

# Bee Image Classification Using Convolutional Neural Networks

September 2023

## 1 Introduction

The objective of our machine learning project is to develop a model that can accurately distinguish images of bees from all other images. Images containing bees are labeled with 1, and images without bees are labeled with 0. This binary classification problem holds significant real-world implications as it has several applications in agriculture, environmental conservation, and beekeeping. For instance, accurately detecting the presence of bees can lead to increased crop yields and reduced pesticide risks, thereby benefiting both agriculture and ecosystems. By accurately classifying bee images, this project contributes to addressing critical ecological and agricultural challenges, benefiting both the environment and society.

The report is organized as follows: Section 2 describes the dataset, features, and labels. Section 3 presents the methodology, including data preprocessing and the selection of models. Section 4 discusses model comparisons, training, validation, and test results, and selects the final method. Finally, Section 5 concludes the report and suggests future improvements. The code can be found in Appendix A and references in Appendix B.

## 2 Problem Formulation

### 2.1 Data Points, Features, and Labels

This project is classified as supervised learning as the model is provided with both the images and their corresponding binary labels. During the training phase, the models learn from labeled data, identifying patterns and features that enable them to distinguish between bee and non-bee images.

- **Data Points:** We have a total of 30,000 images in our dataset, with a distribution of 10,000 bee images and 20,000 non-bee images. With each image represented by a 3D tensor of RGB values, the input is multi-dimensional.
- **Features:** The features are the pixel values of each image. The images depict either bees or various natural scenes that are common in images with bees, such as landscapes, flowers, animals, and other insects. We chose related images to help the algorithm distinguish between bees from their natural background (flowers & landscapes) as well as distinguish bees from other similar insects.
- **Labels:** The labels are binary, with a value of 1 representing bee images and 0 representing non-bee images, such as

## 3 Methods

### 3.1 Dataset Description and Preprocessing

The dataset used for this project has been collected from multiple sources including the Intel Image Classification dataset and The BeeImage Dataset. The full list is too long to be fully listed here. As such, to give due credit, all datasets used for this project, and links to them, are listed in Appendix B.

The datasets combined contain 30,000 images split between bee and non-bee images. Each image varies in dimensions and color. The preprocessing/cleaning steps include resizing images to a uniform size, normalizing pixel values, and performing in-place data augmentation with techniques like image rotation and color manipulation.

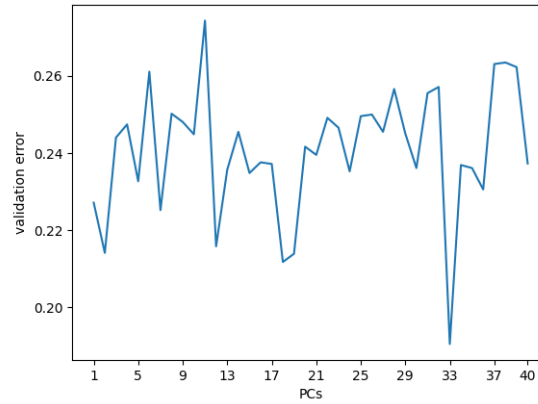
### 3.2 Feature Selection

In this project, the features are the pixel values of the images. We made a deliberate choice not to manually pick specific features because CNN models excel at feature extraction, which eliminates the need for manual feature engineering.

### 3.3 Convolutional Neural Network (CNN)

We chose Convolutional Neural Networks (CNNs) as the first model for this image classification task. CNNs are well-suited for image-related tasks, as they can autonomously learn to extract relevant features from raw pixel data through convolutional layers. Additionally, CNNs have the ability to capture spatial relationships within images. They are effective at recognizing patterns, objects, and features in images by considering how pixels relate to each other in terms of their positions and structures. This ability to analyze the spatial arrangement of features makes CNNs a natural choice for image classification.

### 3.4 Principle Component Analysis (PCA) On A Neural Network



Graph 1.1: A graph of validation error against the number of remaining Principal Components(PCs)

We chose Principal Component Analysis (PCA) in combination with a CNN as our second machine-learning model for image classification. Principal Component Analysis defines new attributes (principal components or PCs) as mutually orthogonal linear combinations of the original features.<sup>1</sup> Using PCA in the preprocessing step allowed us to significantly reduce the data's dimensionality while preserving most information. PCA, essentially, compresses the image, thereby removing noise from the image and allowing the CNN to focus on key details. This could improve the network's ability to differentiate between bee and non-bee images.

Choosing the correct number of principal components for the compressed image is challenging. Too few, and the image loses too much detail for the CNN to effectively distinguish the two types of images. Too many and the advantages of image compression become negligible. To solve this, we ran a simulation by training 14 different CNNs for each number of remaining PCs from 1 to 40 with a limited data set. This allowed us to observe any general trends in validation error with respect to the number of PCs. From graph 1.1 we observe that the validation error is the lowest for 33 PCs. Therefore, we chose 33 PCs for the final model.

### 3.5 Integration of PCA with CNNs

CNN is classified as a supervised machine learning model as the model is provided with both the images and their corresponding binary labels. On the other hand, PCA is an unsupervised technique, used for dimensionality reduction and feature extraction without any labeled data. Hence, PCA's integration with CNNs in our machine learning project may create ambiguity between supervised and unsupervised aspects. In this case, PCA serves as an unsupervised preprocessing step to reduce the dimensionality of the data before being used in a supervised learning context with CNNs. Therefore, both methods combined form a supervised learning model.

---

<sup>1</sup>Shereena&David (2015)

### 3.6 Choice of Loss Function

$$\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$

Where

$$y_i = \text{label}$$
$$\hat{y}_i = \text{prediction}$$

Since this classification problem only involves two classes, we use binary cross-entropy as a loss function for **both** models since it is uniquely suited for binary classification problems: binary cross-entropy is simple and fast to compute. Furthermore, the logarithm in the binary cross-entropy will undo the exponential behavior of the sigmoid activation function. Large gradients are important for gradient-descent algorithms to make sufficient progress in each iteration;<sup>2</sup> the logarithm will avoid the vanishing gradient problem, where, with large values, the gradient approaches 0 and learning slows down.

### 3.7 Model Validation

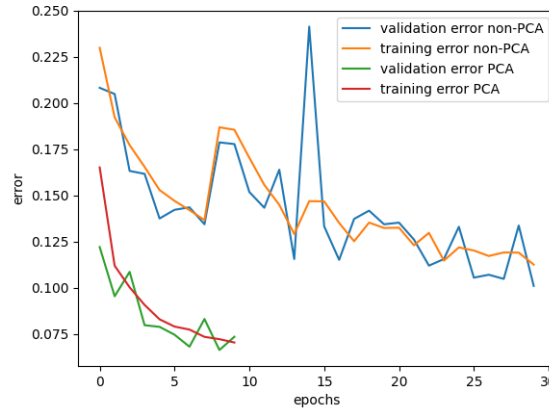
We have split the dataset into three subsets: training, validation, and testing sets. The data split is as follows:

- 70% of the data is allocated to the training set: This large training set is important for the model to learn the underlying patterns in the data.
- 15% of the data is allocated to the validation set: The validation set is used to evaluate the model's performance, tune the hyperparameters, and monitor for overfitting.
- 15% of the data is allocated to the testing set: The testing set is used to evaluate the model's final performance.

By plotting training and validation loss and comparing it with training and validation accuracy, we can look for trends such as overfitting (if validation loss starts increasing while training loss keeps decreasing) or underfitting (if both losses are high).

## 4 Results

### 4.1 Training and validation errors



Graph 2.1: A graph of error against epochs for PCA-based and pure CNNs

---

<sup>2</sup>oW\_

We compared the training and validation errors of each model across 10 epochs of training for the PCA-based CNN and 30 epochs of training for the pure CNN. Each epoch is a forward pass and backpropagation through the CNN with the entire training dataset. The training and validation error consistently decreased with the number of epochs as models learned from the data. From all epochs, we chose the versions with the lowest validation error for both the PCA-based CNN and CNN models. From Figure 1 we can see that the PCA-based CNN model reached the lowest validation error of 6.6% at epoch 9 while the pure CNN model reached the lowest validation error of 10.1% at epoch 30.

The reason for the pure CNN having 3 times more training epochs than the PCA-based CNN is due to hardware limitations: the machine on which both were trained required 2 hours to train the pure CNN through 30 epochs, while it required 5 hours to train the PCA-based CNN through only 10. This is due to applying the costly PCA compression algorithm to each of the approximately 25,000 images in each epoch. Due to time constraints, we could not train either for longer. Even so, the PCA-based CNN far outperformed the pure CNN model already at the second epoch in terms of validation error. Therefore, the final choice of model is the PCA-based CNN.

## 4.2 Test set

The test set used for our evaluation was constructed using augmented versions of the images that were not included in either the training or validation sets. To ensure a balanced representation of both bee and non-bee images, the test set was shuffled, creating an even distribution of these categories. Finally, we obtained a test error of 8.2% from the PCA-based CNN model. This shows the model is capable of accurately classifying images of bees from all other images.

## 5 Conclusion

In summary, our project aimed to develop two distinct models for image classification to differentiate between bee and non-bee images. We created both a Convolutional Neural Network (CNN) and a PCA-based CNN. The PCA-based CNN consistently performed better than the standard CNN, with a 8.2% test error rate. Training and validation errors for both models decreased at approximately the same rate and neither show signs of overfitting. Therefore, we can confidently select the PCA-based CNN as the superior model.

While the results are promising, there is room for improvement. Future directions may include exploring additional data augmentation techniques, experimenting with different CNN architectures, and fine-tuning hyperparameters to further enhance the model's performance. Additionally, collecting more diverse training data could contribute to even more robust classification capabilities.

In conclusion, our findings indicate that the PCA-based CNNs method is the most suitable for accurately classifying bee and non-bee images. This highlights the importance of integrating dimensionality reduction techniques like PCA with deep learning models for image classification.

## 6 Bibliography

- oW\_. (n.d.). What makes binary cross entropy a better choice for binary ...  
<https://datascience.stackexchange.com/questions/53400/what-makes-binary-cross-entropy-a-better-choice-for-binary-classification-than-o>
- Hussain, I. (2021). An ultimate guide to color image compression using PCA in python. Medium. Available at: <https://towardsdatascience.com/dimensionality-reduction-of-a-color-photo-splitting-into-rgb-channels-using-pca-algorithm-in-python-ba01580a1118>
- Shereena, V. B. & David, J. M. (2015, June). SIGNIFICANCE OF DIMENSIONALITY REDUCTION IN IMAGE PROCESSING . <https://core.ac.uk/download/pdf/212147876.pdf>

## 7 Appendix A code

```
1 from os.path import join
2 import datetime
3 from os import listdir
4 from os import getcwd
5 from os.path import isfile, splitext
6 import numpy as np
7 import keras
8 import cv2
9 from tensorflow.keras.preprocessing.image import ImageDataGenerator
10 from matplotlib.pyplot import imread
11 import matplotlib.pyplot as plt
12 import os
13 from sklearn.model_selection import train_test_split
14 import tensorflow as tf
15 from sklearn.metrics import confusion_matrix
16 import pandas as pd
17 import seaborn as sns
18
19
20 def display_conf_matrix (conf_mat):
21
22     plt.figure(figsize=(8,6), dpi=100)
23     sns.set(font_scale = 1.1)
24
25
26     ax = sns.heatmap(conf_matrix, annot=True, fmt='d')
27
28     ax.set_xlabel("Predicted Diagnosis", fontsize=14, labelpad=20)
29     ax.xaxis.set_ticklabels(['Negative', 'Positive'])
30
31     ax.set_ylabel("Actual Diagnosis", fontsize=14, labelpad=20)
32     ax.yaxis.set_ticklabels(['Negative', 'Positive'])
33
34     # set plot title
35     ax.set_title("Confusion Matrix for the Diabetes Detection Model", fontsize=14, pad=20)
36
37     plt.show()
38
39
40
41 # a method to correct an incorrect directory structure for the keras imagedatagenerator
42 def correct_directory_structure(train_percent, test_percent, validation_percent, *dpaths): #
43     dpaths: directories with classes of data
44     import random
45     from math import floor
46     cwd = getcwd() + '\\\\'
47     dir_list = listdir(cwd)
48     #check that we are not overwriting directories train,test,validation
49     assert not any (d in dir_list for d in ['train', 'test', 'validation']), 'cannot correct
50     directory structure, as some files will be overwritten'
```

```

49     assert all (d in dir_list for d in dpaths), 'the listed directories are not in the
        current path'
50
51     [os.mkdir(f"{i}") for i in ['train','test','validation']]
52     [os.mkdir(f"{i}\\{d}") for i in ['train','test','validation'] for d in dpaths]
53
54     data_gen = ((d,listdir(d)) for d in dpaths)
55     for d,i in data_gen:
56         shuffled_files = random.sample(i,len(i))
57
58         train_end = floor(len(shuffled_files)*train_percent)
59         test_end = floor(train_end + (len(shuffled_files) - train_end)*test_percent/(
            test_percent + validation_percent))
60
61         train_data = shuffled_files[0:train_end]
62         test_data = shuffled_files[train_end:test_end]
63         validation_data = shuffled_files[test_end:]
64
65         for m in train_data:
66             try:
67                 os.rename (f"{cwd}{d}\\{m}", f"{cwd}train\\{d}\\{m}")
68             except:
69                 print ("hit and error, moving forward")
70         for m in test_data:
71             try:
72                 os.rename (f"{cwd}{d}\\{m}", f"{cwd}test\\{d}\\{m}")
73             except:
74                 print ("hit and error, moving forward")
75         for m in validation_data:
76             try:
77                 os.rename (f"{cwd}{d}\\{m}", f"{cwd}validation\\{d}\\{m}")
78             except:
79                 print ("hit and error, moving forward")
80
81 image_shape = (400,400,3)
82
83 # returns the iterators for data in directories so the model can be trained on batches of
    data without needing much RAM
84 def get_image_iterators(train_dir = 'train',
85                         val_dir = 'validation',
86                         test_dir = 'test',
87                         mode = 'binary',
88                         batch = 64,
89                         size = image_shape[:2],
90                         processing = None):
91     cwd = getcwd() + '\\\\'
92     image_gen = ImageDataGenerator(
93         rotation_range = 90,
94         width_shift_range = 0.1,
95         height_shift_range = 0.1,
96         shear_range = 5,
97         brightness_range = (0.7,1.),
98         horizontal_flip = True,
99         vertical_flip = True,
100        preprocessing_function = processing
101    )
102    # these already shuffle the data
103    train_set = image_gen.flow_from_directory(train_dir, class_mode = mode, batch_size =
        batch, target_size = size)
104    val_set = image_gen.flow_from_directory(val_dir, class_mode = mode, batch_size = batch,
        target_size = size)
105    test_set = image_gen.flow_from_directory(test_dir, class_mode = mode, batch_size = batch
        , target_size = size)
106    return (train_set, val_set, test_set)
107
108 #correct_directory_structure(0.7,0.15,0.15, 'Bees', 'Not_Bees')
109
110 # for this function to work, the directory structure must be:

```

```

111 # train/
112 # |
113 # |         class_1/
114 # |         |         image1.jpg
115 # |         |         image2.jpg
116 # |         |         ...
117 # |
118 # |         class_2/
119 # |         |         image1.jpg
120 # |         |         image2.jpg
121 # |         |         ...
122 # test/
123 # |
124 # |         class_1/
125 # |         |         image1.jpg
126 # |         |         image2.jpg
127 # |         |         ...
128 # |
129 # |         class_2/
130 # |         |         image1.jpg
131 # |         |         image2.jpg
132 # |         |         ...
133 # ...
134
135 train, validation, test = get_image_iterators(batch = 32)
136
137 class SaverCallback (tf.keras.callbacks.Callback):
138     def __init__ (self, model, test_data, save_path, batch = 10): # test_data must be either
139         a keras sequence or dataset
140         self.best = test_data if type(test_data) == float else model.evaluate(test_data,
141         callbacks = [self], batch_size = batch)[1] #init val accuracy
142         self.save_path = save_path
143         self.model = model
144     def on_epoch_end (self, epoch, logs = None):
145         assert logs.get('val_accuracy') != None
146         if self.best < logs['val_accuracy']:
147             self.best = logs['val_accuracy']
148             self.model.save(self.save_path)
149
150 class accuracy_recorder (keras.callbacks.Callback):
151     def __init__ (self, init_compression = None):
152         self accuracies = pd.DataFrame (columns = ['compression', 'epoch', 'val_acc', '
153         val_loss', 'acc', 'loss'])
154         self.compression_factor = init_compression
155     def on_epoch_end (self, epoch, logs=None):
156         assert logs.get('val_accuracy') != None, 'no validation set for val_acc_recorder'
157         self accuracies = self accuracies._append({'compression': self.compression_factor,
158         'epoch' : epoch,
159         'val_acc': logs['val_accuracy'],
160         'val_loss': logs['val_loss'],
161         'acc': logs['accuracy'],
162         'loss': logs['loss']}, ignore_index = True)
163
164     def set_compression (self, new_compression):
165         assert new_compression > 0, 'cannot have pca with n_components = 0'
166         self.compression_factor = new_compression
167
168 model = tf.keras.Sequential([
169     tf.keras.layers.Conv2D(32, (3,3), padding='same', activation=tf.nn.relu,
170     rs.MaxPooling2D((2, 2), strides=2),
171     tf.keras.layers.Flatten(), input_shape=image_shape),
172     tf.keras.layers.MaxPooling2D((2, 2), strides=2),
173     tf.keras.layers.Conv2D(64, (3,3), padding='same', activation=tf.nn.relu),
174     tf.keras.layers.Dense(128, activation=tf.nn.relu),
175     tf.keras.layers.Dense(1, activation=tf.nn.sigmoid)
176 ])

```

```

176 model.compile (
177     optimizer = 'adam',
178     loss = 'binary_crossentropy',
179     metrics= ['accuracy']
180 )
181
182 CNN_model_path = join(getcwd(), 'models', 'CNNmodel')
183 saver = SaverCallback(model, 0., CNN_model_path)
184 recorder = accuracy_recorder()
185
186 model.fit(train,
187     validation_data = validation,
188     callbacks = [recorder, saver],
189     epochs = 30,
190     batch_size = 10,
191     initial_epoch = 20)
192
193 csv_data_path = join(getcwd(), 'CSVdata', 'CNNtraining2.csv')
194 recorder accuracies.to_csv(csv_data_path)
195 model.evaluate (test)
196
197 # let's plot the data
198 val_error = 1 - recorder.accuracies['val_acc']
199 error = 1 - recorder.accuracies['acc']
200 epochs = recorder.accuracies['epoch']
201 plt.plot(epochs, val_error, label='validation error non-PCA')
202 plt.plot (epochs, error, label = 'training error non-PCA')
203 plt.xlabel('epochs')
204 plt.ylabel('error')
205 plt.legend()
206 plt.show()
207
208 import matplotlib.pyplot as plt
209 from sklearn.decomposition import PCA
210 from scipy.stats import stats
211 import matplotlib.image as mpimg
212
213 compression = 33
214 def compress_images(img):
215
216     # Split into channels
217     blue, green, red = cv2.split(img)
218
219     # Scale the data between 0 and 1 for all channels
220     r = red / 255
221     b = blue / 255
222     g = green / 255
223
224     # Fit and transform data in PCA to reduce dimensionality
225     red_pca = PCA(n_components=compression)
226     red_pca.fit(r)
227     red_trans = red_pca.transform(r)
228
229     blue_pca = PCA(n_components=compression)
230     blue_pca.fit(b)
231     blue_trans = blue_pca.transform(b)
232
233     green_pca = PCA(n_components=compression)
234     green_pca.fit(g)
235     green_trans = green_pca.transform(g)
236
237     # Reconstruct the images
238     r_arr = red_pca.inverse_transform(red_trans)
239     b_arr = blue_pca.inverse_transform(blue_trans)
240     g_arr = green_pca.inverse_transform(green_trans)
241
242     # Merge channels into one
243     compressed_image = (cv2.merge((b_arr, g_arr, r_arr)))

```



```

244         return np.array(compressed_image)
245
246
247 modelPCA = tf.keras.Sequential([
248     tf.keras.layers.Conv2D(32, (3,3), padding='same', activation=tf.nn.relu,
249         input_shape=image_shape),
250     tf.keras.layers.MaxPooling2D((2, 2), strides=2),
251     tf.keras.layers.Conv2D(64, (3,3), padding='same', activation=tf.nn.relu),
252     tf.keras.layers.MaxPooling2D((2, 2), strides=2),
253     tf.keras.layers.Flatten(),
254     tf.keras.layers.Dense(128, activation=tf.nn.relu),
255     tf.keras.layers.Dense(1, activation=tf.nn.sigmoid)
256 ])
257
258 modelPCA.compile (
259     optimizer = 'adam',
260     loss = 'binary_crossentropy',
261     metrics= ['accuracy']
262 )
263
264 train_cmp, validation_cmp, test_cmp = get_image_iterators(processing = compress_images)
265 PCA_model_path = join(getcwd(), 'models', 'PCAmodel')
266 saver = SaverCallback(modelPCA, 0., PCA_model_path)
267 recorderPCA = accuracy_recorder()
268 modelPCA.fit(train_cmp,
269     validation_data = validation_cmp,
270     epochs = 40,
271     callbacks = [saver, recorderPCA])
272
273
274 train_cmp, validation_cmp, test_cmp = get_image_iterators(
275     processing = compress_images, batch = 8)
276
277 from functools import reduce
278 stack_features_and_labels = lambda fl1, fl2: (np.vstack((fl1[0], fl2[0])), np.hstack((fl1
279     [1], fl2[1])))
280
281 # take some data and load it into ram so it's fast to train with it
282 def limit_data (data, included_batches): # requires an iterator
283     holder = next(data)
284     for i in range (1, included_batches):
285         holder = stack_features_and_labels (holder, next(data))
286     return holder
287
288 limited_train_data = limit_data(train_cmp,3)
289 limited_test_data = limit_data(test_cmp,1)
290 limited_validation_data = limit_data(validation_cmp,3)
291
292 trend_tester_model = tf.keras.Sequential([
293     tf.keras.layers.Conv2D(32, (3,3), padding='same', activation=tf.nn.relu,
294         input_shape=image_shape),
295     tf.keras.layers.MaxPooling2D((2, 2), strides=2),
296     tf.keras.layers.Conv2D(64, (3,3), padding='same', activation=tf.nn.relu),
297     tf.keras.layers.MaxPooling2D((2, 2), strides=2),
298     tf.keras.layers.Flatten(),
299     tf.keras.layers.Dense(128, activation=tf.nn.relu),
300     tf.keras.layers.Dense(1, activation=tf.nn.sigmoid)
301 ])
302
303 trend_tester_model.compile (
304     optimizer = 'adam',
305     loss = 'binary_crossentropy',
306     metrics= ['accuracy']
307 )
308 validation accuracies = accuracy_recorder()
309 save_path = join(getcwd(), 'trend_tester_temp')
310 trend_tester_model.save_weights(save_path)

```

```

311 #what compression we start at
312 start_compression = 40
313 #how many times we train a new model under one compression
314 reps = 14
315
316 for i in range (start_compression, 0, -1):
317     validation_accuracies.set_compression(i)
318     compression = i
319     for _ in range(0, reps):
320         trend_tester_model.load_weights(save_path)
321         trend_tester_model.fit(limited_train_data[0],
322                               limited_train_data[1],
323                               validation_data = limited_validation_data,
324                               batch_size= 8,
325                               epochs = 20,
326                               callbacks = [validation_accuracies])
327
328 csv_data_path = join(getcwd(), 'CSVdata', 'trend_tester4.csv')
329 validation_accuracies.to_csv(csv_data_path)
330 dat = pd.read_csv(join(getcwd(), 'CSVdata', 'trend_tester4.csv'))._append(pd.read_csv(join(
331     getcwd(), 'CSVdata', 'trend_tester3.csv')))
332 dat.groupby('compression')['val_acc'].mean().plot()

```

## 8 Appendix B datasets

- Bhathena, Jehan. “Weather Image Recognition.” Kaggle, 26 Nov. 2021, [www.kaggle.com/datasets/jehanbhathena/weather-dataset](https://www.kaggle.com/datasets/jehanbhathena/weather-dataset)
- DePie. “Butterfly Image Classification.” Kaggle, 24 June 2023, [www.kaggle.com/datasets/phuchthai02/butterfly-image-classification](https://www.kaggle.com/datasets/phuchthai02/butterfly-image-classification)
- Gerry. “Birds 525 Species- Image Classification.” Kaggle, 20 Apr. 2023, [www.kaggle.com/datasets/gpiosenska/100-bird-species](https://www.kaggle.com/datasets/gpiosenska/100-bird-species)
- Rey, George. ”Bee or wasp?” Kaggle, <https://www.kaggle.com/datasets/jerzydziewierz/bee-vs-wasp>
- Yang, Jenny. “The Beeimage Dataset: Annotated Honey Bee Images.” Kaggle, 16 Sept. 2018, [www.kaggle.com/datasets/jenny18/honey-bee-annotated-images](https://www.kaggle.com/datasets/jenny18/honey-bee-annotated-images)