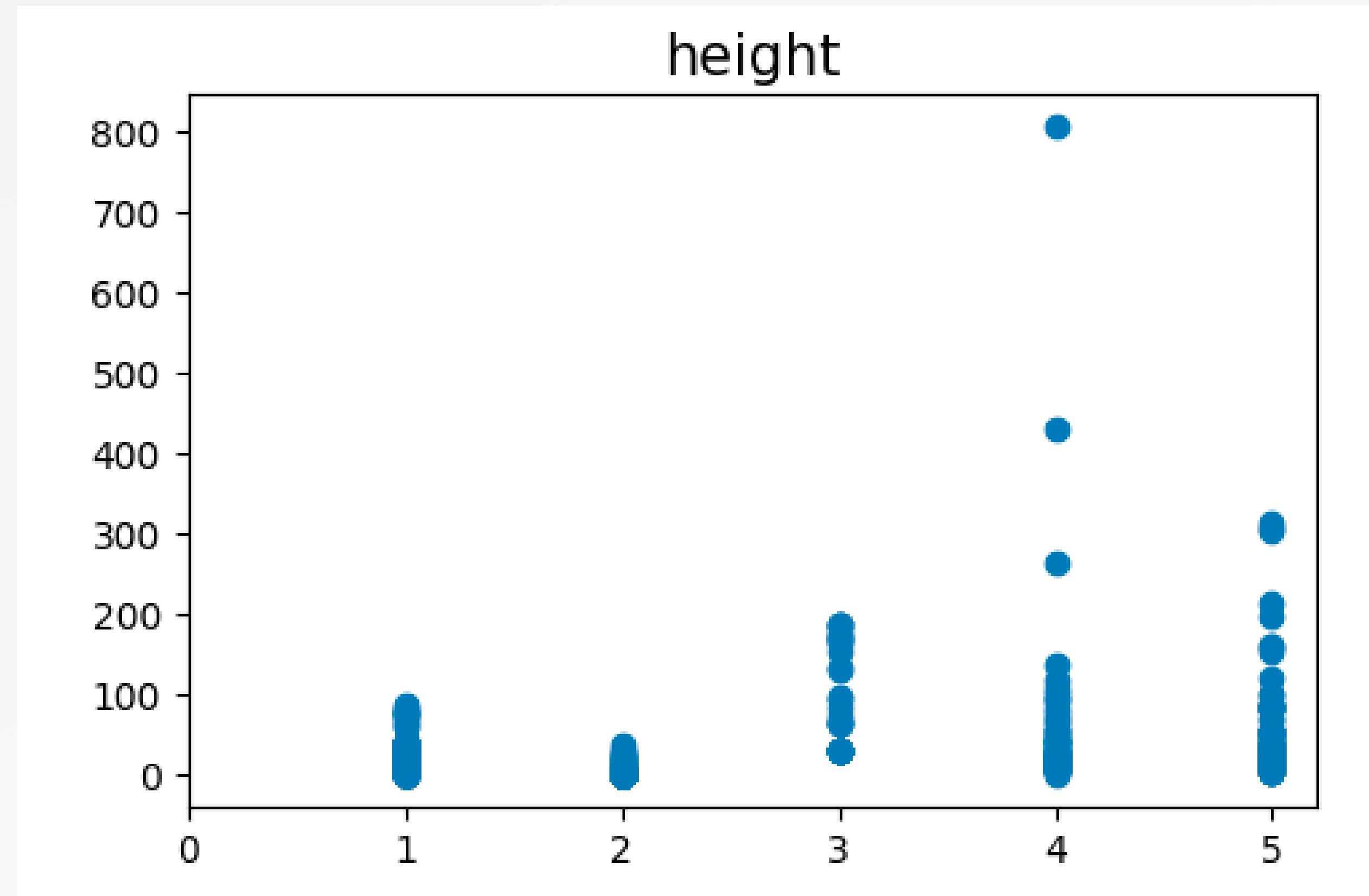
DISTRIBUTION OF ALL ATTRIBUTES



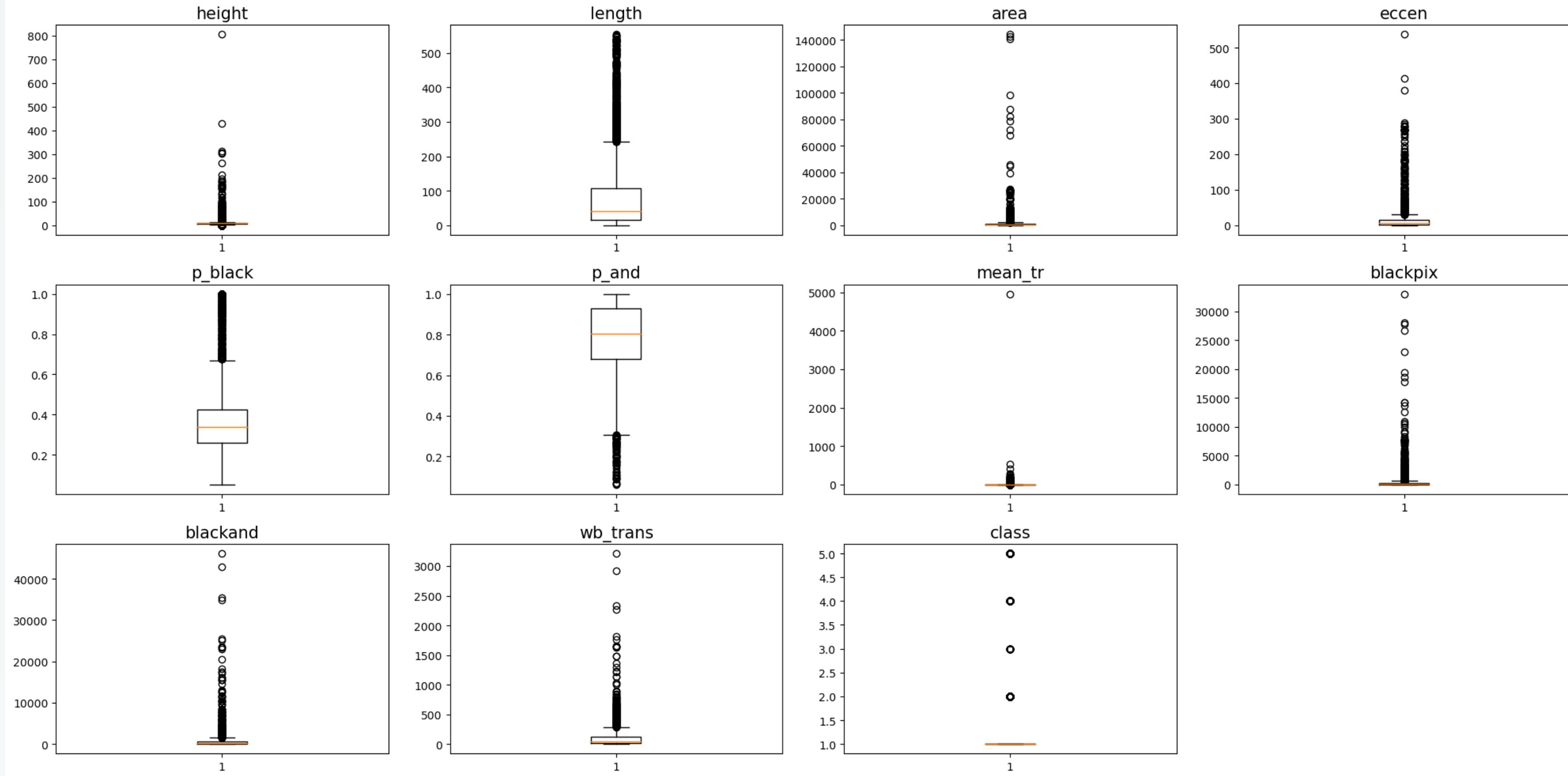
RELATIONSHIP BETWEEN OTHER ATTRIBUTES WITH TARGET ATTRIBUTE



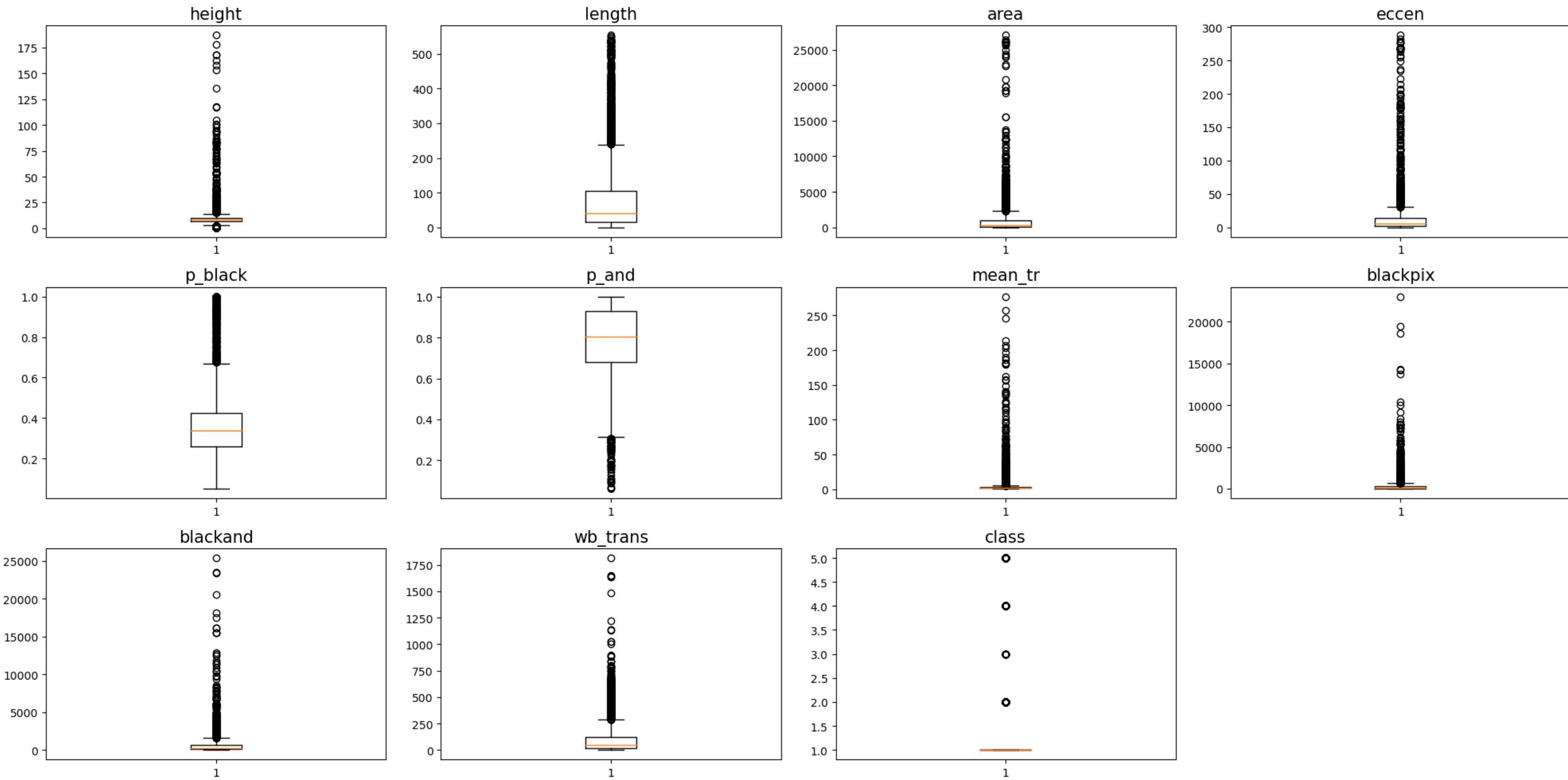
RELATIONSHIP BETWEEN OTHER ATTRIBUTES WITH TARGET ATTRIBUTE



BOXPLOT OF ALL ATTRIBUTES (OUTLIER DETECTION)



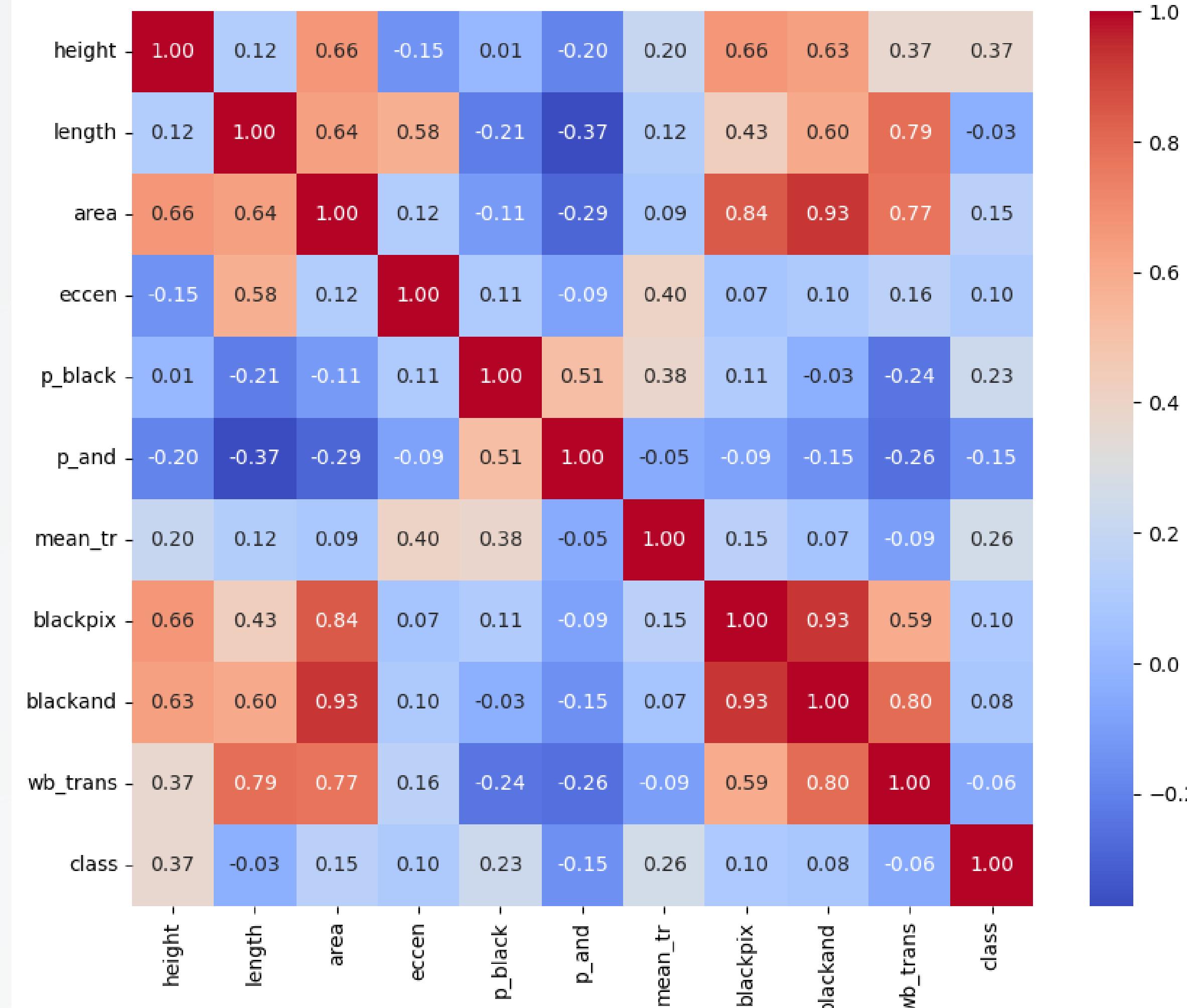
AFTER REMOVING OUTLIERS



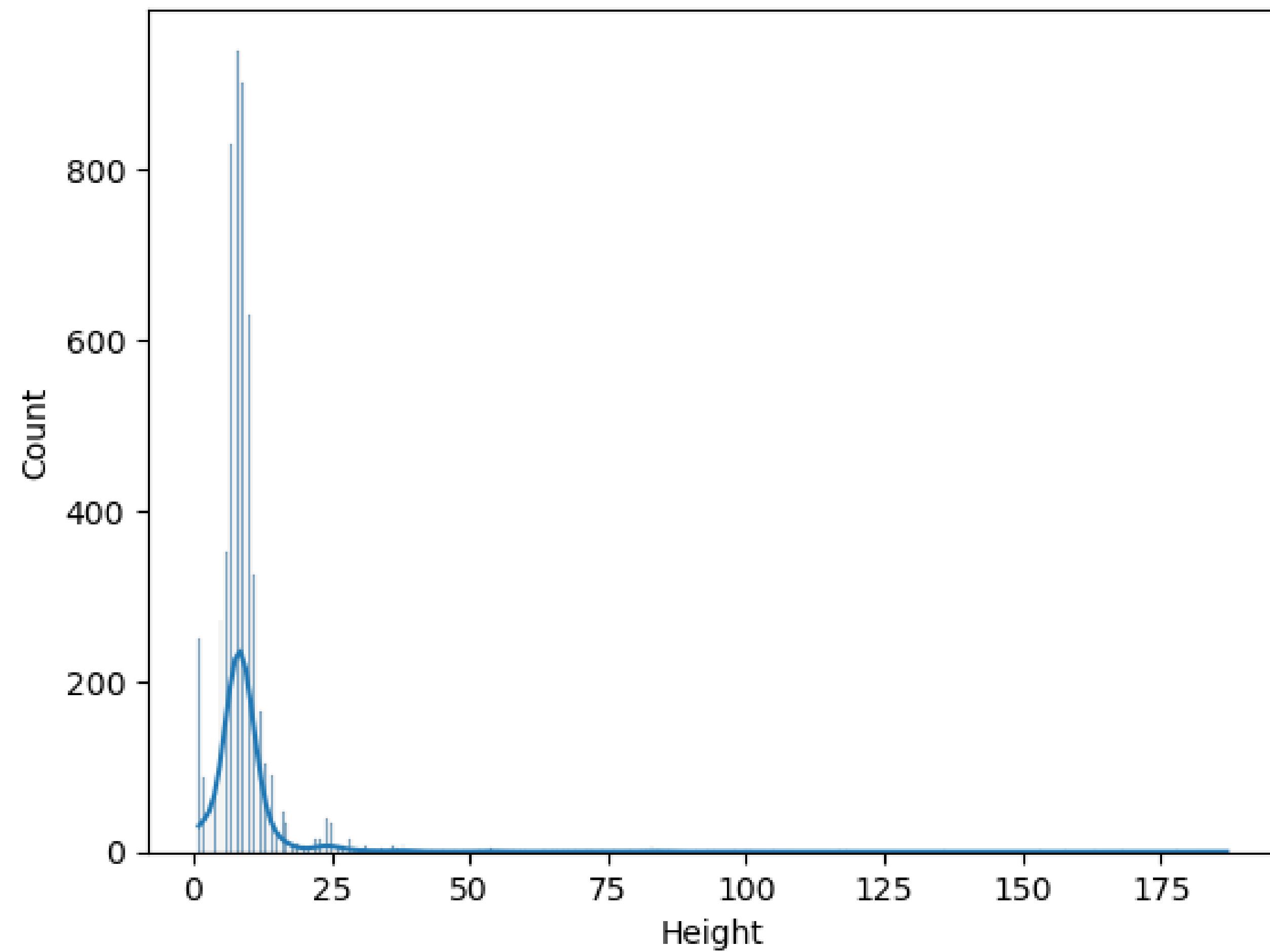
DATASET RESULTS

- Data cleaning: There are no missing values in the data set.
- Outlier Detection: By setting conditions, some extreme values are removed: ep, 'height' < 250, 'area' < 35000, to improve the accuracy of analysis and the generalization ability of the model.
- Conversion and encoding: No conversion or encoding is necessary in this simulation because all data is numeric.
- Normalization: All features except the target variable are standardized, that is, the mean is subtracted and then divided by the standard deviation, so that the mean of the feature becomes 0 and the standard deviation becomes 1.

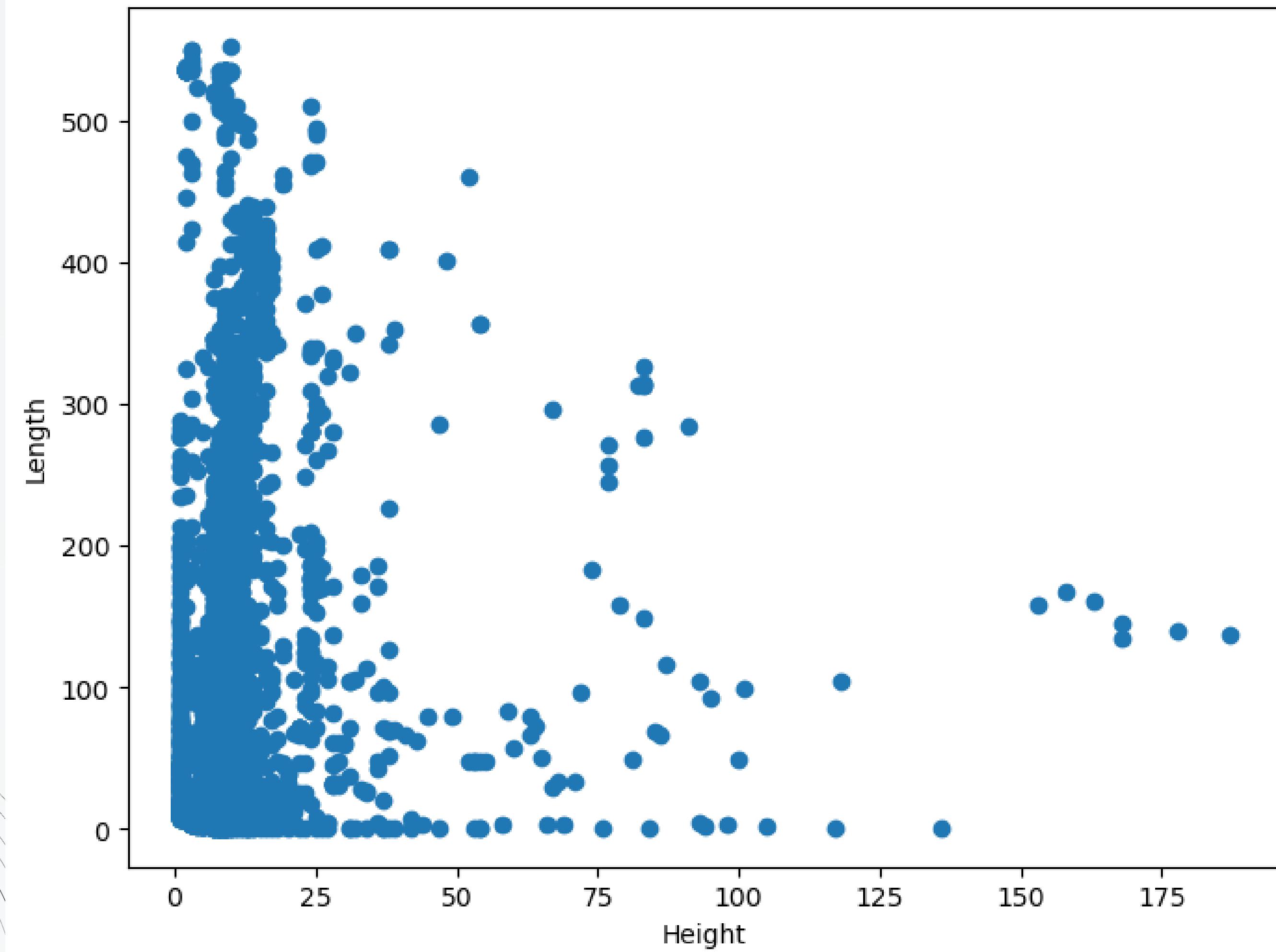
Feature Correlation Matrix



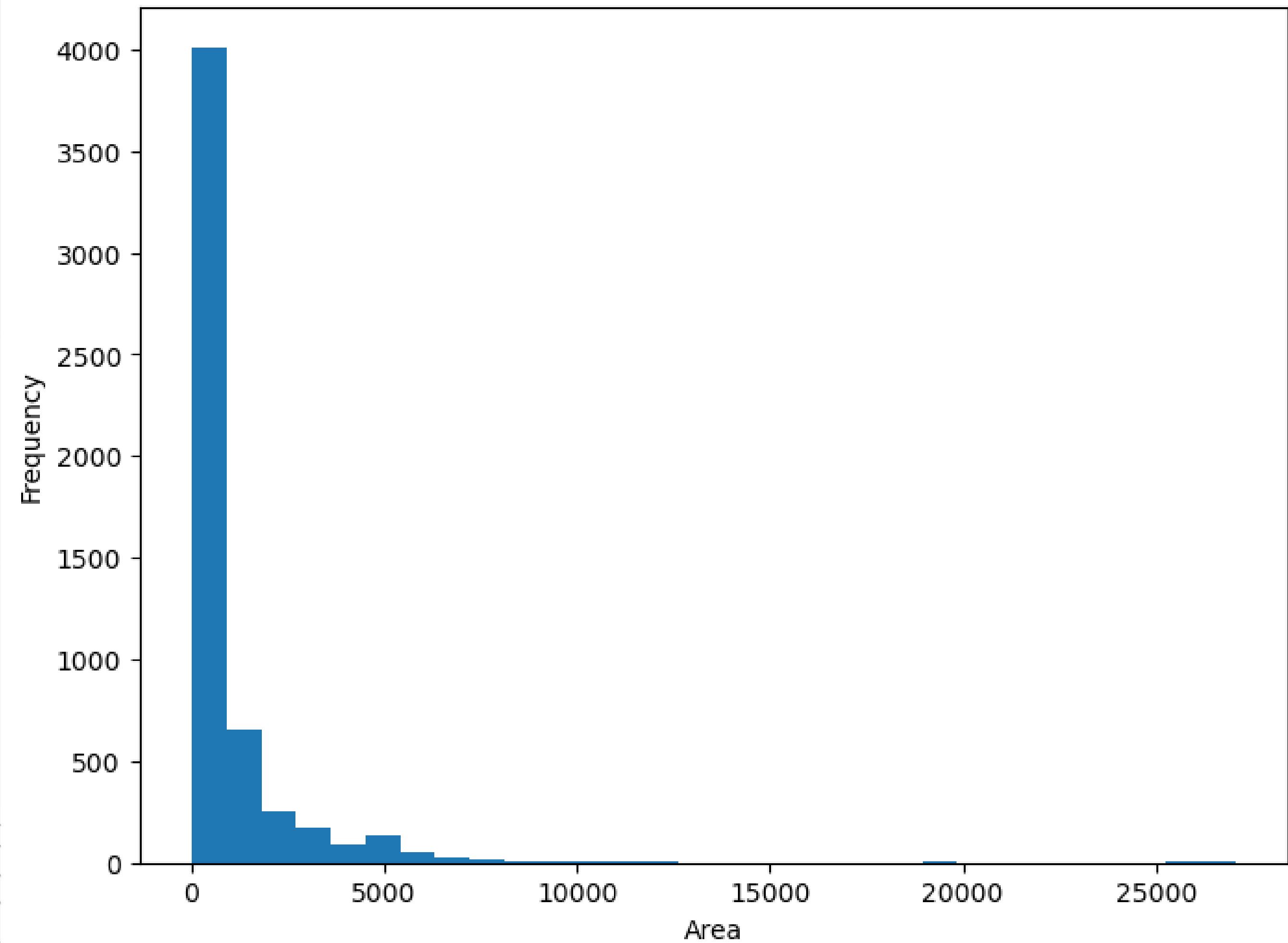
Distribution of Block Heights



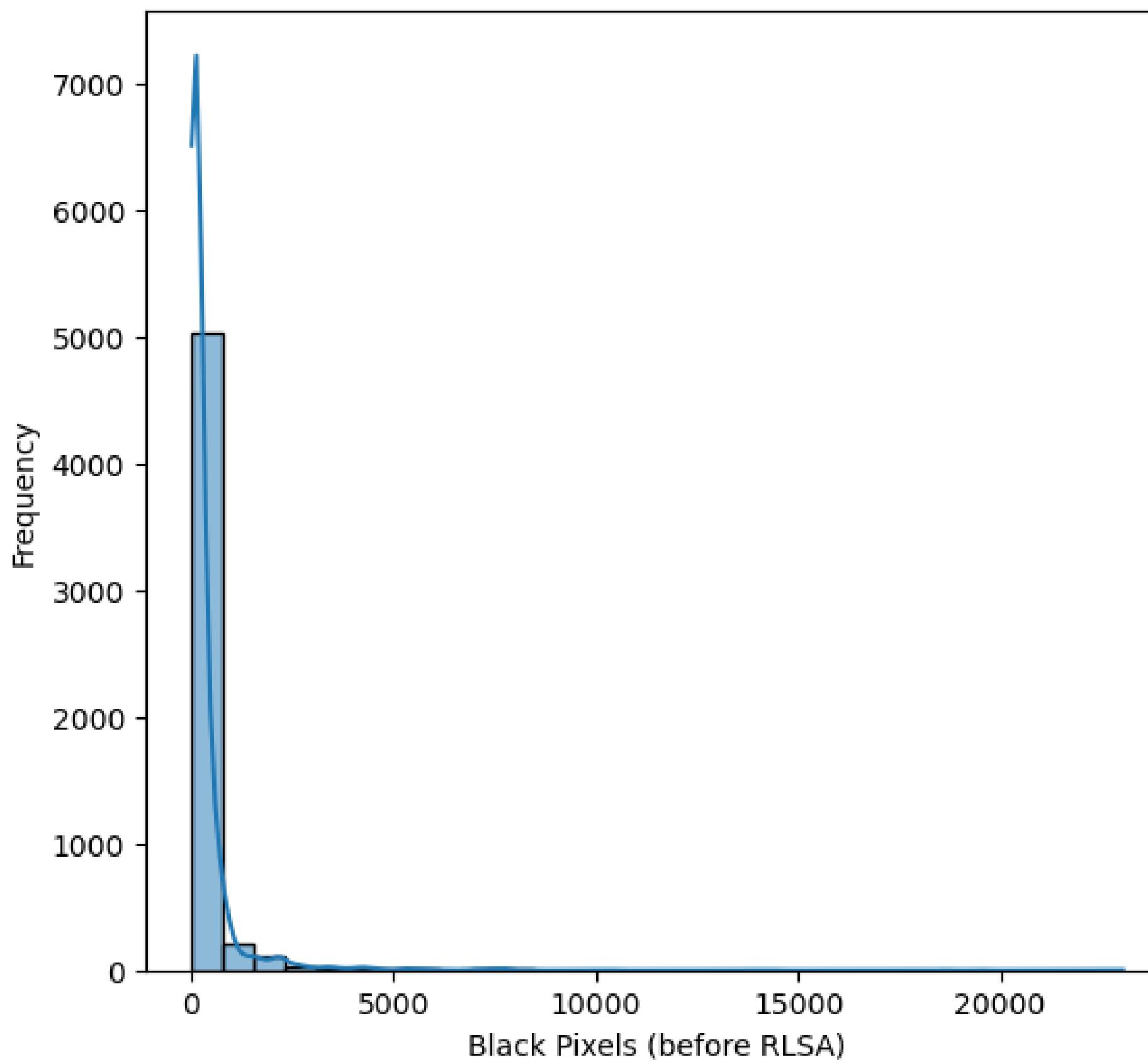
Height vs Length of Blocks



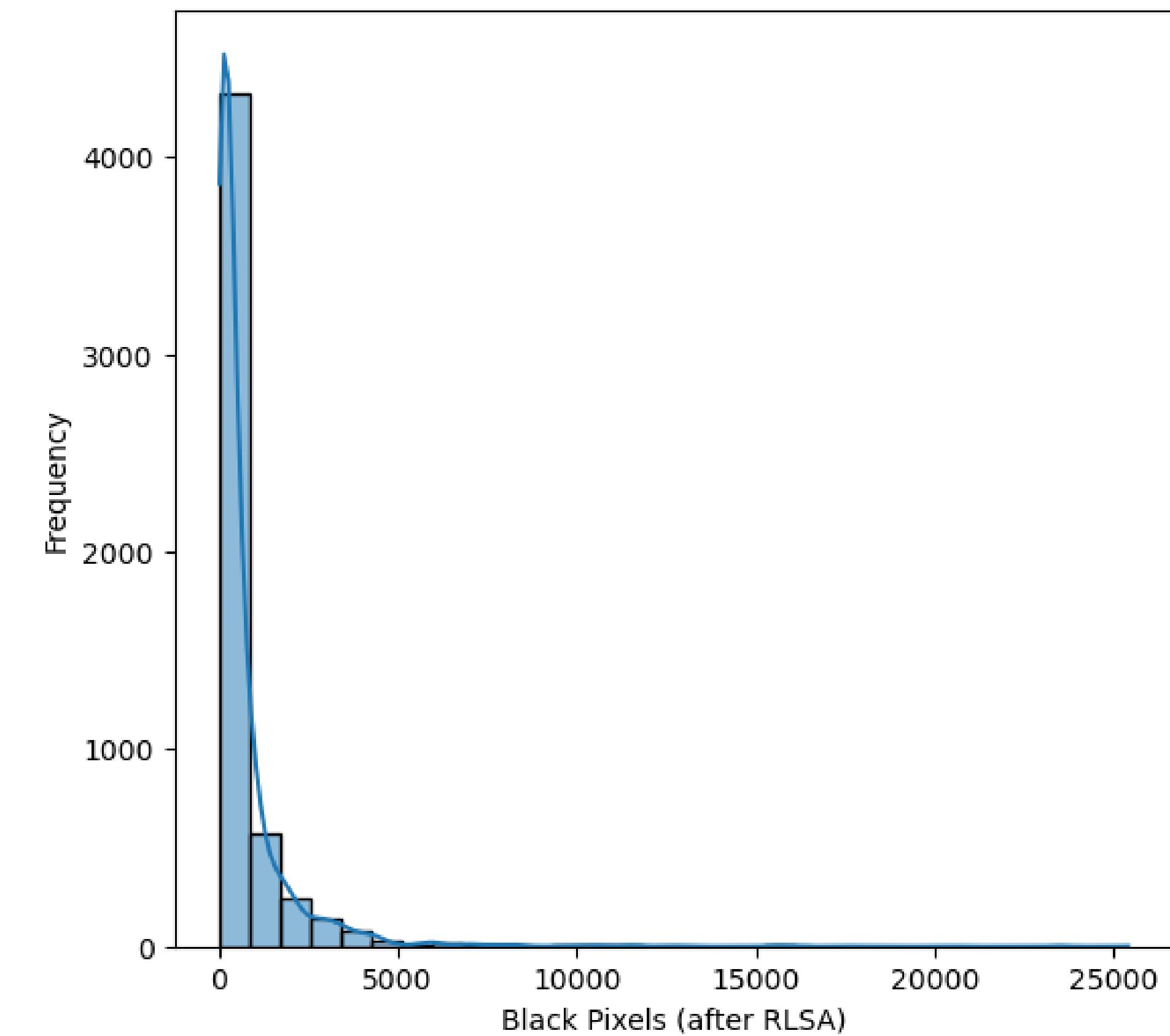
Distribution of Block Area



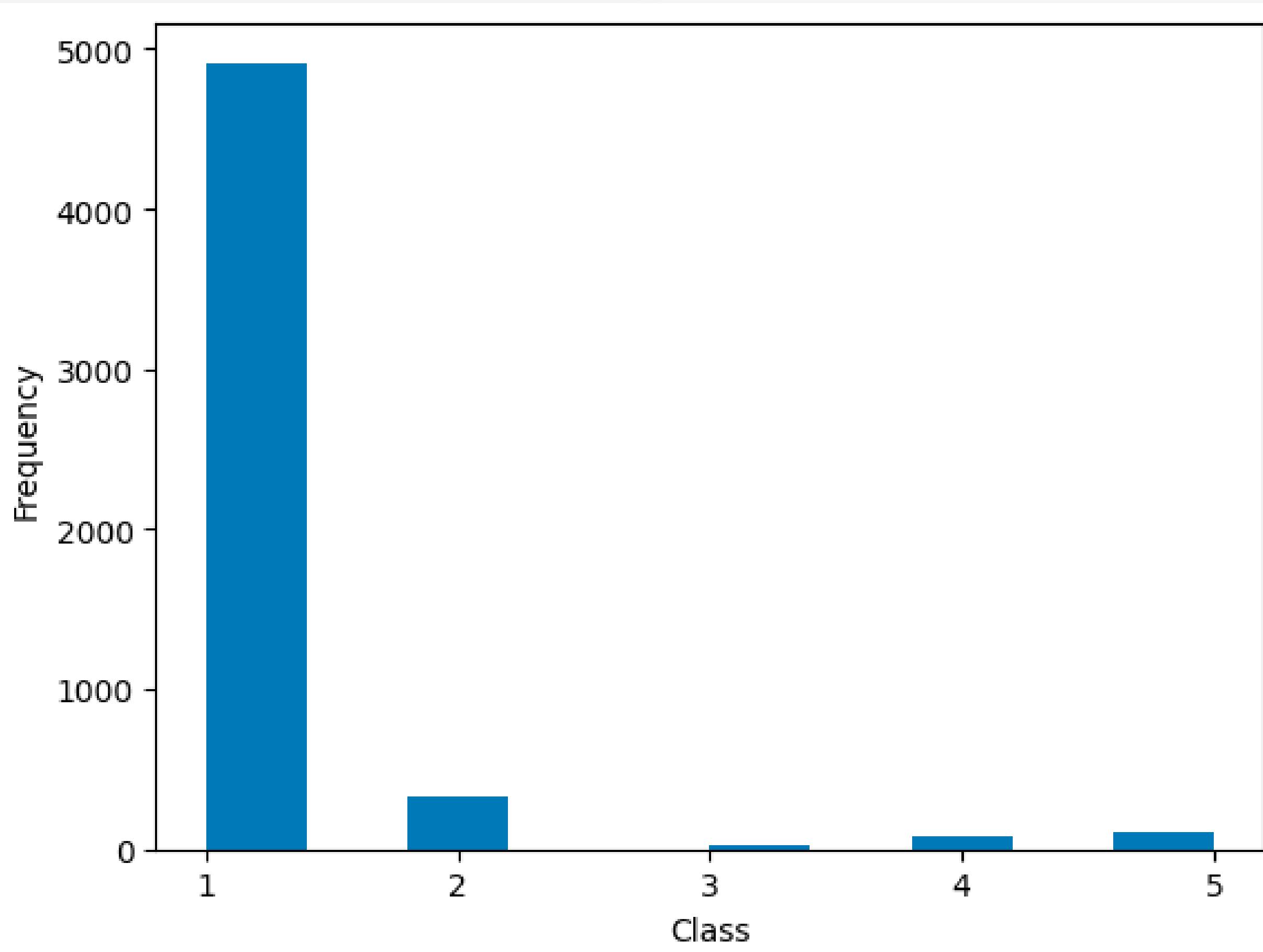
Distribution of Black Pixels Before RLSA



Distribution of Black Pixels After RLSA



VISUALIZATIONS



OBSERVATION OF THE DATASET

- The dataset is highly imbalanced, most of examples belong to class 1.
- Most blocks have a small height, suggests that blocks of smaller height are more common.
- The 'wb_trans' has a high correlation with 'blackpix' and 'blackand'.
- The area of the block has a high correlation with height, length, blackpix, and blackand.
- The blackpix has a high correlation with blackand.

BASELINE MODEL OF LOGISTIC REGRESSION

The dataset involves **continuous and integer** features, and the goal is classification.

For binary classification problems, this is a good starting point because it is simple and easy to implement, while the output probabilities provide a clear classification decision.

RANDOM FOREST MODEL

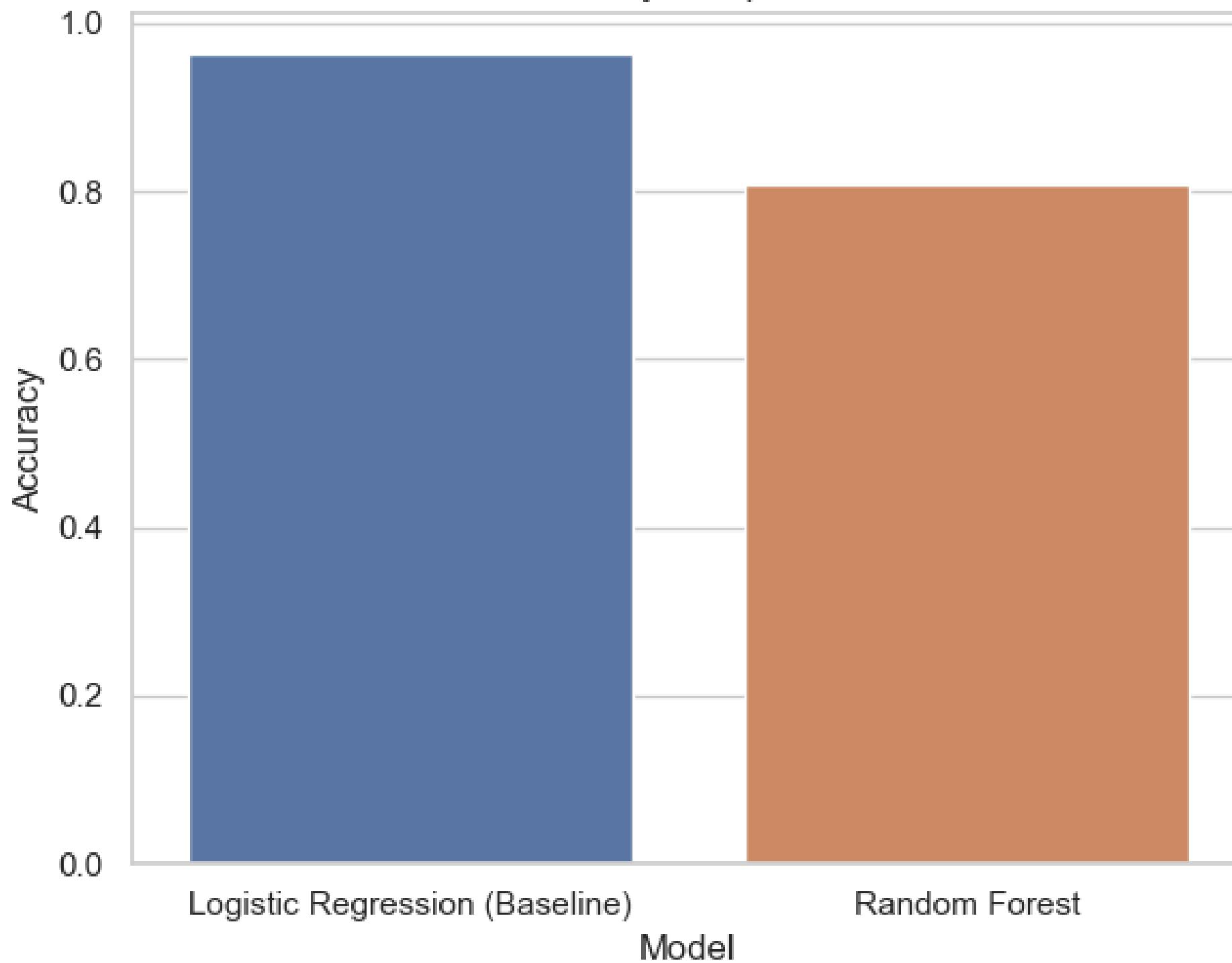
Random Forest is an ensemble learning method that improves the performance and robustness of the overall model by combining multiple decision trees.

- (1) **High performance:** Random forests generally perform well on various types of data sets. It is very effective for processing high-dimensional, non-linear and complex relationship data sets.
- (2) **Feature importance:** Random forests provide an intuitive way to evaluate feature importance. By observing how each feature is used across multiple decision trees, we can determine which features are critical to the model's predictions.
- (3) **Handle imbalanced data:** Random forests can handle imbalanced data sets and can balance predictions between classes through a voting mechanism.

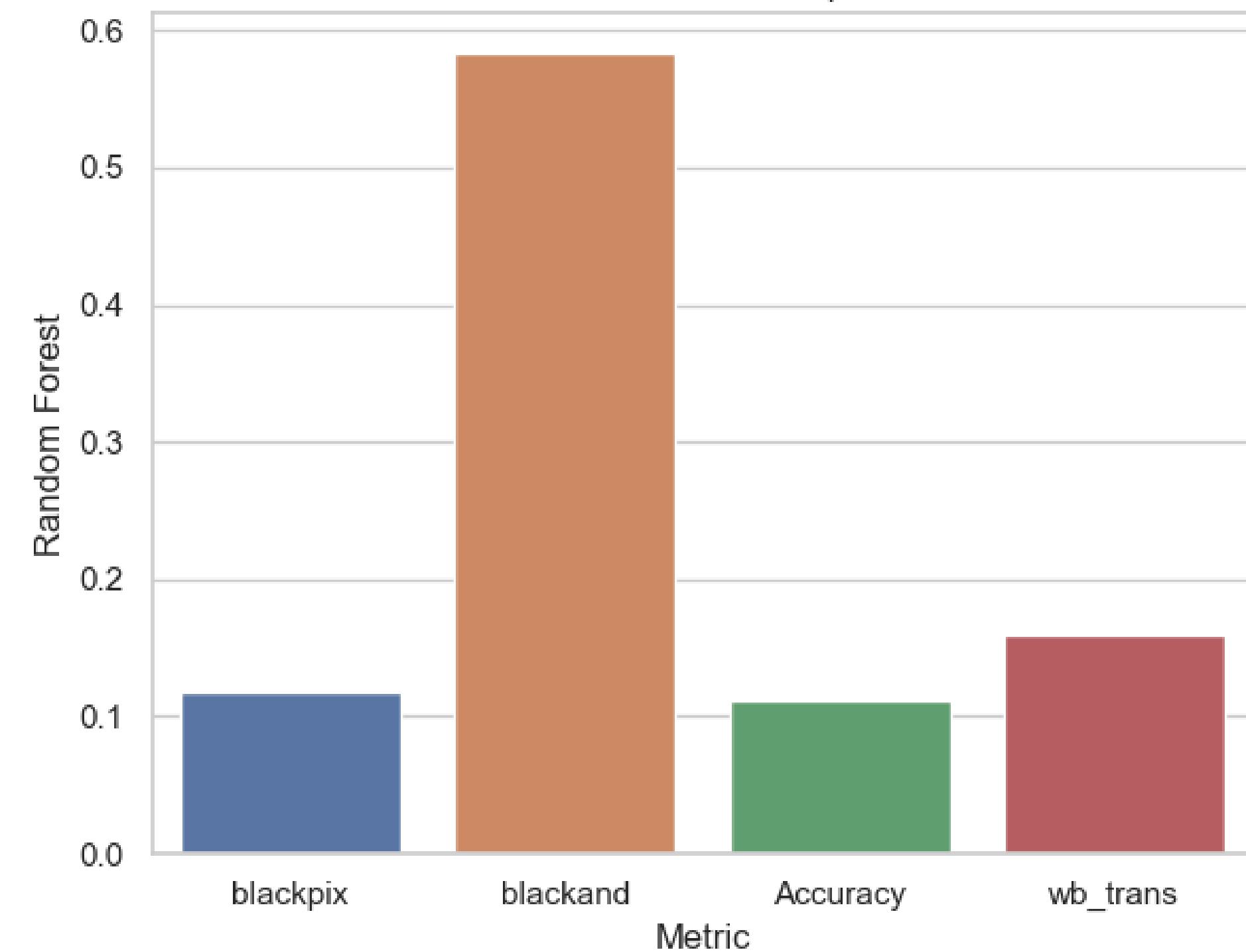
(1) PREDICTING A SINGLE TARGET (REGRESSION PROBLEM)

```
Accuracy: 0.8065994500458296
Accuracy (Random Forest): 0.8065994500458296
Feature Importance:
eccen: 0.2413
area: 0.1598
length: 0.0962
blackpix: 0.0905
blackand: 0.0879
mean_tr: 0.0876
p_black: 0.0828
p_and: 0.0791
wb_trans: 0.0579
class: 0.0168
```

Accuracy Comparison



Additional Metrics Comparison



(2) MULTI-OBJECTIVE REGRESSION PROBLEM

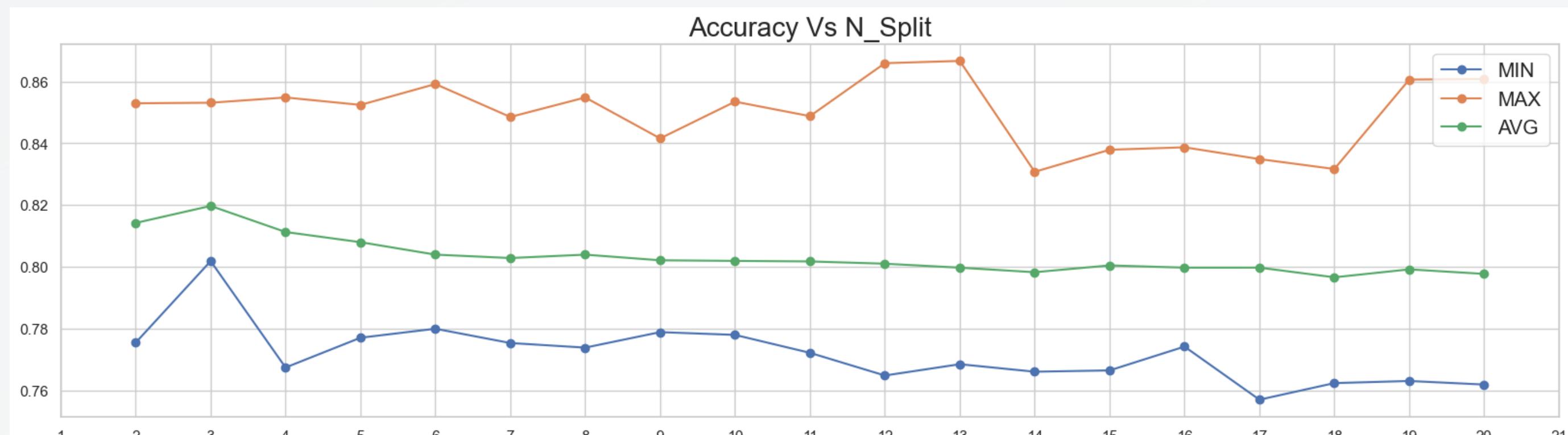
Mean Squared Error: 82160.83161707914

not good...

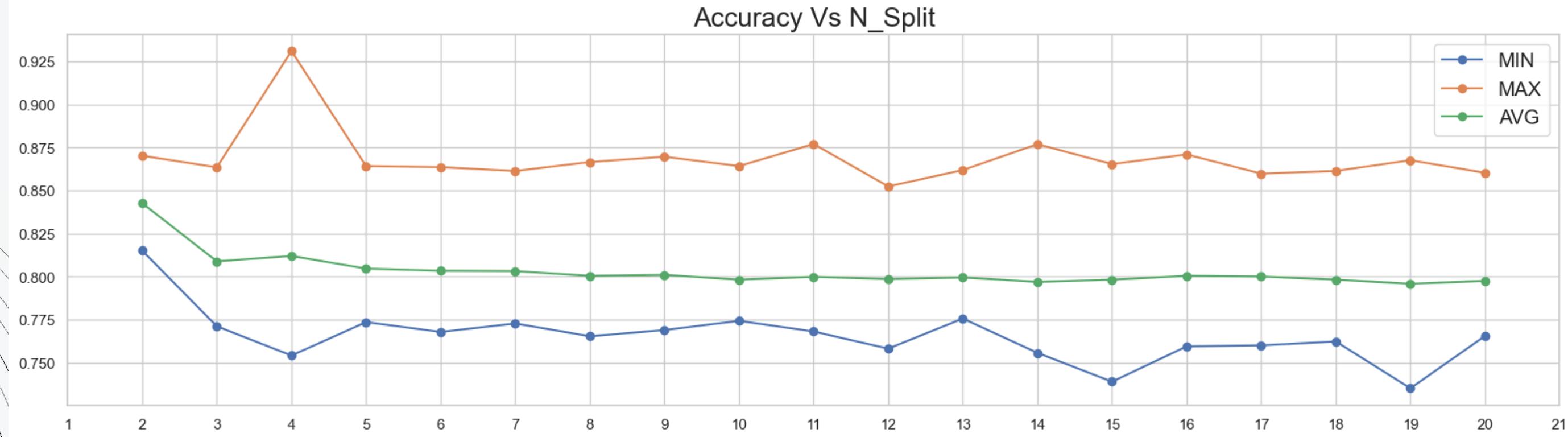
The background features two sets of abstract black line art. On the left, a series of thin lines radiate from a central point, creating a fan-like or spiral effect that tapers towards the top left. On the right, a more complex, dense cluster of lines forms a large, rounded, organic shape that tapers towards the bottom right.

MODEL EVALUATION

CALCULATE THE CROSS VALIDATION SCORE ON DIFFERENT SPLIT POINT IN K-FOLD



K-FOLD



STRATIFIED K-FOLD

RESULTS AND DISCUSSION

- Using GaussianNB, the model achieves an average accuracy of approximately 88%.
- In the 5-fold and 10-fold cross-validation, the interval between the minimum and maximum accuracy values is not large, indicating that the model has a certain stability.
- Try to create new features or transform existing features to improve the performance of the model. Handle imbalanced data sets through oversampling or undersampling.

REFERENCE

- Dataset:

<https://archive.ics.uci.edu/dataset/78/page+blocks+classification>

- Article:

1. Fushiki, Tadayoshi. "Estimation of prediction error by using K-fold cross-validation." *Statistics and Computing* 21 (2011): 137-146.
2. Anguita, Davide, et al. "The 'K' in K-fold Cross Validation." *ESANN*. 2012.

**THANK'S FOR
WATCHING**

