

EDA

PLANT VILLAGE DATASET

ABOUT THE DATASET

The dataset contains images of plant leaf samples collected under two circumstances.

1. When the plant is healthy
2. When the plant is suffering from a disease

The dataset is highly diverse encompassing various species and their condition during different types of diseases.

The dataset includes a total of 14 species of plants and 38 different classes of images.



INITIAL ANALYSIS

The dataset consists of 14 different species of plants.

```
['Apple', 'Blueberry', 'Cherry', 'Corn', 'Grape',
'Orange', 'Peach', 'Pepper', 'Potato', 'Raspberry',
'Soybean', 'Squash', 'Strawberry', 'Tomato']
```

The dataset takes into account multiple disease conditions for each plant which amount to around 38 different classes in which the dataset is divided.

```
['Apple__Apple_scab' 'Apple__Black_rot' 'Apple__cedar_apple_rust'
'Apple__healthy' 'Blueberry__healthy'
'Cherry_(including_sour)__Powdery_mildew'
'Cherry_(including_sour)__healthy'
'Corn_(maize)__Cercospora_leaf_spot_Gray_leaf_spot'
'Corn_(maize)__Common_rust_-' 'Corn_(maize)__Northern_Leaf_Blight'
```

```
'Corn_(maize)__healthy' 'Grape__Black_rot'
'Grape__Esca_(Black_Measles)'
'Grape__Leaf_blight_(Isariopsis_Leaf_Spot)' 'Grape__healthy'
'Orange__Haunglongbing_(Citrus_greening)' 'Peach__Bacterial_spot'
'Peach__healthy' 'Pepper,_bell__Bacterial_spot'
'Pepper,_bell__healthy' 'Potato__Early_blight' 'Potato__Late_blight'
'Potato__healthy' 'Raspberry__healthy' 'Soybean__healthy'
```

```
'Squash__Powdery_mildew' 'Strawberry__Leaf_scorch'
'Strawberry__healthy' 'Tomato__Bacterial_spot' 'Tomato__Early_blight'
'Tomato__Late_blight' 'Tomato__Leaf_Mold' 'Tomato__Septoria_leaf_spot'
'Tomato__Spider_mites Two-spotted_spider_mite' 'Tomato__Target_Spot'
'Tomato__Tomato_Yellow_Leaf_Curl_Virus' 'Tomato__Tomato_mosaic_virus'
'Tomato__healthy']
```

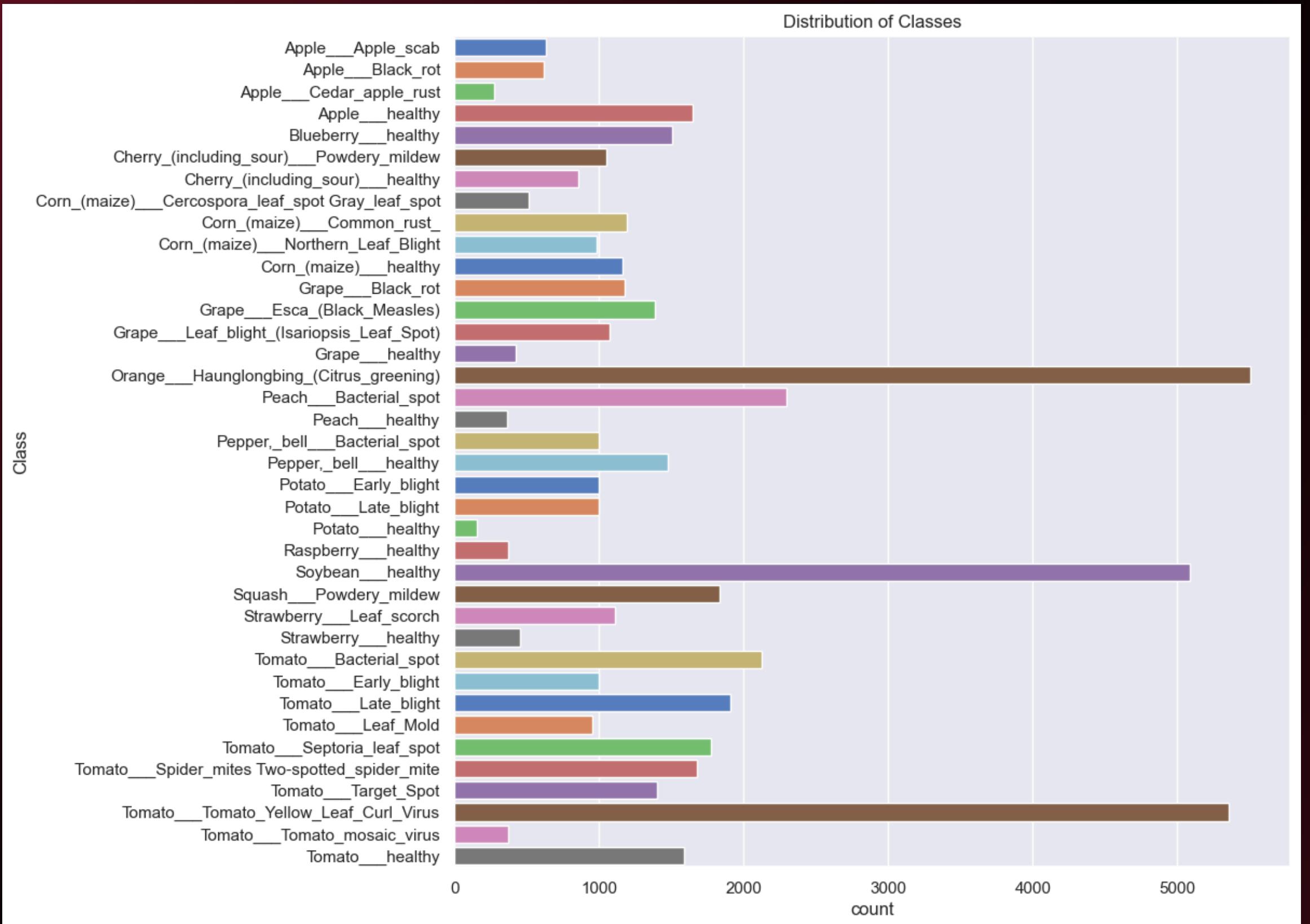
ASPECT RATIO AND RESOLUTION

The images in dataset are uniform in size and resolution. The images share a common resolution and aspect ratio.

ASPECT RATIO : 1.0

RESOLUTION : 256 x 256 = 65536 pixels

CLASS DISTRIBUTION



INFERENCE :

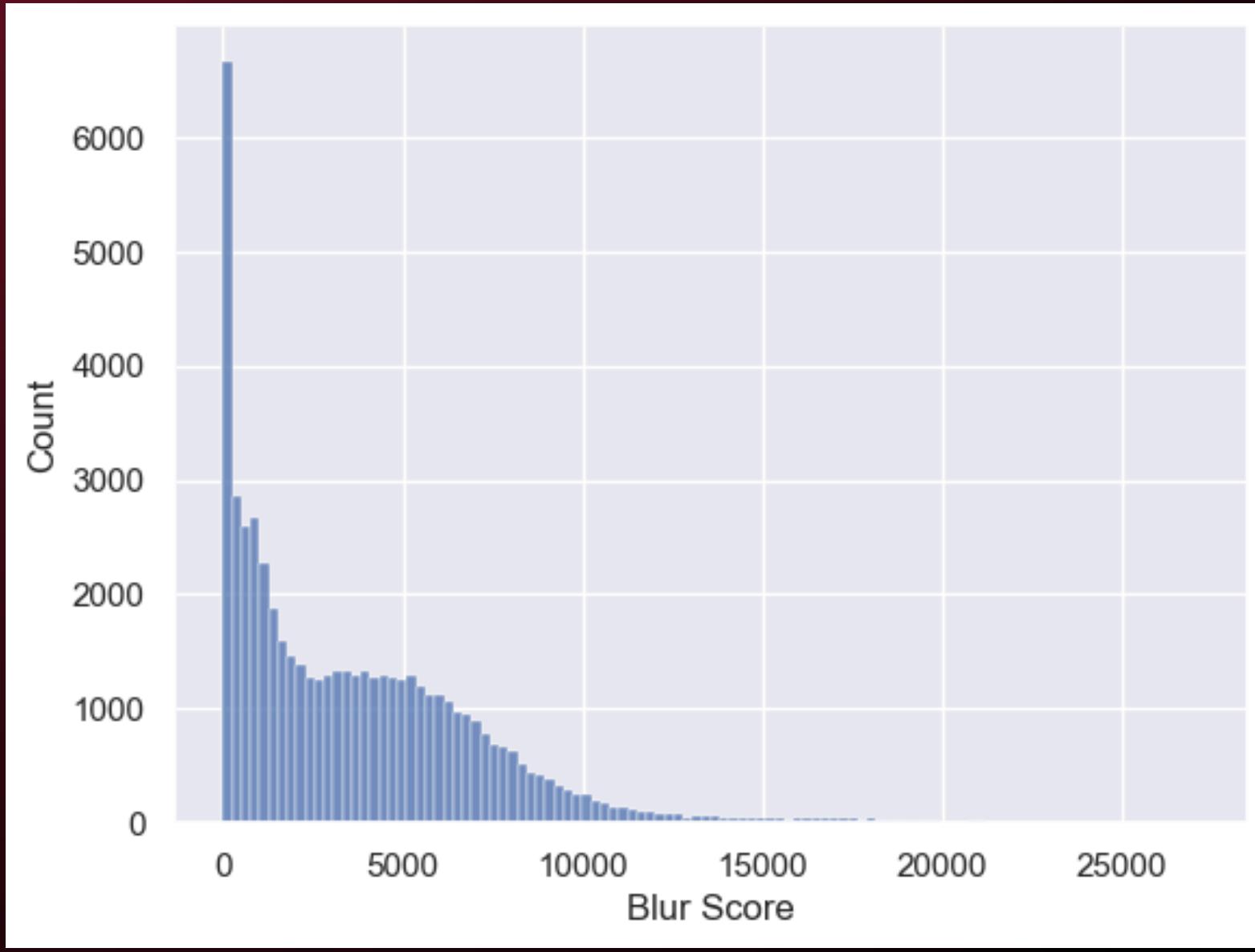
The dataset presents a classic case of class imbalance.

Certain classes are very highly represented having more than 5000 samples.

On the other hand certain classes are very poorly represented having less than 500 samples.

This class imbalance is very unfavorable as this makes the model biased towards the dominant classes.

BLURRY IMAGES



The histogram suggests that many images have very low blur score indicating that a substantial portion of dataset comprises of blurry photos. Such photos if used will yield poor results.

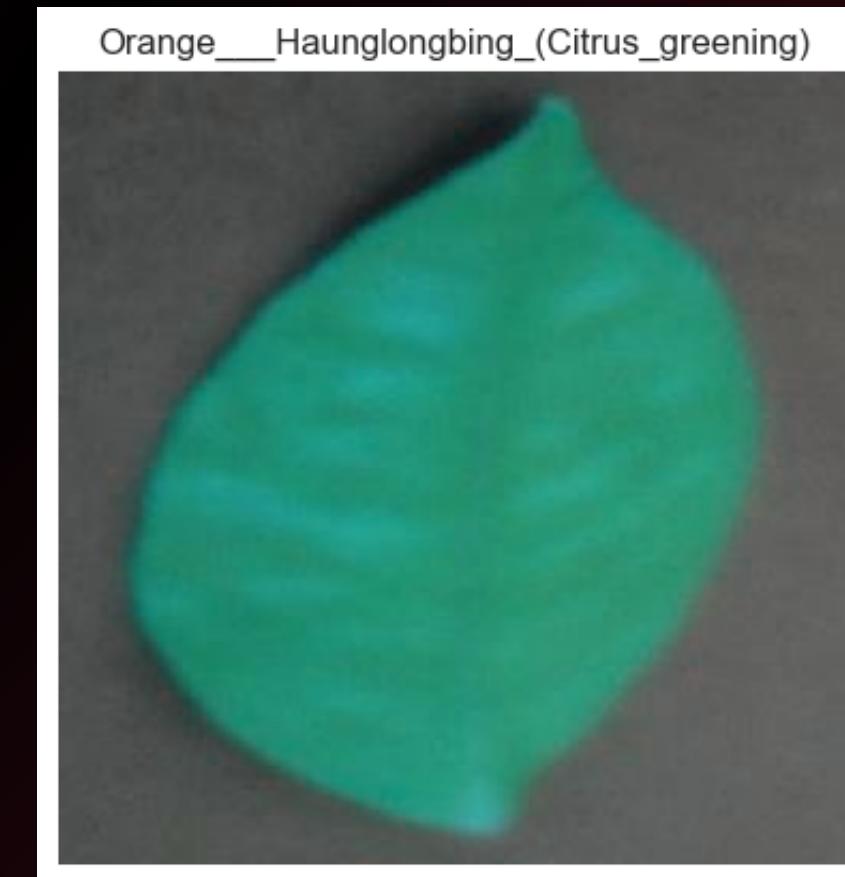
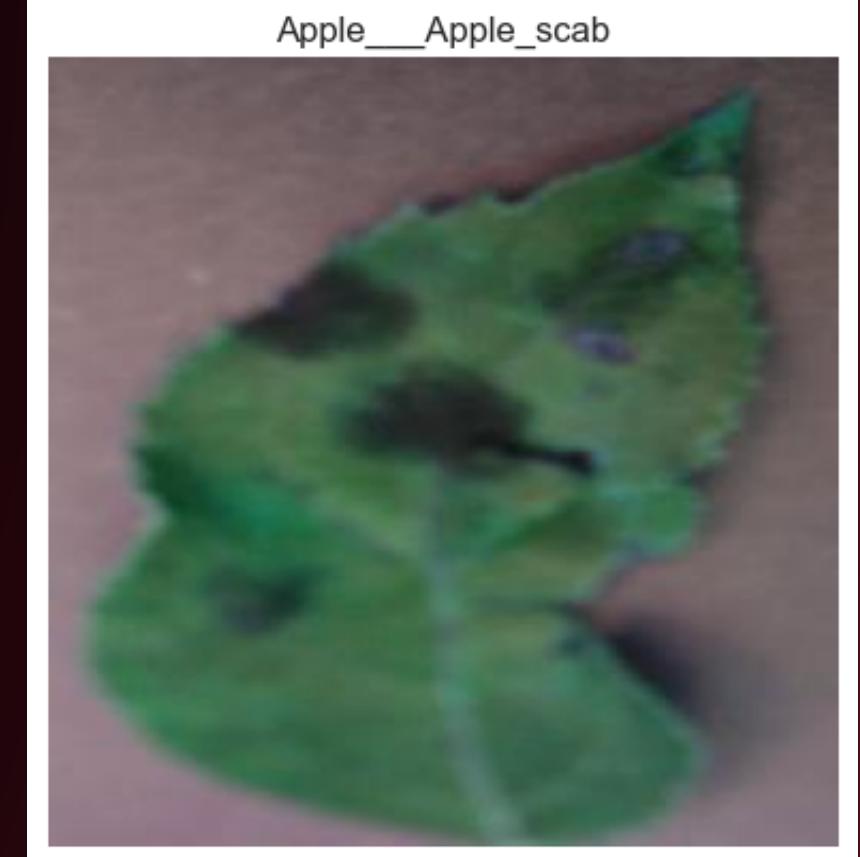


Fig. Examples of blurry images

VARIATION IN QUALITY

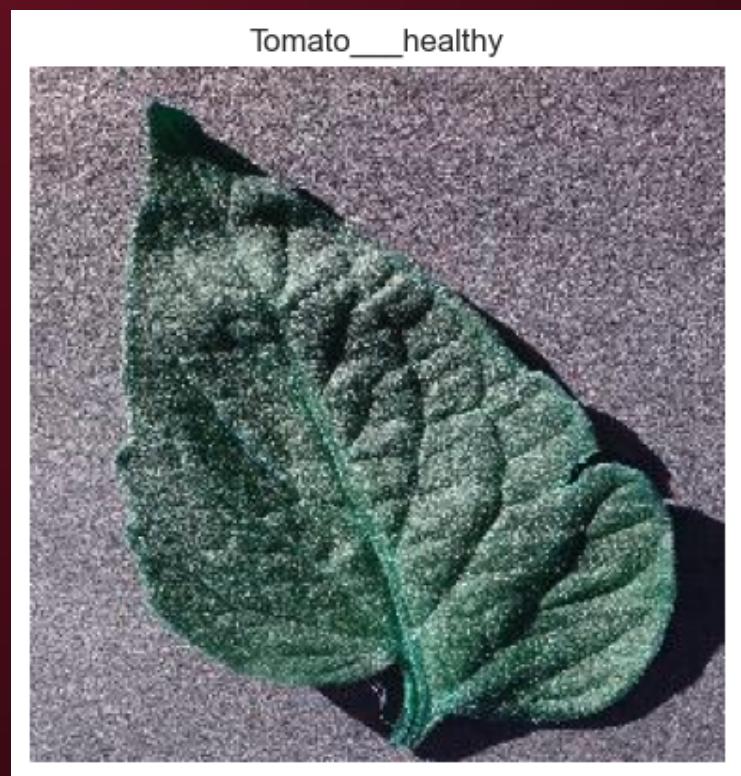


Fig. Sharpest Images

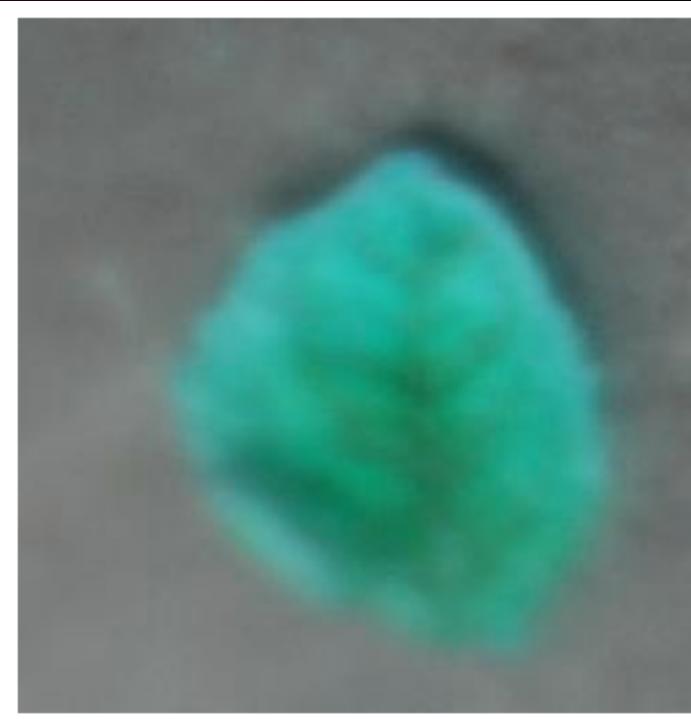


Fig. Most blurry images

Fig. Average Sharpness

VARIATION IN BRIGHTNESS

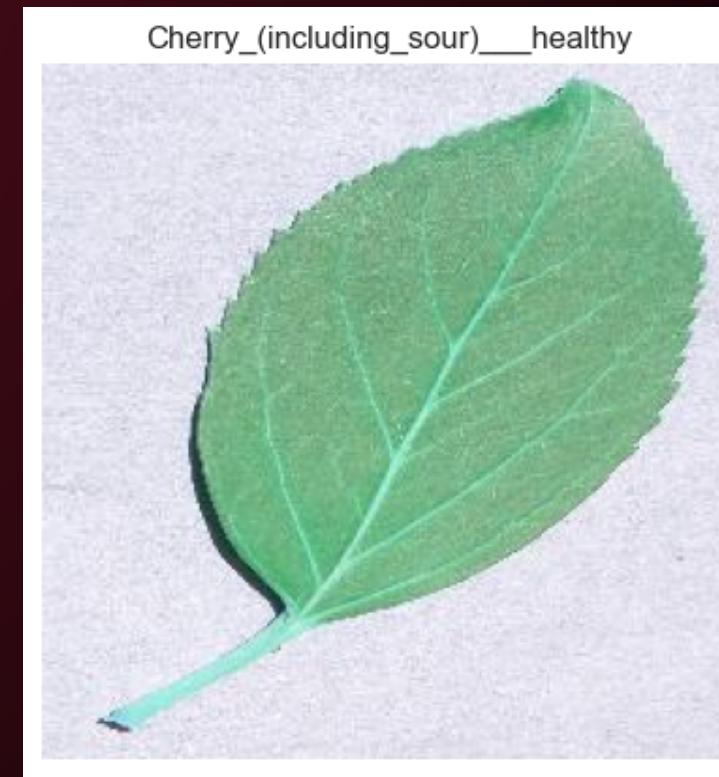
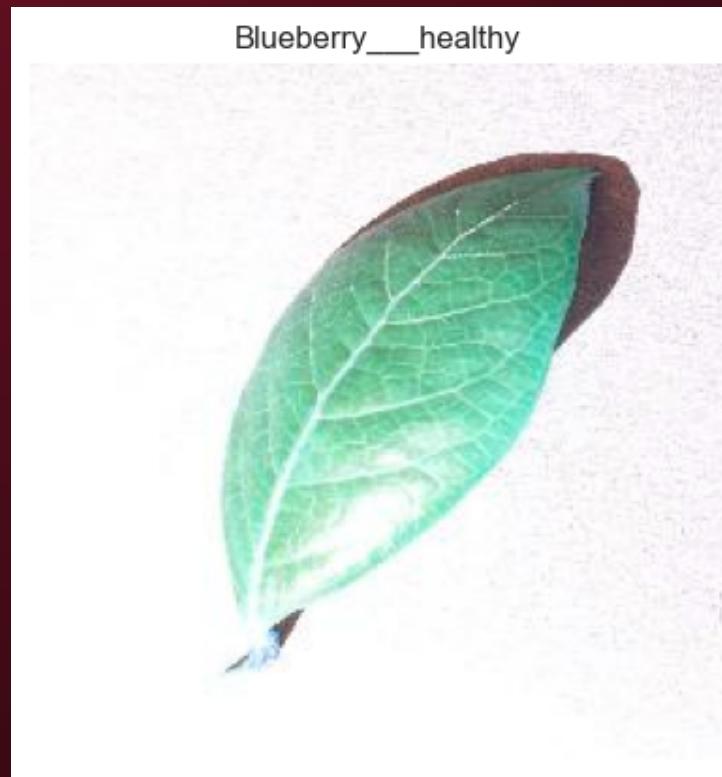


Fig. Brightest Images

Fig. Average Brightness

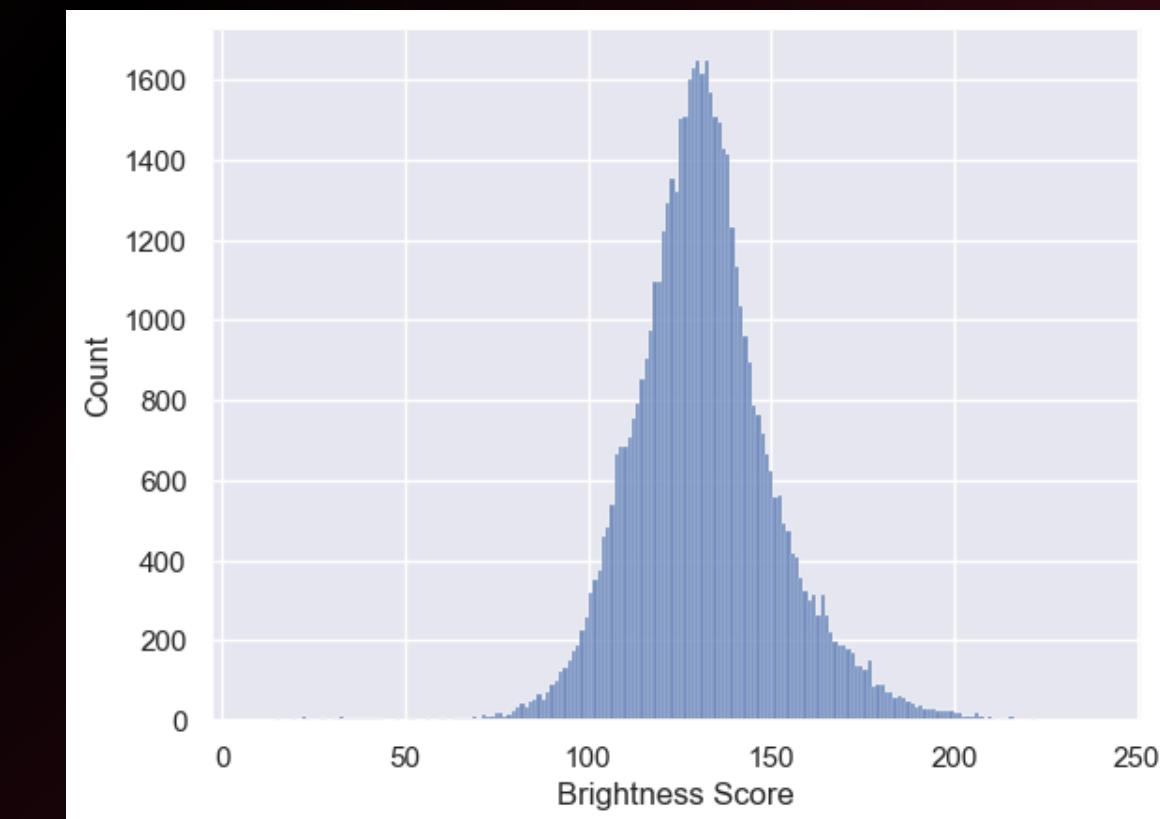
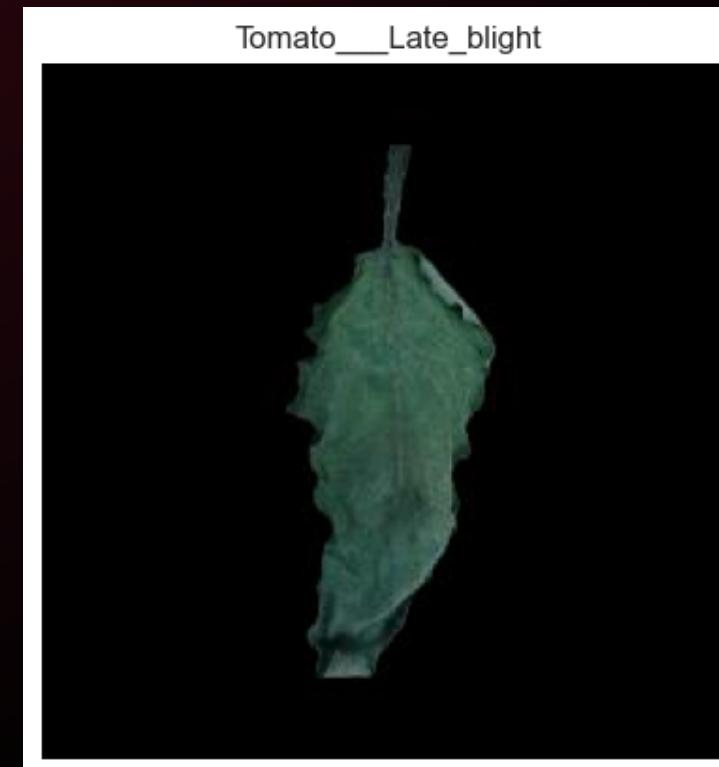
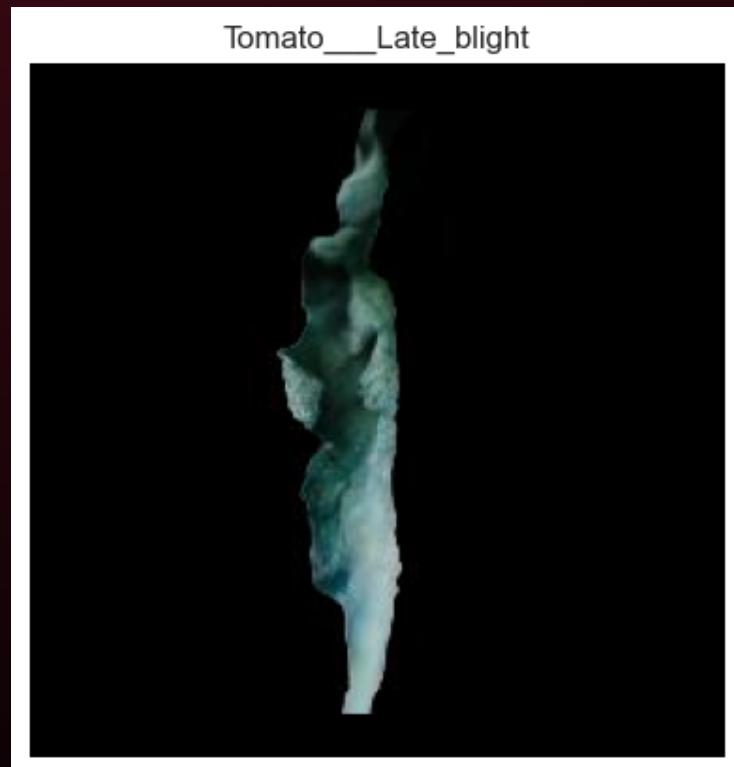


Fig. Dullest Images

The plot suggests that most of the images have brightness around the mean. Short left tail suggests that darkness is rare in the dataset and the images are more or less well lit.

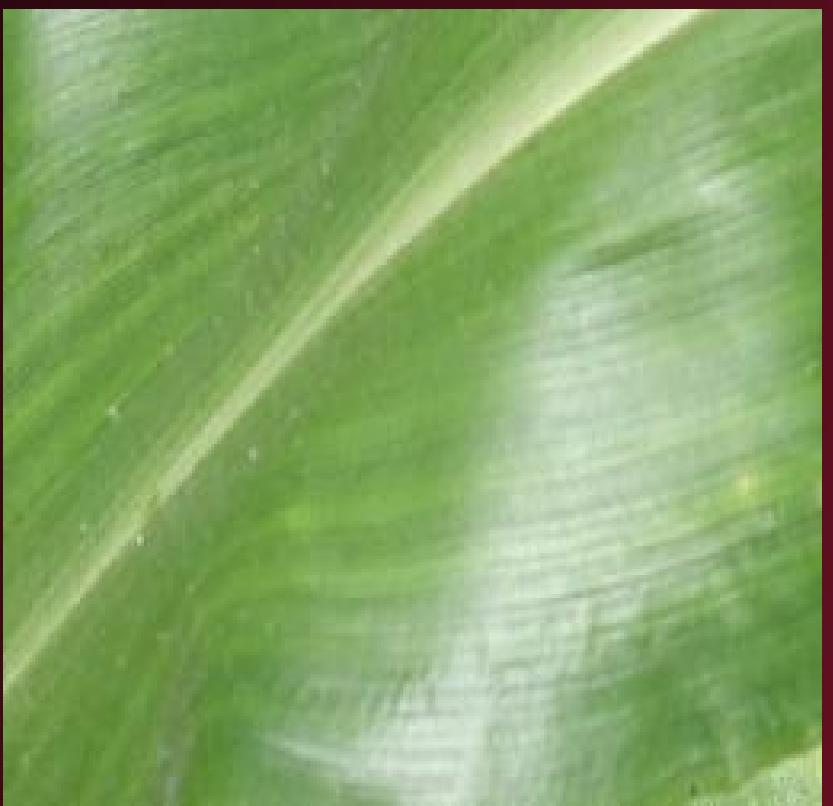
VARIATION IN BACKGROUND

Background in an image is an important aspect that must not be overlooked. If the images have varying backgrounds, it becomes a proxy variable for the model.

This means that, instead of learning what the leaf looks like (texture, veins, disease spots), the model might learn to cheat by just looking at the background color.

Our dataset consists of majority of images with light background.

About 2.5 % of images have black background. A certain portion of images contain leaf in the entire region.



THANK YOU