

Biostatistics (STA5195) Lab 4 Report

Rice Seed Classification

Christopher Rutherford

September 29, 2020

Abstract

Computer algorithms have been used to process and “see” images to learn what objects are contained in images. [1] More often than not, these machine and deep learning algorithms are used to identify and classify large quantities of images. We can simplify these algorithms in order to classify images of rice seeds as proper or broken. By assessing the brightness of each pixel and converting the image to pure black and white, we present a method for classifying these rice seeds.

1 Introduction

Training machine and deep learning algorithms for image classification is often computationally and time intensive, particularly when dealing with larger datasets. Instead of using a specific algorithm, we will instead analyze a set of 397 images of rice seeds - some broken, and the rest normal (“proper”).

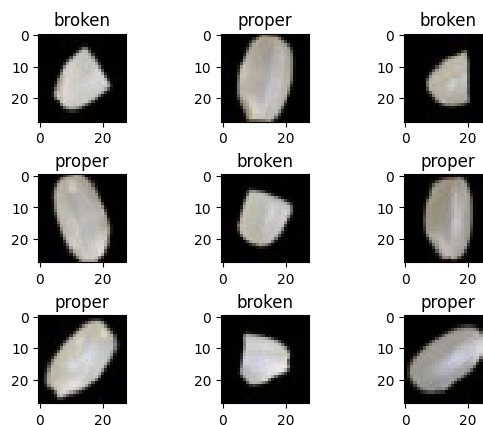


Figure 1: Assortment of rice seed images

2 Theory

Because the rice seed images are quite simple images - a white/gray seed on a black background - we can take advantage of this when working with these images. We convert each image from RGB to grayscale so that each pixel is only defined by how bright/dark it is. The value 0 is associated with completely black pixels while 255 is associated with white pixels. We use scikit-image's `threshold_otsu` function to map anything in between 0 and 255 to either be black if it is closer to 0, and white if it is closer to 255. An example of how this conversion works:



Figure 2: Conversion: RGB \rightarrow grayscale \rightarrow black & white

3 Procedures

To begin, we import the libraries and functions that we will be using throughout this lab:

```
1 import numpy as np
2 import math, os, sys
3 import itertools
4
5 import matplotlib.pyplot as plt
6 plt.style.use('default')
7 from scipy import ndimage
8
9 #functions for image processing
10 from skimage import measure, morphology
11 from skimage.io import imsave, imread
12 from skimage.color import rgb2gray
13 from skimage.filters import threshold_otsu
14 from skimage.transform import resize
15
16 from sklearn import svm, datasets
17 from sklearn.metrics import confusion_matrix
18 import pandas as pd
```

And now we need to download the dataset that we are going to be using:

```
1 !#we will use the rice seed image dataset
2 !apt-get install subversion > /dev/null
```

```

3 !svn export https://github.com/totti0223/
  deep_learning_for_biolologists_with_keras/trunk/notebooks/data/
  image image > /dev/null

```

We can make sure the images imported properly by plotting one of them:

```

1 #let's visualize a single file
2 image = imread("image/train/proper/100.jpg") #read in the image
3 plt.figure(figsize=(3,3)) #plot the image on a 3x3 figure
4 plt.imshow(image)

```

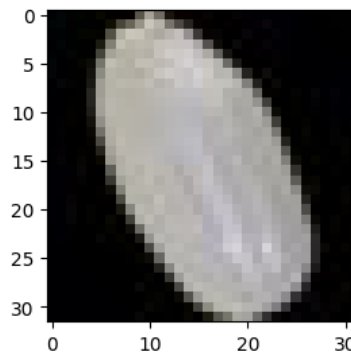


Figure 3: Proper rice seed

Since we have the images that we need, we can start preparing our data for our algorithm by splitting it into train/test splits. First we prepare the training set:

```

1 #load training dataset
2 X_train = [] #create empty lists for training data
3 y_train = []
4
5 #cycle through the images
6 for root, dirs, files in os.walk("image/train"):
7     files = [x for x in files if x.endswith(".jpg")]
8     for file in files:
9         image_path = os.path.join(root, file)
10
11         image = imread(image_path)/255.
12         image = resize(image,(28,28)) #make all the images 28x28
13         X_train.append(image)
14         category = os.path.split(root)[-1]
15         if category == "proper":
16             y_train.append(0) #0 for proper rice seeds
17         else:
18             y_train.append(1) #1 for broken rice seeds
19 X_train = np.array(X_train)
20 y_train = np.array(y_train)

```

A good portion of this code is fail-safes to avoid importing the wrong files, but all we are really doing is adding images to `X_train`, and making the corresponding `y_train` value 0 for proper rice seeds and 1 for broken seeds. We also do the same for the testing set:

```

1 #load test dataset
2 X_test = [] #create empty lists for testing data
3 y_test = []
4
5 for root, dirs, files in os.walk("image/test"):
6     files = [x for x in files if x.endswith(".jpg")]
7     for file in files:
8         image_path = os.path.join(root, file)
9
10        image = imread(image_path)/255.
11        image = resize(image,(28,28))
12        X_test.append(image)
13        category = os.path.split(root)[-1]
14        if category == "proper":
15            y_test.append(0)
16        else:
17            y_test.append(1)
18 X_test = np.array(X_test)
19 y_test = np.array(y_test)
20
21 print("train dataset shape is:", X_train.shape,y_train.shape)
22 print("test dataset shape is:", X_test.shape,y_test.shape)

```

train dataset shape is: (377, 28, 28, 3) (377,)

test dataset shape is: (20, 28, 28, 3) (20,)

Before converting the image to pure black and white, we convert it from RGB to grayscale:

```

1 #gray conversion
2 gray = rgb2gray(image)
3 print(gray.shape)
4 plt.figure(figsize=(3,3))
5 plt.imshow(gray, cmap=plt.cm.gray)
6 plt.title("gray converted")

```

Then, by using the threshold function, pixels that are darker than the specified threshold are made black, while those above the threshold are made white.

```

1 #binary conversion
2 # white/gray pixels above a certain threshold are made white
3 # black/darker gray pixels below the threshold are made black
4 threshold = threshold_otsu(gray)
5 binary = gray > threshold
6 plt.figure(figsize=(3,3))
7 plt.imshow(binary, cmap=plt.cm.gray)

```

We can then calculate the area of the rice seed in each image by counting how many white squares are in the image.

```

1 # calculates the area and lengths of major and minor axes
2 label_im, nb_labels = ndimage.label(binary)
3 regionprops = measure.regionprops(label_im, intensity_image=gray)
4 regionprop = regionprops[0]
5
6 print("area is", regionprop.area)
7 print("major axis length is", regionprop.major_axis_length)
8 print("minor axis length is", regionprop.minor_axis_length)

```

area is 355
major axis length is 28.54477065932296
minor axis length is 15.906083837305355

For the sake of convenience, we can wrap this up into a function and test it out.

```

1 #bundling the above into a function
2 def quantify_area(image):
3     gray = rgb2gray(image)
4     threshold = threshold_otsu(gray)
5     binary = gray > threshold
6     label_im, nb_labels = ndimage.label(binary)
7     regionprops = measure.regionprops(label_im,
8                                       intensity_image=gray)
9     regionprop = regionprops[0]
10    area = regionprop.area
11    return area
12
13 #test
14 area = quantify_area(image)
15 print(area)

```

355

We now create new train/test splits that contain the values of the area of each rice seed's image.

```

1 X_train_area = []
2 for image in X_train:
3     area = quantify_area(image)
4     X_train_area.append(area)
5
6 X_test_area = []
7 for image in X_test:
8     area = quantify_area(image)
9     X_test_area.append(area)

```

```

1 #check the calculated data area value of training dataset
2 plt.scatter(range(len(X_train_area)), X_train_area, c=y_train, cmap=
    "jet")
3 plt.xlabel("rice seed images")
4 plt.ylabel("area (px)")

```

```
5 plt.show()
```

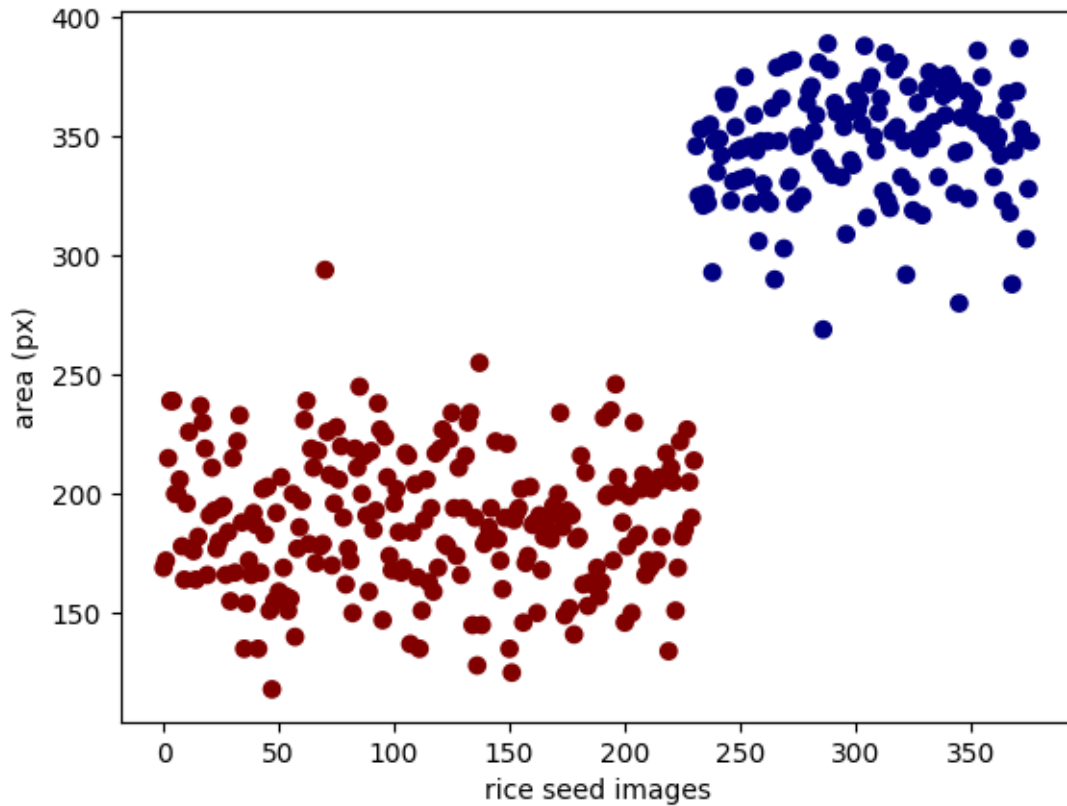


Figure 4: Area of rice seeds in training set (blue = regular, red = broken)

4 Analysis

Now that we have a visual of how our rice seeds are distributed in terms of white pixel area, we can create our own algorithm to classify the seeds. Based on the above plot, we add a horizontal line at $x=260$ to separate the two groups.

```
1 #define an area threshold that can separate the two classes
2 area_threshold = 260
3
4 #classify whether the image is a proper seed or a broken seed
   according to the area_threshold value
5 train_y_pred = []
6 for area in X_train_area:
7     if area > area_threshold:
8         train_y_pred.append(0)
9     else:
10         train_y_pred.append(1)
11
```

```

12 #plot scatter with threshold line
13 plt.figure(figsize=(5,3))
14 plt.scatter(range(len(X_train_area)),X_train_area,c=y_train,cmap=plt
    .cm.coolwarm)
15 plt.axhline(y=area_threshold)
16 plt.title("blue:proper seed, red: broken seed")
17 plt.show()
18
19 #calculate confusion matrix
20 cnf = confusion_matrix(y_train, train_y_pred)
21
22 #confusion matrix in figure
23 plt.figure(figsize=(3,3))
24 plot_confusion_matrix(cnf, classes=["proper","broken"])
25
26 plt.show()

```

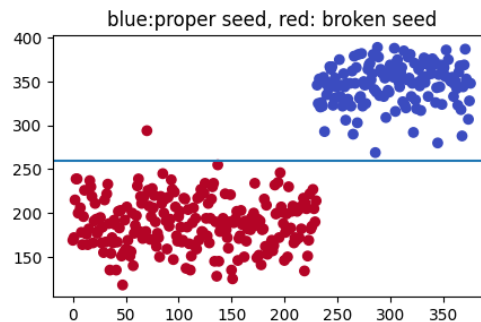


Figure 5: Rice seed areas with threshold line

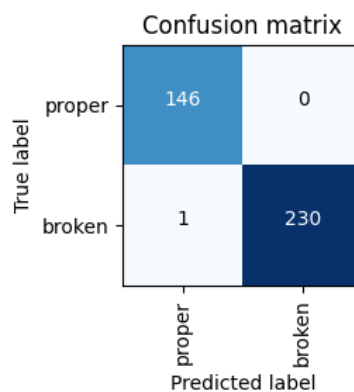


Figure 6: Confusion matrix for training set

Our line seems to do about as well as it could for a simple horizontal line - only one broken seed is incorrectly classified as proper. Simply adding a slope to this would fix this issue, but we're keeping things quite simple.

Now we evaluate this on the test set:

```
1 #evaluate it with the test dataset
2 test_y_pred = []
3 for area in X_test_area:
4     if area > area_threshold:
5         test_y_pred.append(0)
6     else:
7         test_y_pred.append(1)
8
9 #plot scatter with threshold line
10 plt.figure(figsize=(5,3))
11 plt.scatter(range(len(X_test_area)),X_test_area,c=y_test,cmap=plt.cm
12             .coolwarm)
13 plt.axhline(y=area_threshold)
14 #plt.plot([100,0],[100,350],'k-',lw=2)
15 plt.show()
16
17 #calculate confusion matrix
18 cnf = confusion_matrix(y_test, test_y_pred)
19
20 #confusion matrix in figure
21 plt.figure(figsize=(3,3))
22 plot_confusion_matrix(cnf, classes=["proper","broken"])
23 plt.show()
```

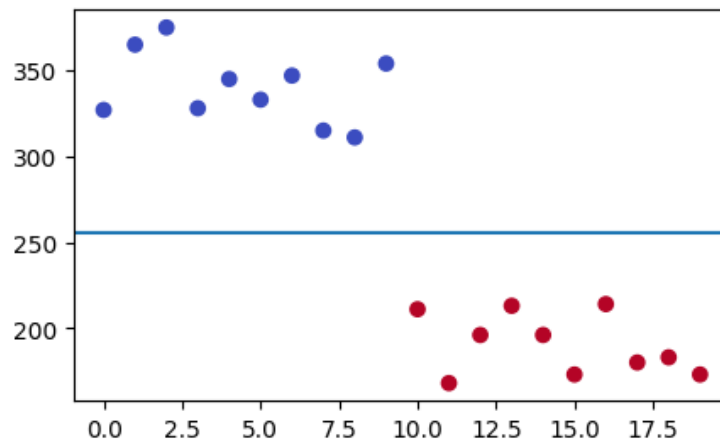


Figure 7: Rice seed areas with threshold line (testing set)

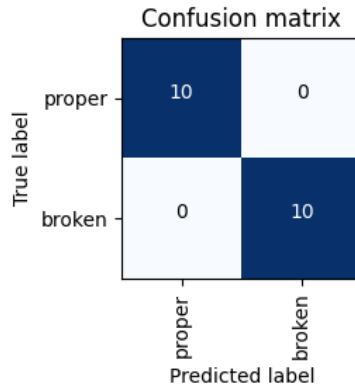


Figure 8: Confusion matrix for testing set

With such little error in the training set, it's no surprise that this managed to perform well on the testing set as well.

We can verify how this works on another, single image from the test set and a threshold of 350:

```

1 #build a classifier
2
3 def manual_classifier(image, area_threshold):
4     gray = rgb2gray(image)
5     threshold = threshold_otsu(gray)
6     binary = gray > threshold
7     label_im, nb_labels = ndimage.label(binary)
8     regionprops = measure.regionprops(label_im, intensity_image=gray
9 )
10    regionprop = regionprops[0]
11    area = regionprop.area
12    if area > area_threshold:
13        return 0
14    else:
15        return 1
16
17 # get a image from test dataset
18 #value must be lower than the size of the test dataset(20-1)
19 n = 10
20 image = X_test[n]
21 label = y_test[n]
22 area_threshold = 350
23 prediction = manual_classifier(image, area_threshold)
24
25 plt.imshow(image)
26 print("correct label is: ", label)
27 print("predicted label is: ", prediction)

```

correct label is: 1
predicted label is: 1

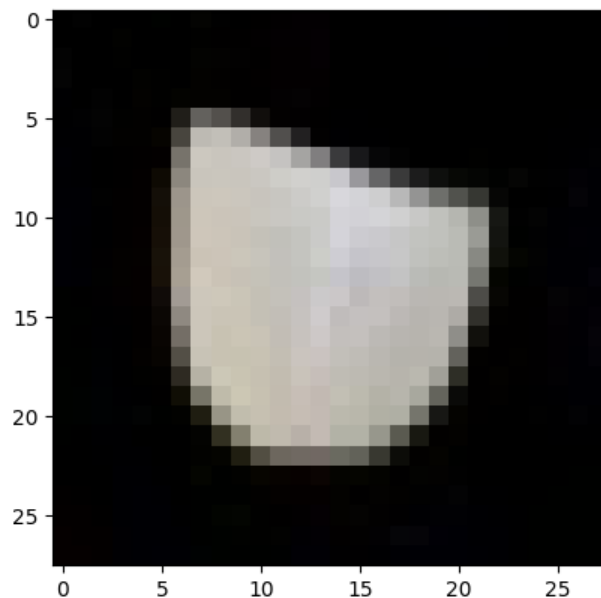


Figure 9: Seed image used for single test

With an area threshold of 350 pixels, our classifier was able to correctly classify this image.

5 Conclusions

Without using any sort of high-level machine learning or deep learning algorithms, we were able to manually create a very accurate rice seed classifier by using a simple threshold. If the image contained more white pixels than the threshold, it would be considered proper, and containing less than the threshold would be considered broken. The naive classifier performed quite well, and it would likely be able to classify the seeds perfectly just by adding a second parameter: the slope. Regardless, this is a testament to how simplicity can sometimes be just as useful as very complex models.

References

- [1] *Image analysis*, available at https://en.wikipedia.org/wiki/Image_analysis.