

# **Indoor Scene Recognition with Data Mining Algorithms**

**Classification methods to predict indoor scene classes**

## GENERAL DESCRIPTION OR PROPERTIES:

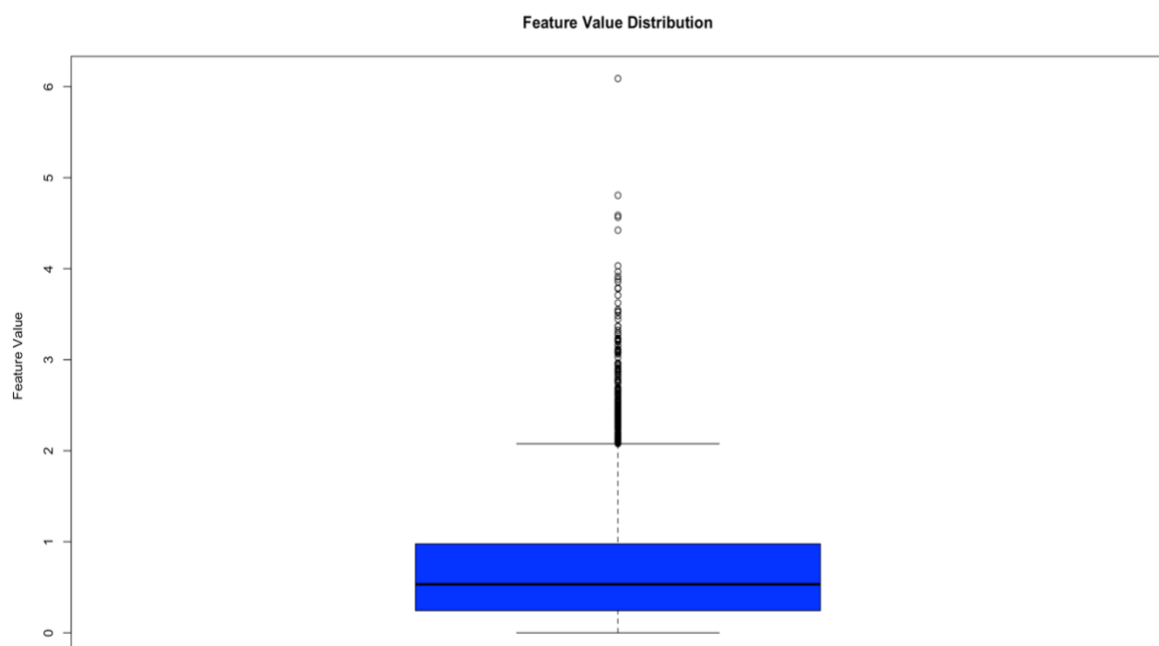
To begin with, the MIT indoor dataset is a public scene recognition dataset that contains 67 indoor scene categories, such as "airport\_inside", "artstudio", "bar", "bookstore", and so on. The dataset consists of a training set with 5360 images and a test set with 1340 images. Each image is represented by a 512-dimensional feature descriptor that was extracted using a ResNet-18 deep learning model pre-trained on the ImageNet dataset. Each image is labeled with one of the 67 different scene classes mentioned above. The feature descriptors are provided in the form of CSV files, with each row representing an image and each column representing a feature.

To get a better understanding of the dataset, we have plotted some basic statistics of the feature descriptors.

## BOX PLOT:

Primarily, to understand the distribution and statistical summary of the feature descriptors for each class label, we have used a box plot. By examining the box plot, we can understand the central tendency, variation, and potential outliers within each class.

- The line in the center of the box represents the median value of the feature distribution. It indicates the midpoint or the 50th percentile of the data. In this case, the median appears to be around 0.5.
- The box in the plot represents the interquartile range, which is the range between the first quartile (25th percentile) and the third quartile (75th percentile) of the data. It provides information about the spread of the central part of the distribution. The larger the box, the greater the spread of the feature values.

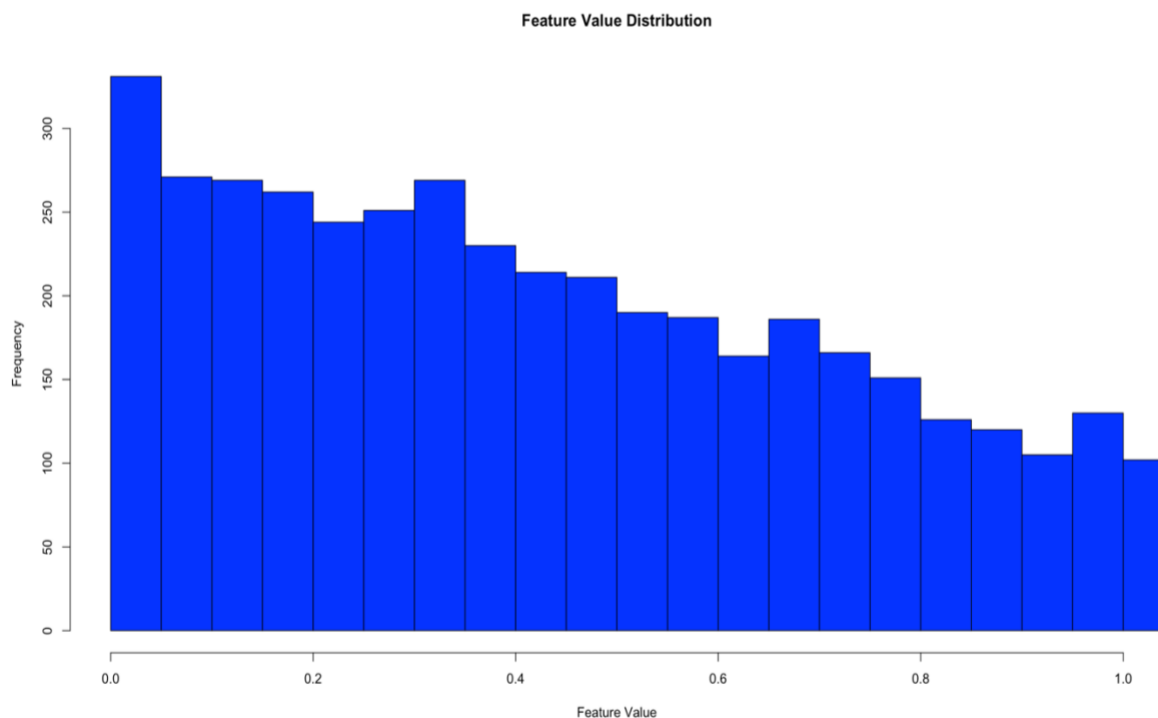


- The whiskers extend from the box and represent the range of non-outlier data points. They can provide insights into the range and variability of the feature values. In this case, the whiskers are relatively short, indicating a limited spread of the data.
- The plot also shows individual points, known as outliers, which are data points that fall outside the whiskers. However, in this case, there are no visible outliers, suggesting that the feature values are relatively consistent and do not deviate significantly from the central distribution.

## HISTOGRAM:

The histogram of the feature values provides insights into the distribution of values within the dataset. In the histogram, the x-axis represents the range of feature values, and the y-axis represents the frequency or count of occurrences. Here are some insights that can be derived from the histogram:

- The histogram appears to have a roughly bell-shaped distribution, with a peak around 0.5. This suggests that the feature values are centered around a certain range, indicating a common characteristic shared by the indoor scene images.
- The histogram provides information about the range and spread of the feature values. In this case, the feature values appear to be concentrated in a relatively narrow range, suggesting that the ResNet-18 model has produced consistent and discriminative feature descriptors for the indoor scene images.



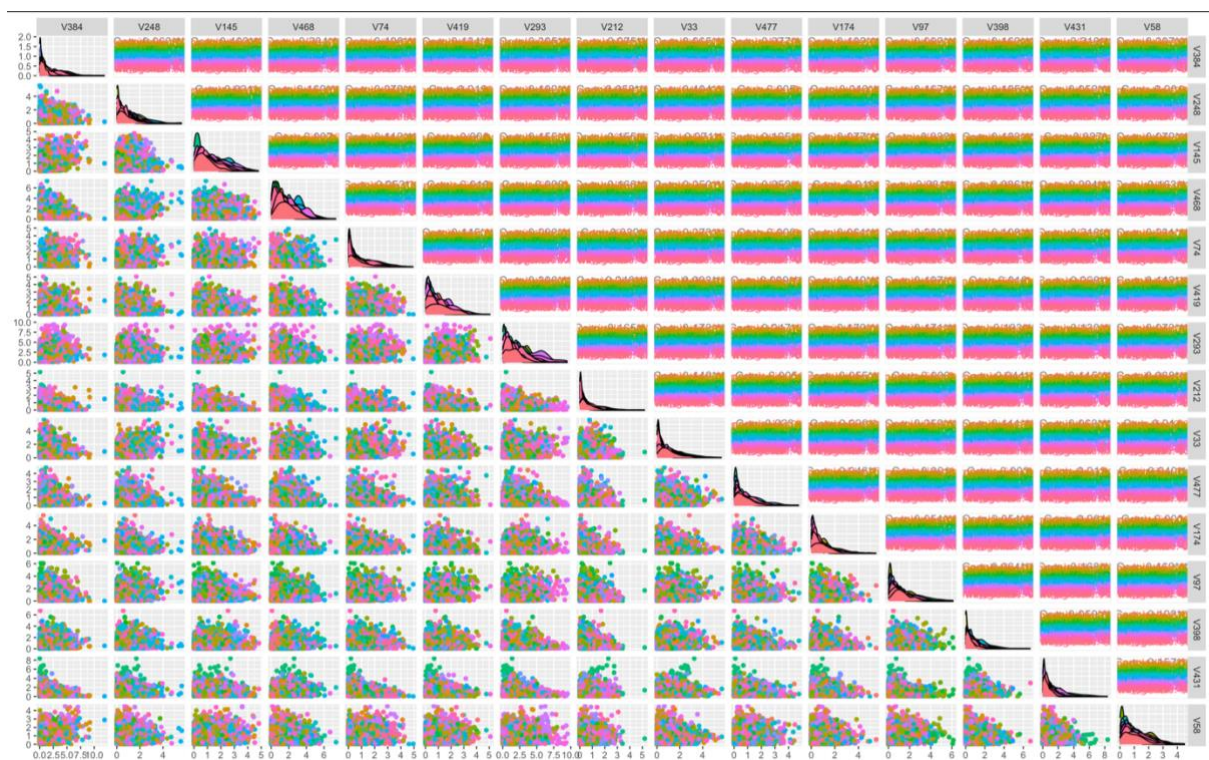
- The histogram can also reveal the presence of outliers, which are data points that significantly deviate from most of the distribution. In this case, there are no apparent outliers, as there are no bars in the histogram that extend far beyond most of the data.

## SCATTERPLOT AND CORRELATION MATRIX:

To understand the correlation between the pairs of attributes, a scatter plot matrix has been used, which provides a visual representation of the relationships better. By examining the scatter plots, you can observe patterns, clusters, or trends in the data, and assess the strength and direction of relationships between variables. After a lot of trial and errors, the following attributes pairs have displayed a good relationship:

V384, V248, V145, V468, V74, V419, V293, V212, V33, V477, V174, V97, V398, V431, V58

This scatter plot matrix provides us a visual representation of the relationships between the attributes selected previously.



- From the scatter plot it can be noticed that the selected attributes present a reasonable correlation with data points forming an almost straight line, but substantially inclined towards the negative as some of the data points start at high y-values on the y-axis and progress down to low values, indicating that the variables have a subtle negative correlation. However, the selected attributes showcase the best correlation among other combinations that have been tried.

- As previously mentioned in other plots, the dataset does not seem to have any prominent outliers as also seen in the scatter plot.
- The overall spread or tightness of data points in a scatter plot can provide an indication of the strength of the relationship between the variables. The clusters in the plot are narrow and well-defined suggesting a strong correlation.

In addition, a correlation matrix was also made for the previously mentioned ideal pair of attributes, to support the attribute correlations which would eventually impact the classification models. The matrix provides a numerical representation of the pairwise correlations between the selected attributes.

- At the first glance at the matrix, positive correlation can be observed as most of the coefficient values of the attribute variables are positive towards +1. The magnitude of the correlation coefficients also indicates a stronger correlation.
- These identified correlated attributes can be important in feature selection, as they may provide redundant information or indicate the influence of certain attributes on the classification task.

