



# **PAGE BLOCKS CLASSIFICATION**

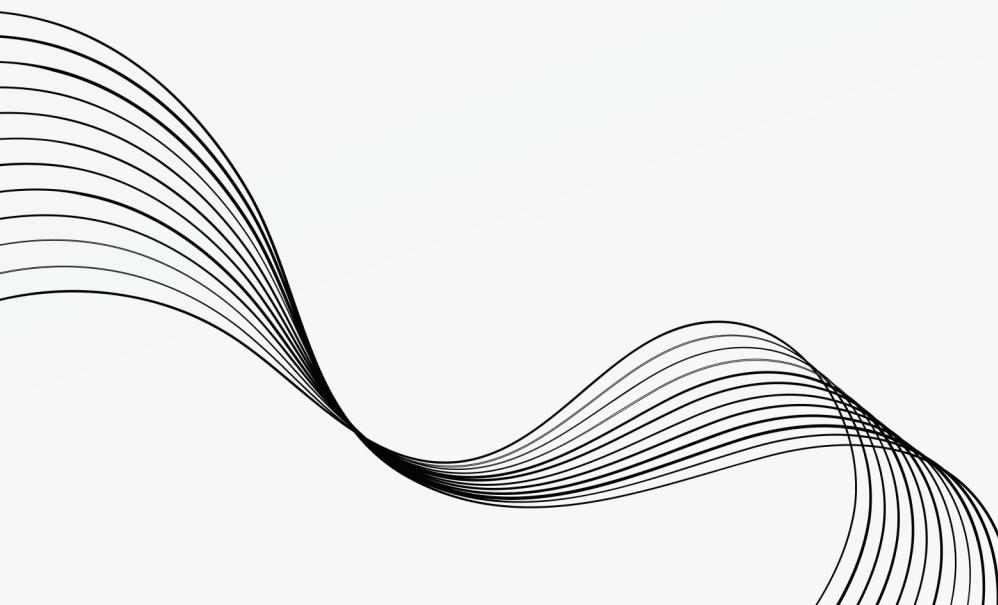
**Presented by Shangzhi LOU, Wenbo  
SUI and Limin TIAN**

# CONTENT

- 01** DEFINITION OF THE PROBLEM
- 02** RESEARCH METHODOLOGY
- 03** EXPLORATORY DATA ANALYSIS
- 04** MODEL BUILDING AND EVALUATION
- 05** RESULTS
- 06** DISCUSSION

# DEFINITION OF THE PROBLEM

The problem consists of classifying all the blocks of the page layout of a document that has been detected by a segmentation process. The goal is to classify each block into one of five classes.

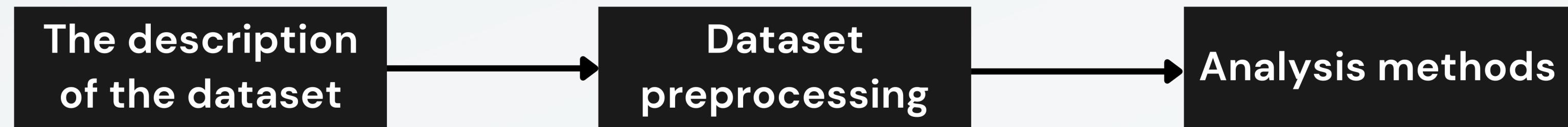


# DEFINITION OF THE PROBLEM

## The different variables

- **height:** Height of the block.
- **length:** Length of the block.
- **area:** Area of the block (**height \* length**);
- **eccen:** Eccentricity of the block (**length / height**);
- **p\_black:** Percentage of black pixels within the block (**blackpix / area**);
- **p\_and:** Percentage of black pixels after the application of the Run Length Smoothing Algorithm (RLSA) (**blackand / area**);
- **mean\_tr:** Mean number of white-black transitions (**blackpix / wb\_trans**);
- **blackpix:** Total number of black pixels in the original bitmap of the block.
- **blackand:** Total number of black pixels in the bitmap of the block after the RLSA.
- **wb\_trans:** Number of white-black transitions in the original bitmap of the block.

# RESEARCH METHODOLOGY



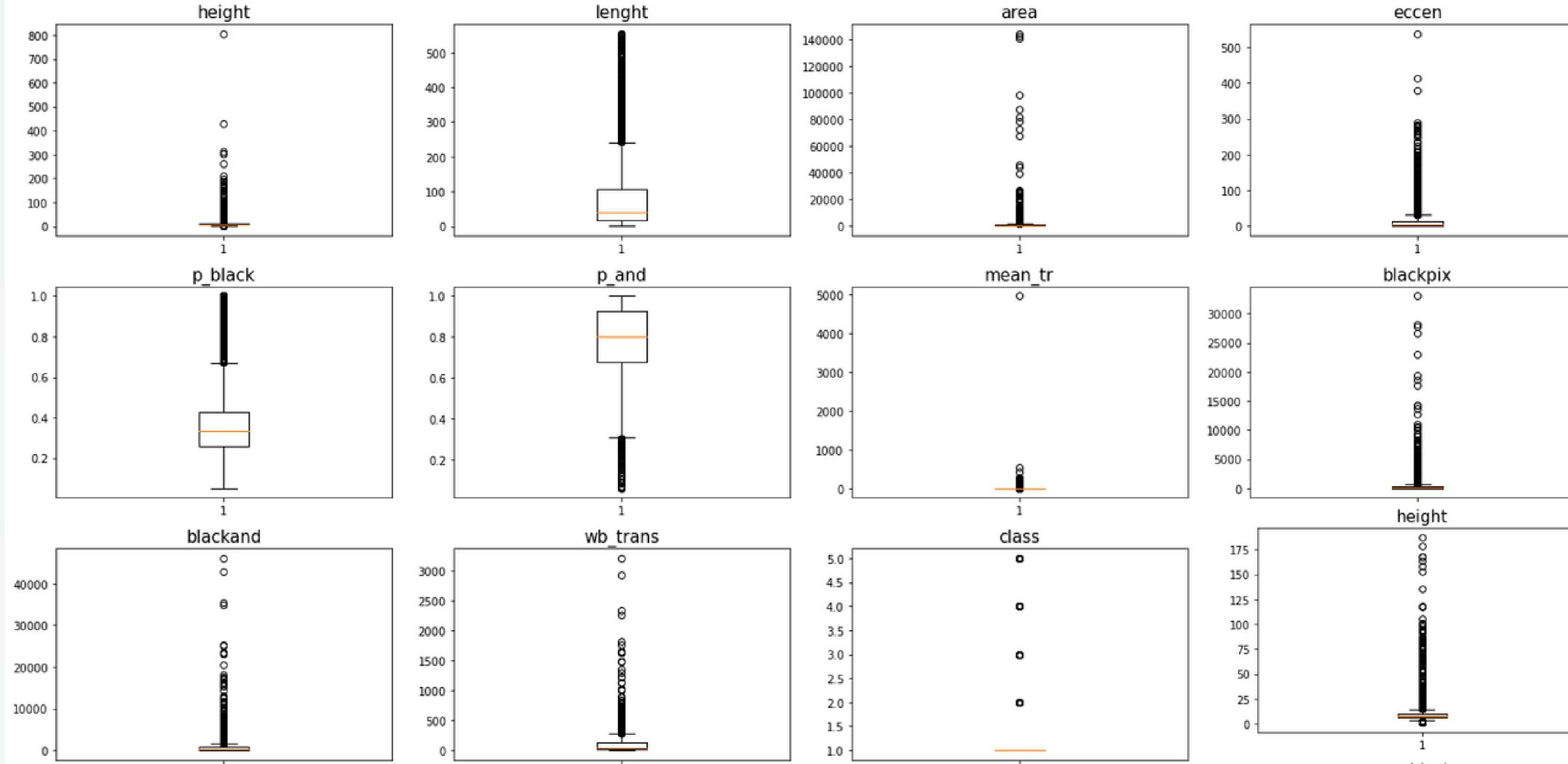
- 5473 examples, 5 classes.
- All attributes are numeric.
- The data is formatted in a way that is readable by C4.5.

- Check the missing values
- Outlier Detection
- Normalization ( z-score )

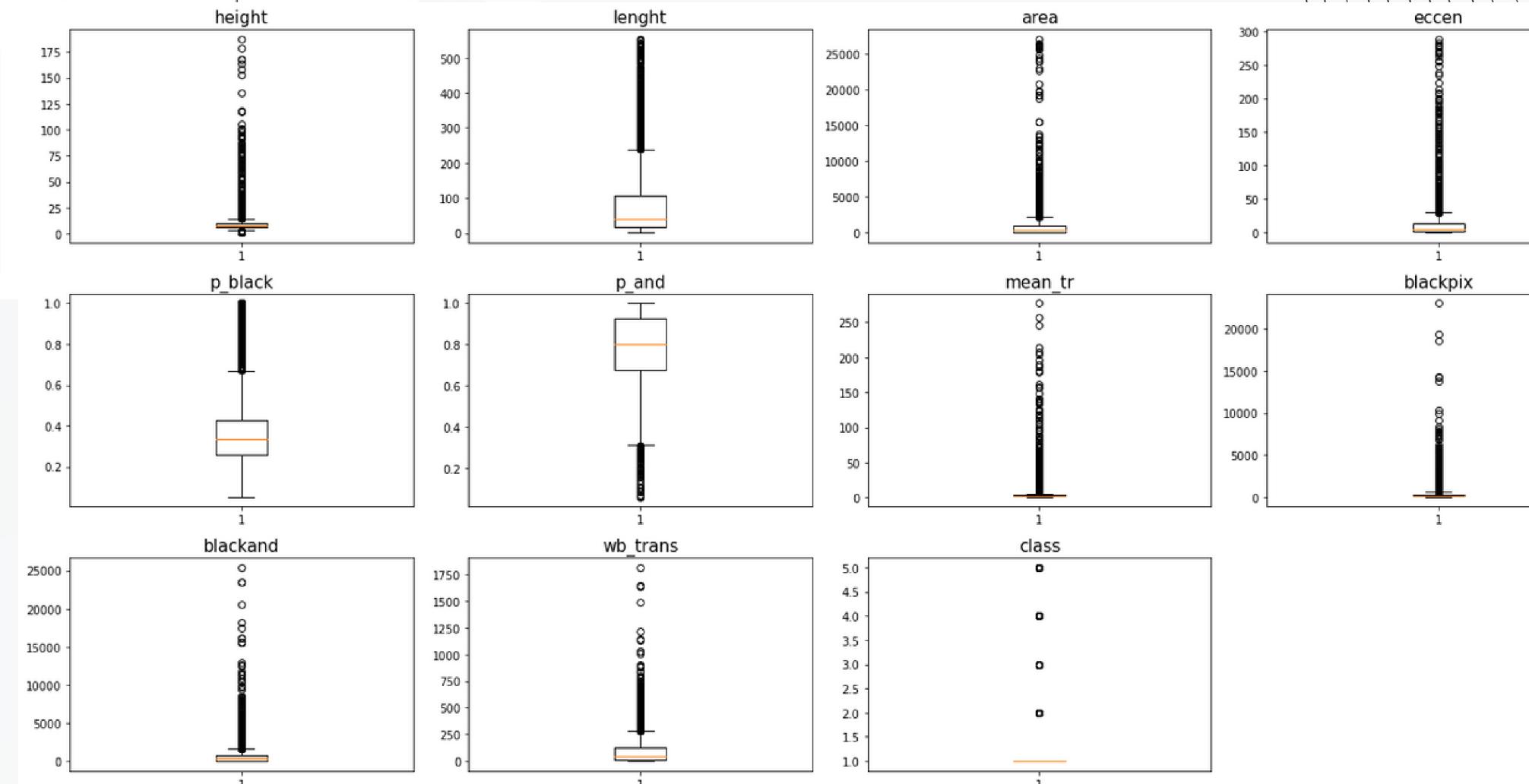
- GaussianNB
- StratifiedKFold

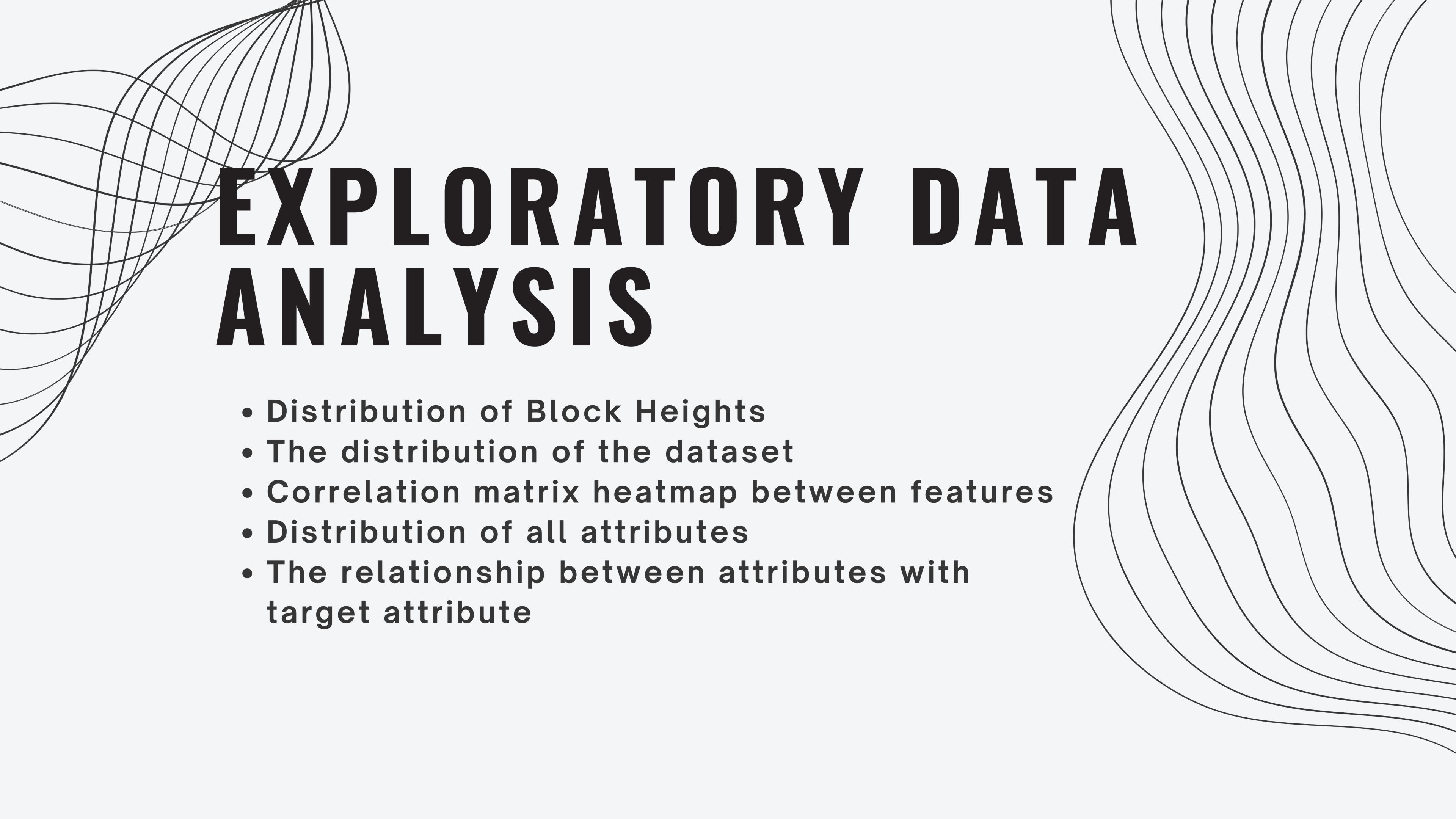


# VISUALIZATIONS



## Outlier Detection

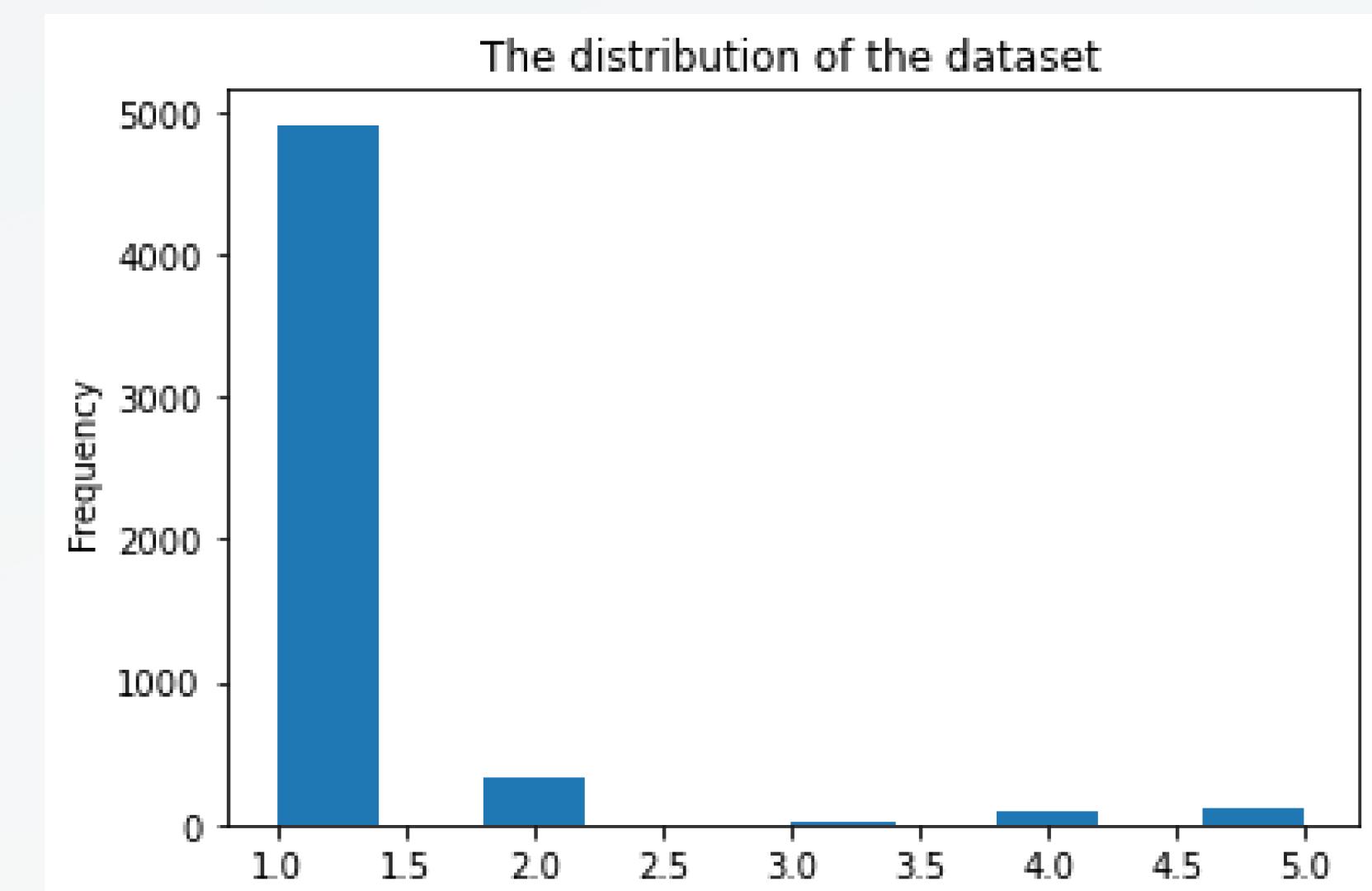
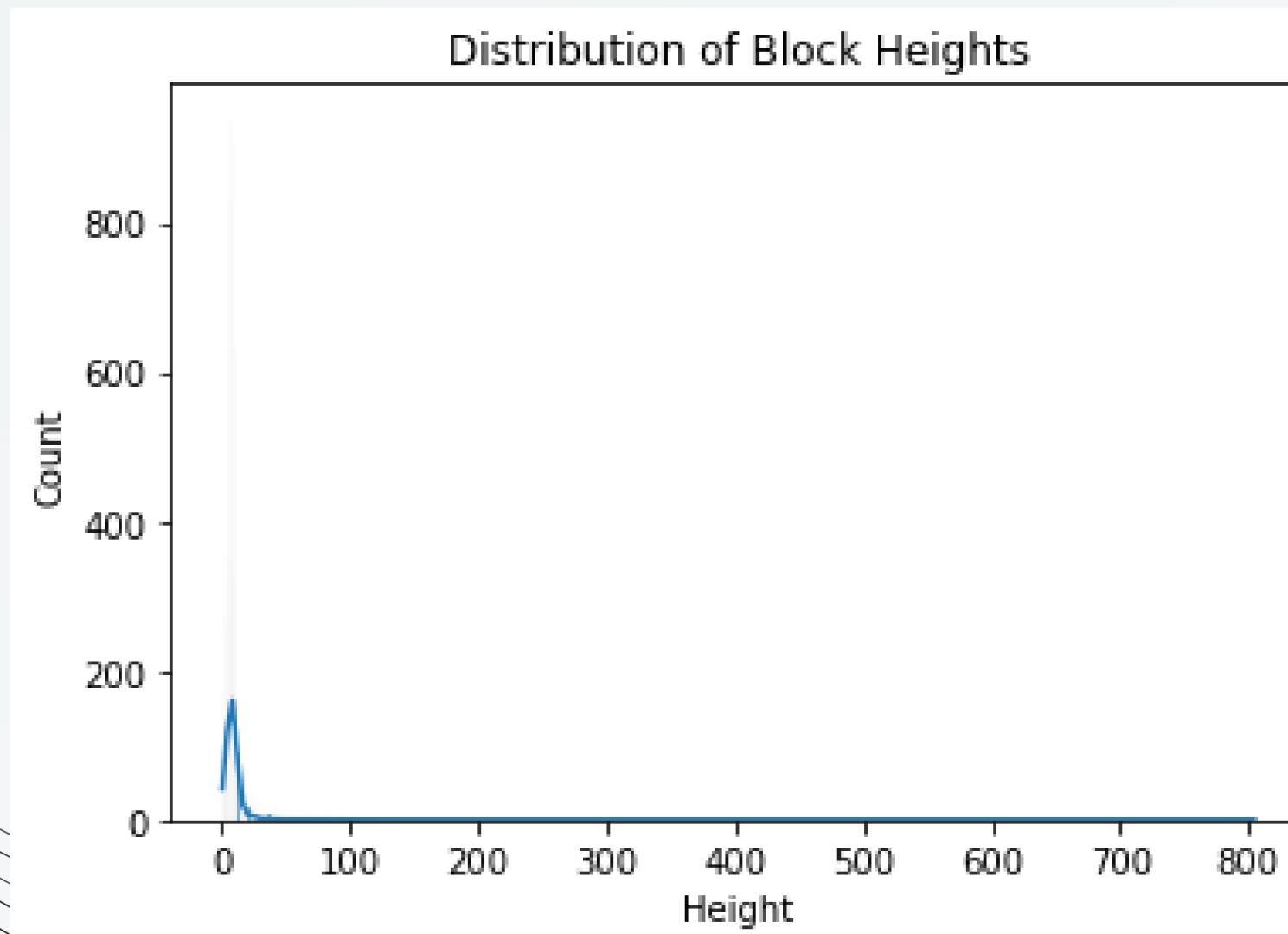




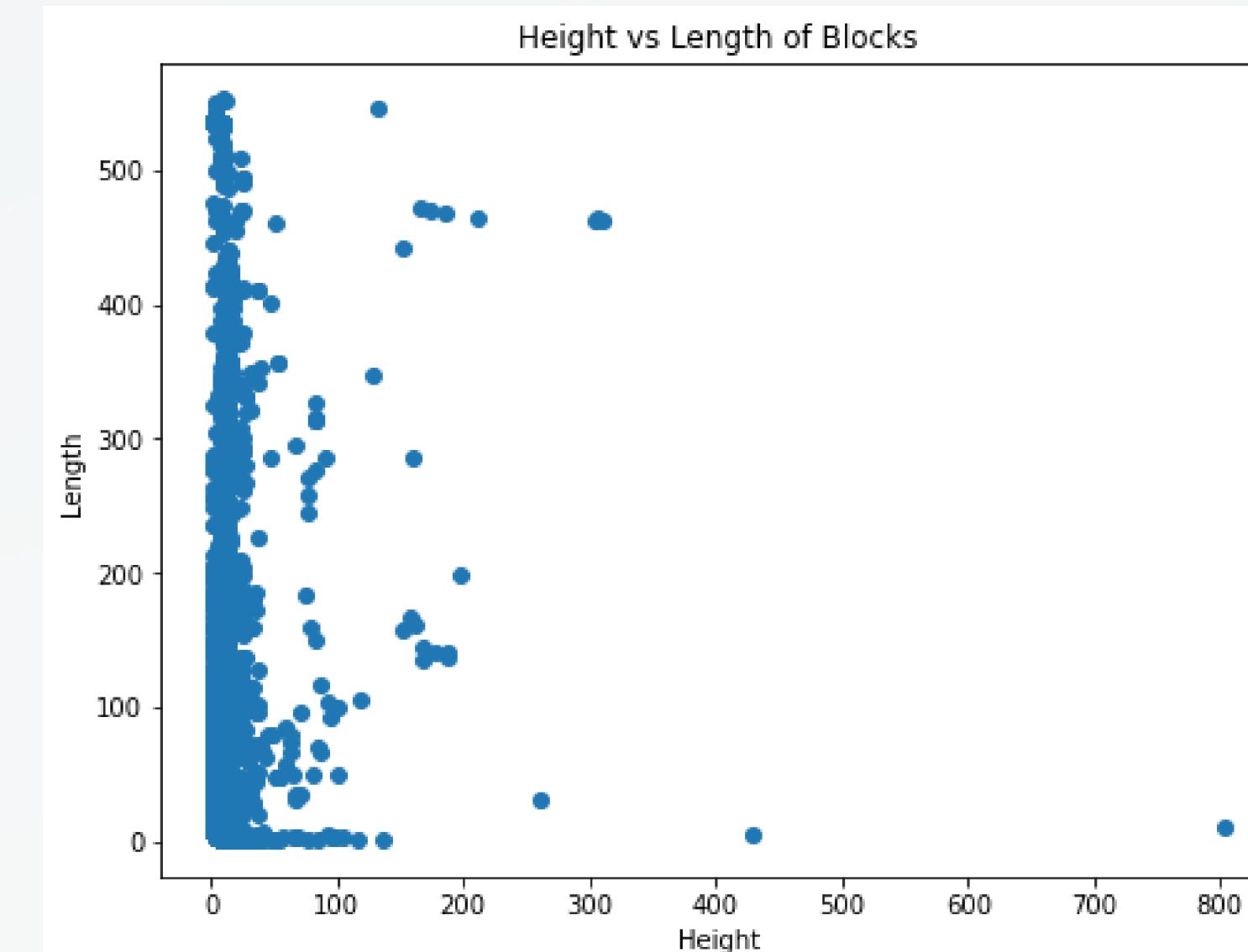
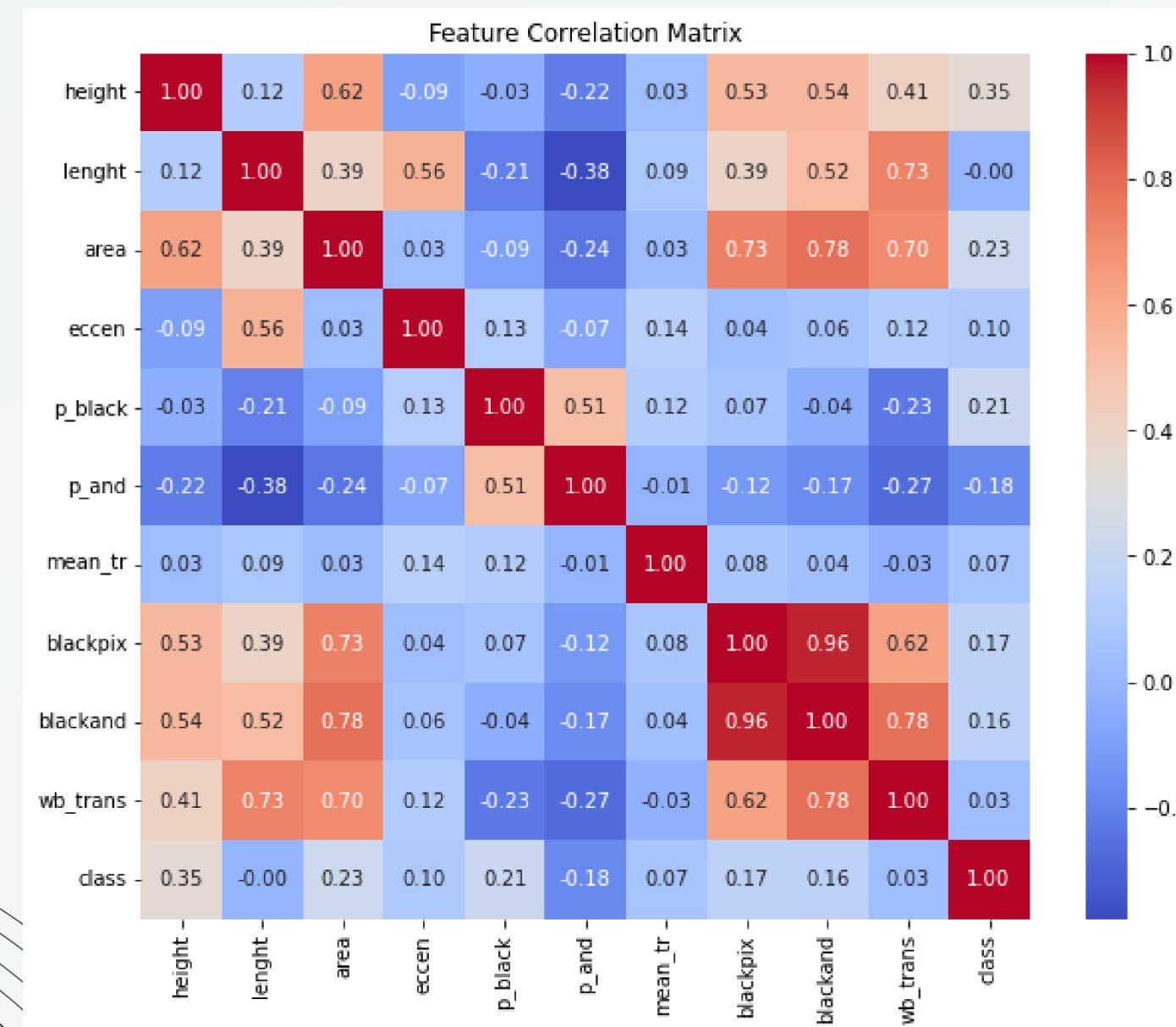
# **EXPLORATORY DATA ANALYSIS**

- Distribution of Block Heights
- The distribution of the dataset
- Correlation matrix heatmap between features
- Distribution of all attributes
- The relationship between attributes with target attribute

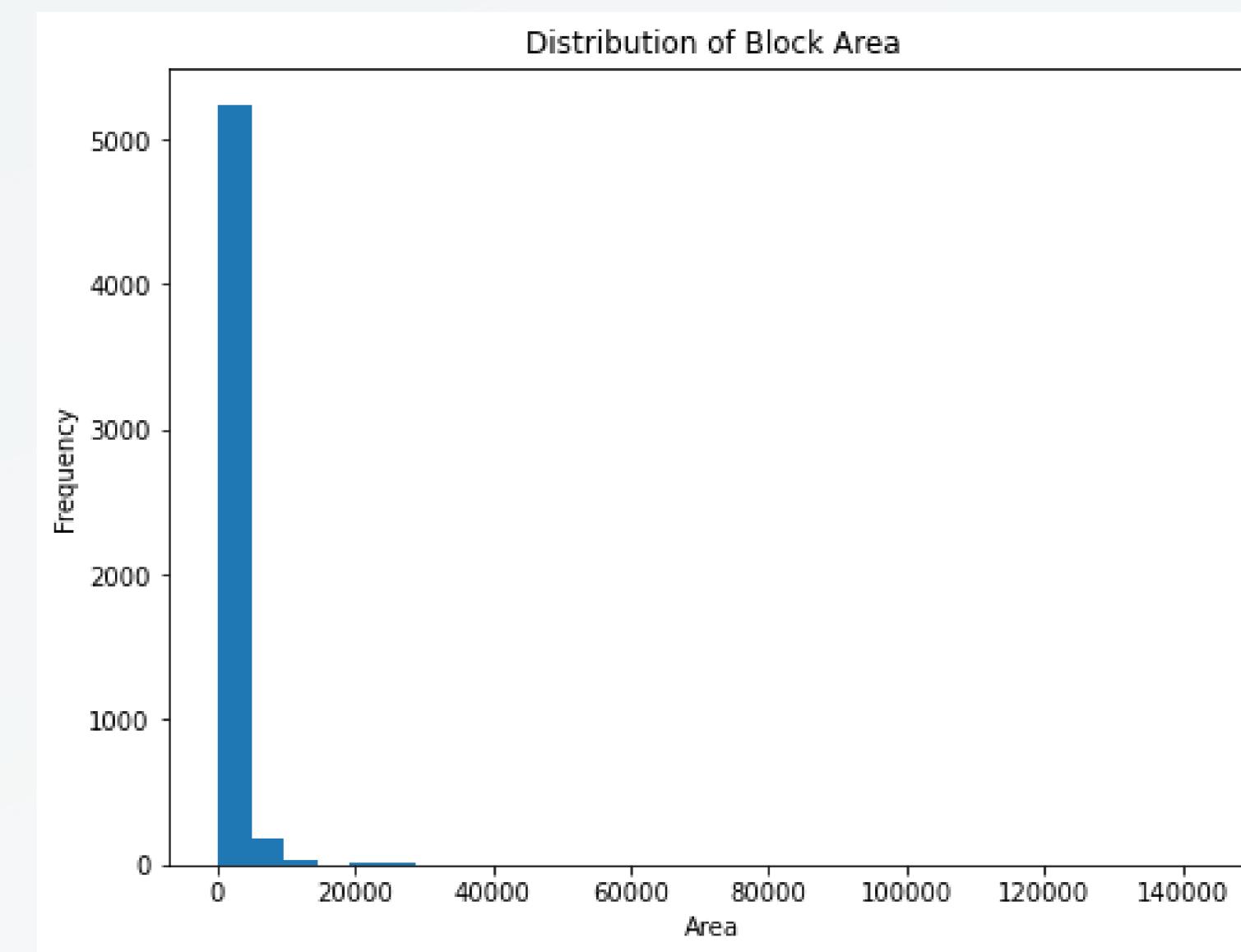
# VISUALIZATIONS



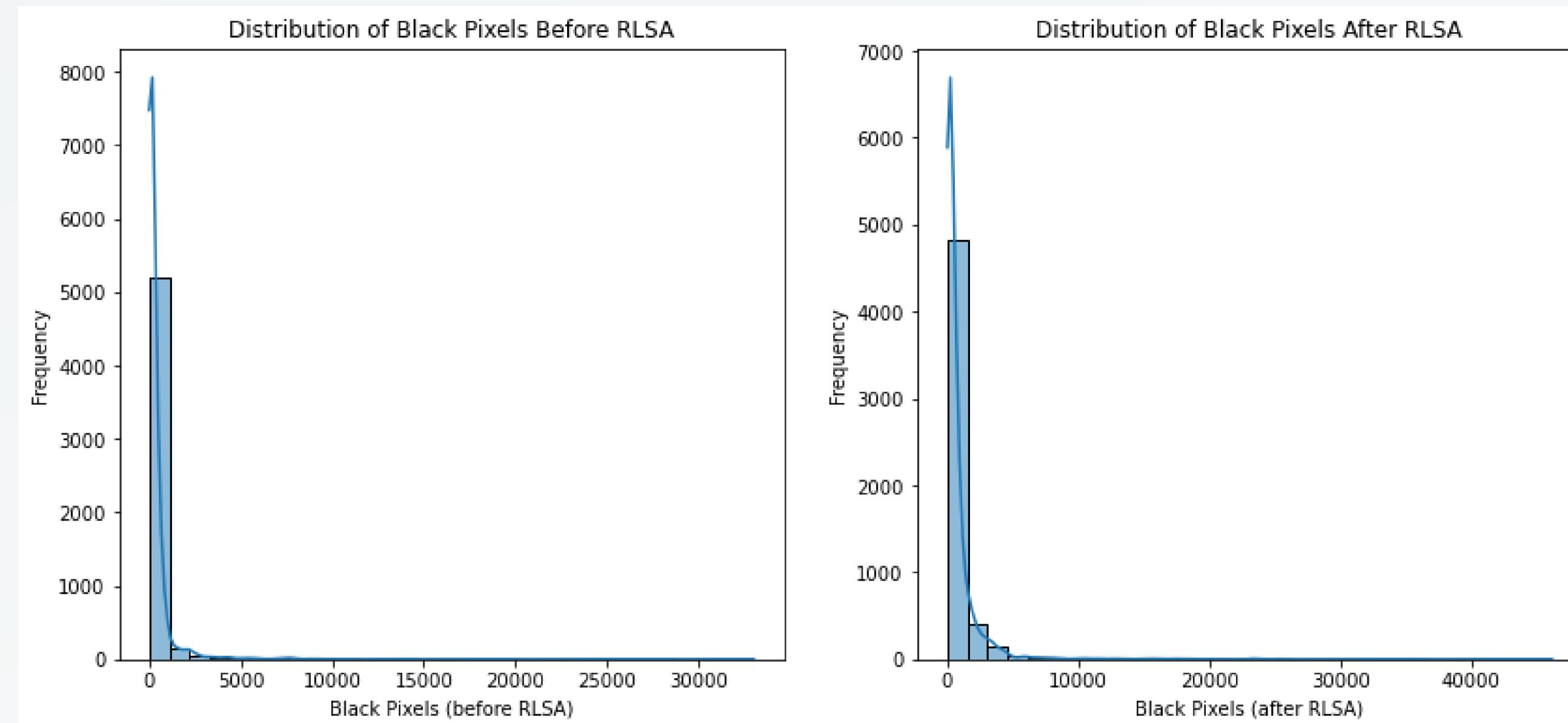
# VISUALIZATIONS



# VISUALIZATIONS

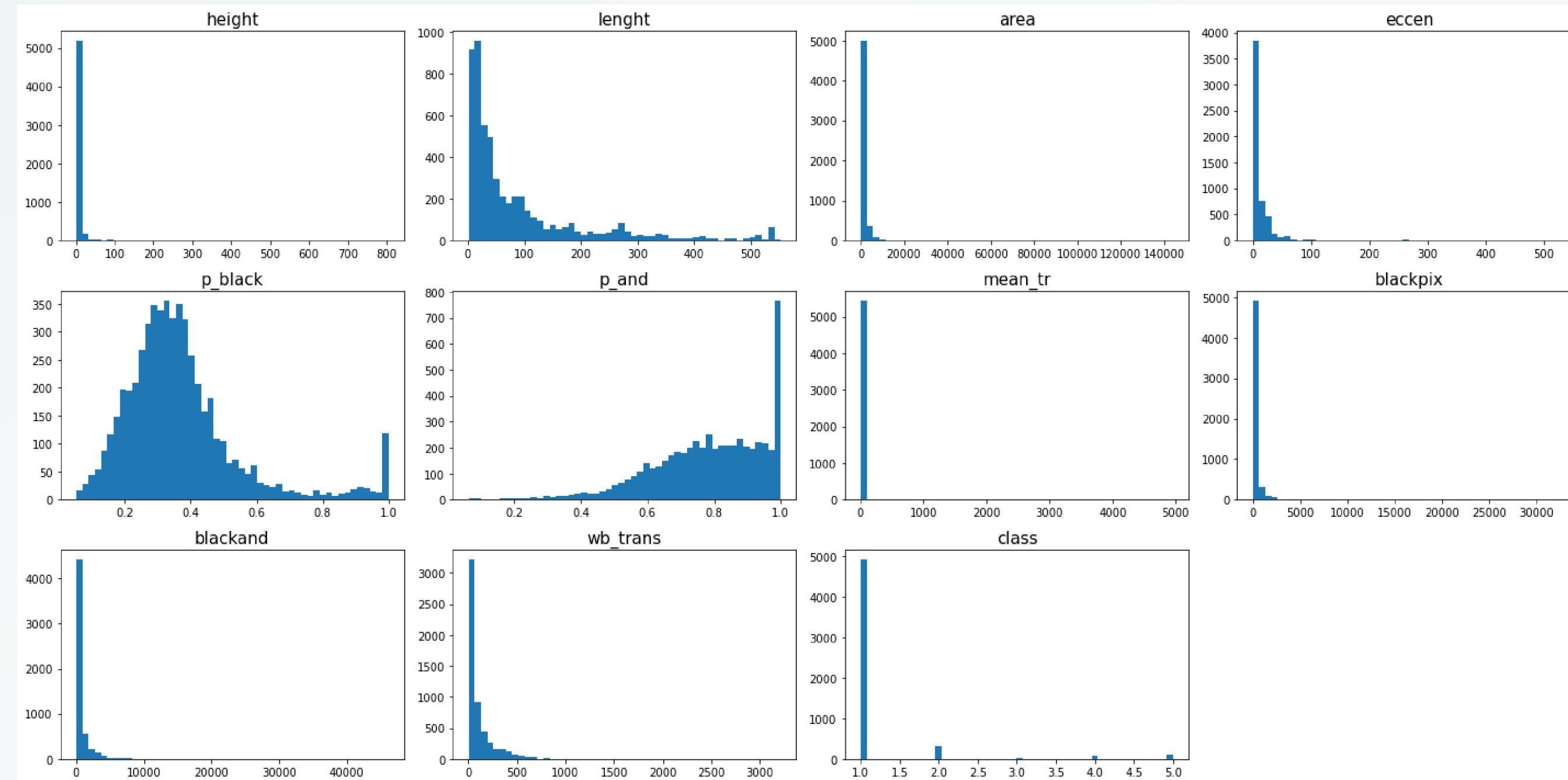


# VISUALIZATIONS



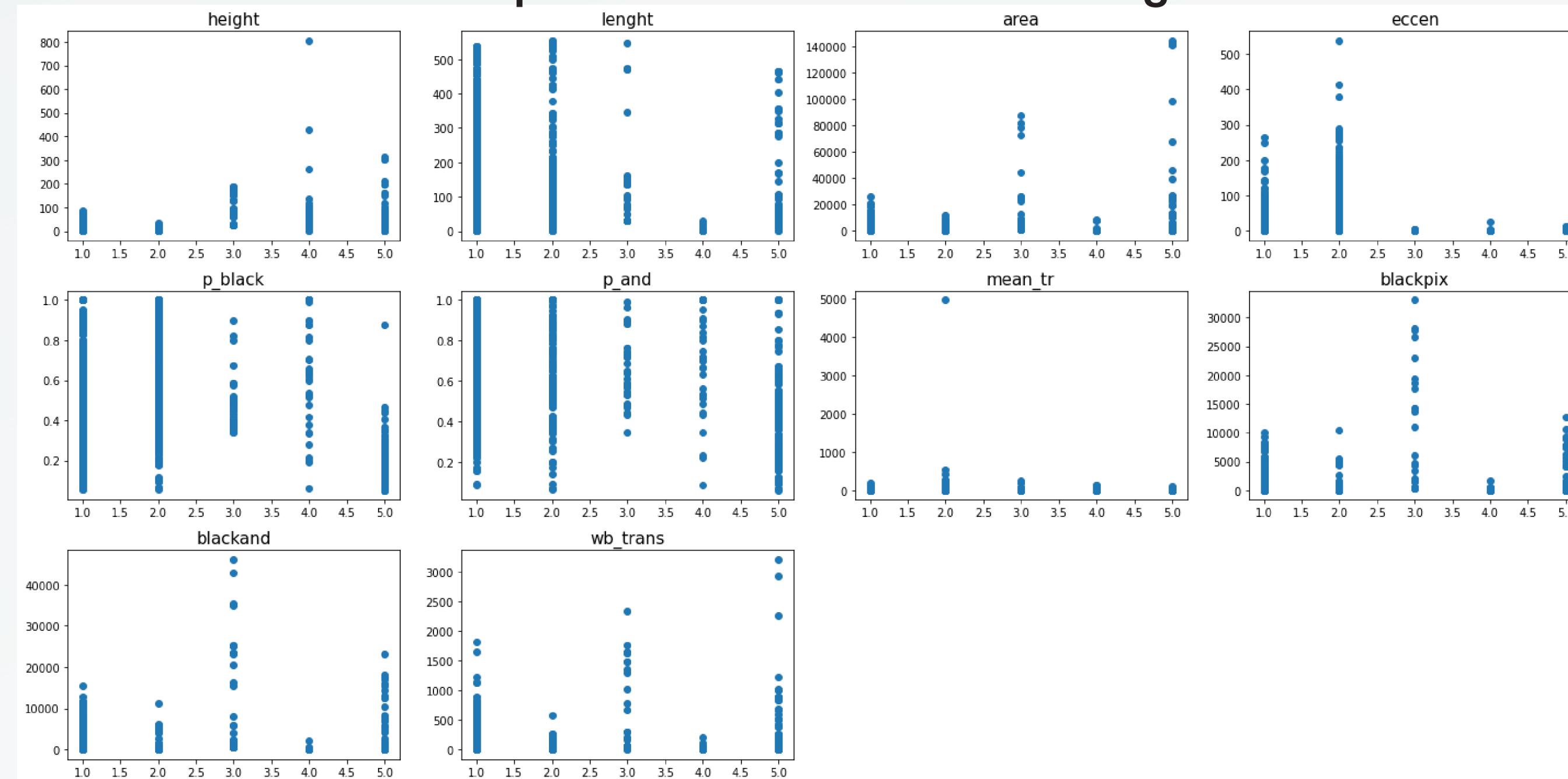
# VISUALIZATIONS

## Distribution of all attributes



# VISUALIZATIONS

The relationship between attributes with target attribute



# OBSERVATION OF THE DATASET

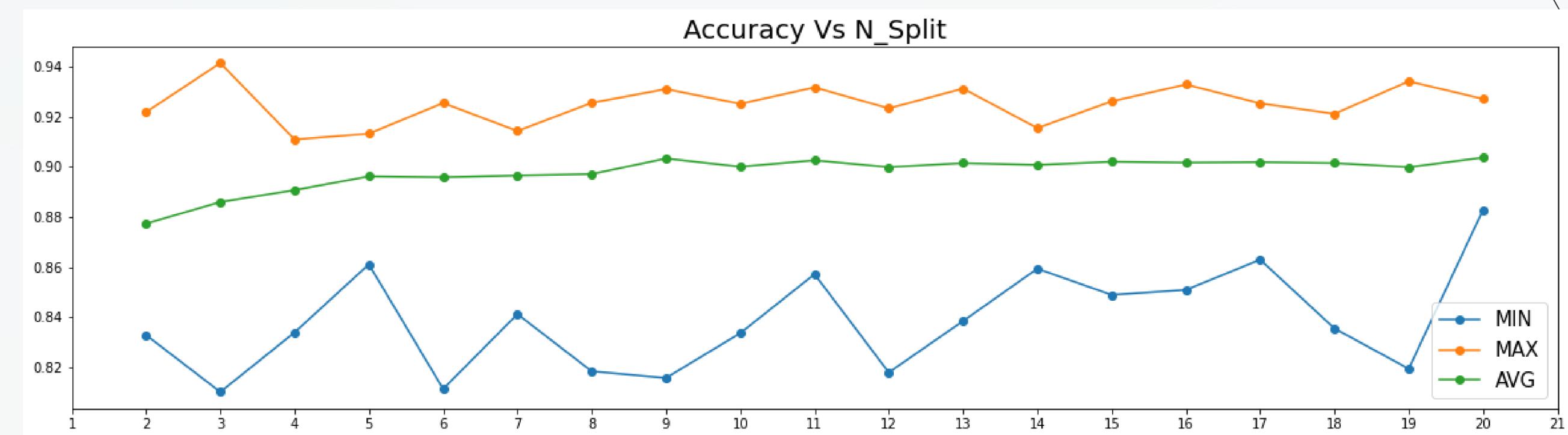
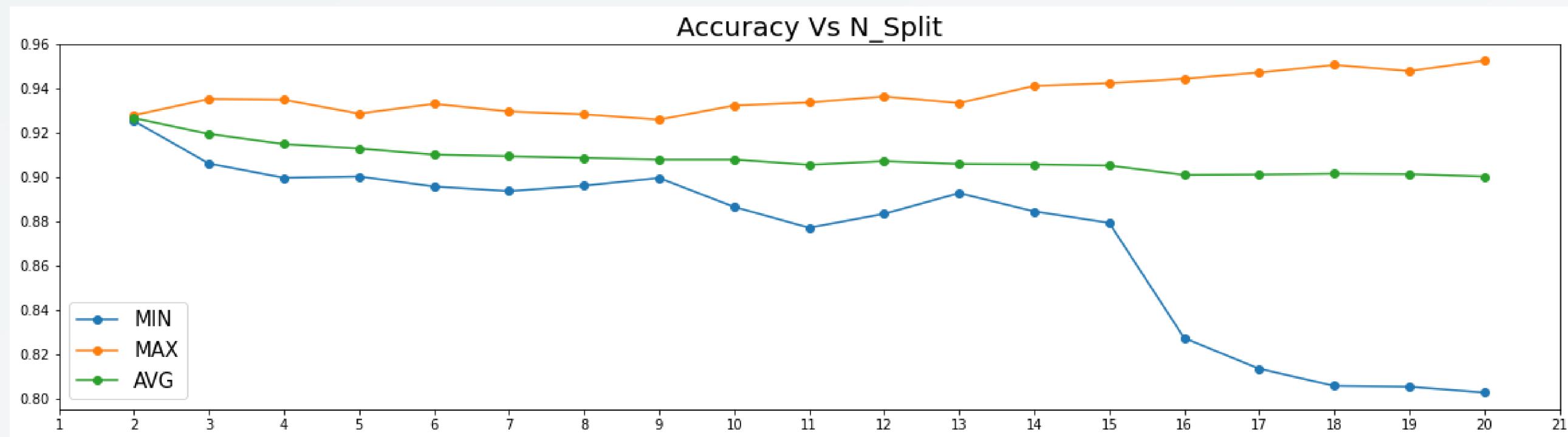
- The dataset is highly imbalanced, most of examples belong to class(1).
- Most blocks have a small height, suggests that blocks of smaller height are more common.
- The area of the block has a high correlation with height, length, blackpix, and blackand.
- The blackpix has a high correlation with blackand.

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 corner. On the right, a more complex, dense cluster of lines forms a shape reminiscent of a stylized leaf or a brain's gyral pattern, extending from the bottom right towards the center.

# **MODEL BUILDING AND EVALUATION**

# VISUALIZATIONS

CALCULATE THE CROSS VALIDATION SCORE ON DIFFERENT SPLIT POINT IN 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

**THANK'S FOR  
WATCHING**

