## Data Collection and Preprocessing Phase
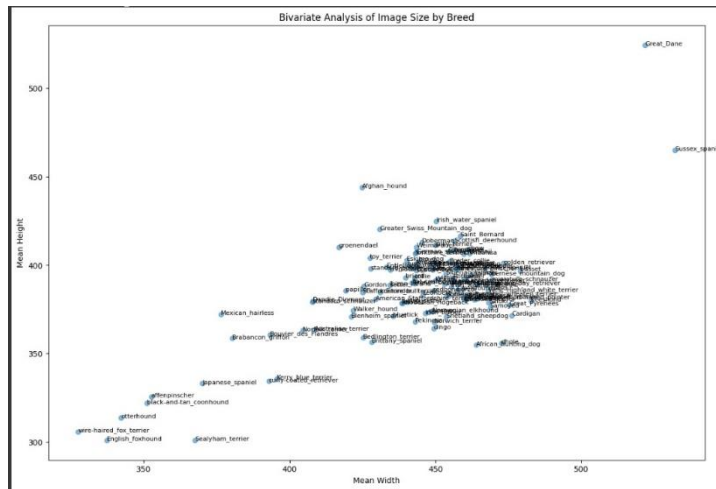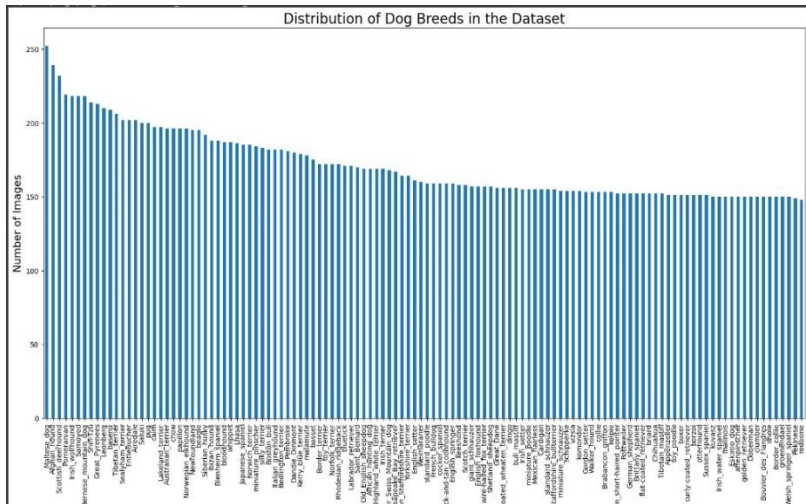
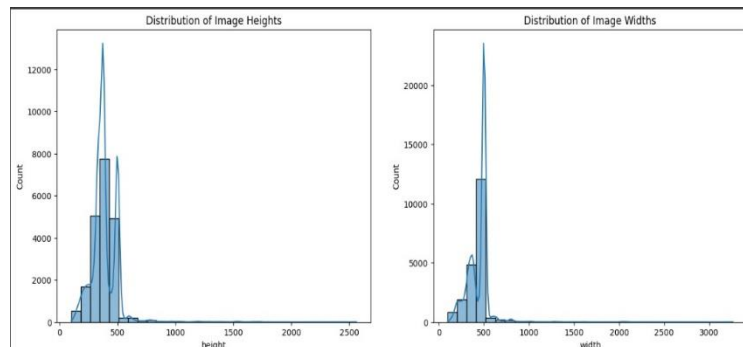| | |
|---|---|
| Date | 27th June 2024 |
| Team ID | SWTID1720073336 |
| Project Title | Dog breed identification using transfer learning |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Report:**

In the dog breed identification project, dataset variables will be statistically analyzed to identify patterns and outliers. Python will be employed for preprocessing tasks such as normalization and feature engineering. This will ensure that the data is standardized and suitable for machine learning algorithms. Through data cleaning processes, missing values and outliers will be addressed, maintaining high data quality for subsequent analysis and modeling. These steps form a strong foundation, enabling the extraction of meaningful insights and the development of accurate predictive models for dog breed identification.
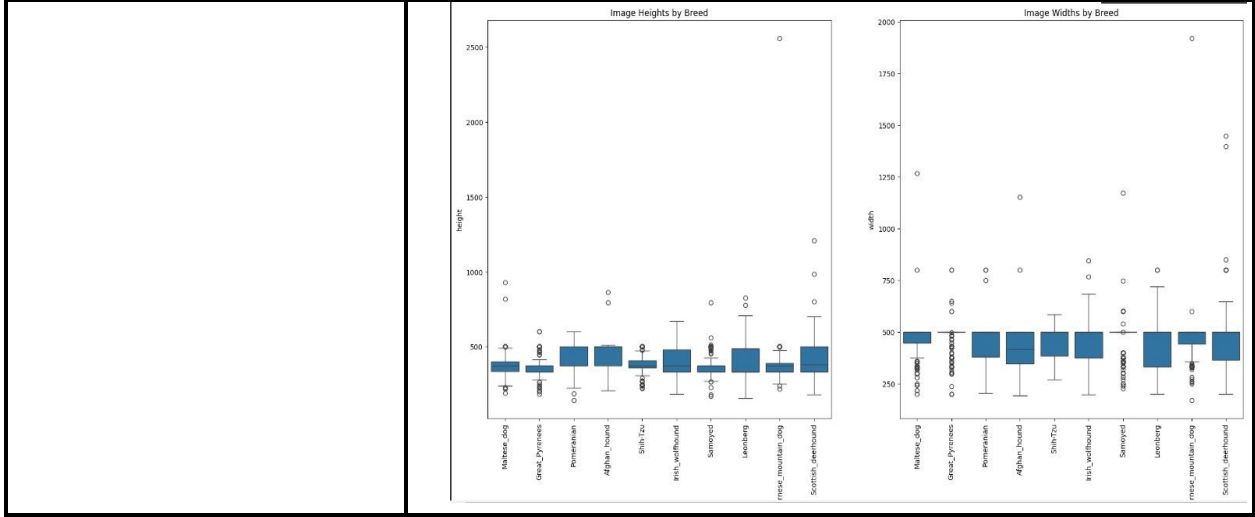
| Section | Description |
|---|---|
| Data Overview | Dimension:<br>Rows: 20580 x<br>columns: 2<br>Descriptive statistics:<br><br>`Image_Path Breed`<br>`0 /content/images/Images/n02107142-Doberman/n021... Doberman`<br>`1 /content/images/Images/n02087046-toy_terrier/n... toy_terrier`<br>`2 /content/images/Images/n02093647-Bedlington_te... Bedlington_terrier`<br>`3 /content/images/Images/n02096177-cairn/n020961... cairn`<br>`4 /content/images/Images/n02111889-Samoyed/n0211... Samoyed`<br>`Number of rows: 20580`<br>`Number of columns: 2` |
| Univariate Analysis | |

| | |
|---|---|
| | Distribution of Dog Breeds in the Dataset |
| Bivariate Analysis | Bivariate Analysis of Image Size by Breed |
| Multivariate Analysis | Distribution of Image Heights / Distribution of Image Widths |

| Outliers and Anomalies | - |
|---|---|

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Unzipping Data | ```python
import zipfile
import os

# Unzip
with zipfile.ZipFile('/content/stanford-dogs-dataset.zip', 'r') as zip_ref:
    zip_ref.extractall('/content/stanford-dogs-dataset')


unzipped_dir = '/content/stanford-dogs-dataset'
for root, dirs, files in os.walk(unzipped_dir):
    for name in files:
        print(os.path.join(root, name))
``` |
| Data Transformation | ```python
import pandas as pd

df = pd.DataFrame({
    'Image_Path': X,
    'Breed': y
})

print(df.head())

print(f"Number of rows: {df.shape[0]}")
print(f"Number of columns: {df.shape[1]}")
``` |