

Building Data Sets Using k-Means Clustering and Investigation of Training Accuracy in Convolutional Neural Networks

Yuichi MIYATA Yoko UWATE Yoshifumi NISHIO
(Tokushima University)

1. Introduction

In recent years, aerial photography became easier than before by using the camera loaded in a drone. Also, convolutional neural network (CNN) is one of deep learning and is the network often used for image recognition. With the development of CNN, drones are being researched for applications in various fields such as agriculture. Actually, wild animals such as deer and boars are rapidly increasing in Japanese forests. Agricultural crops damaged in nearby farms. By using a drone, we expected to wild animals management. However, the drone camera's battery and memory are limited. It is important to use the collected data effectively.

In this study, we used clustering to make more efficient data sets from the collected image data.

2. Proposed method

Clustering is a method of classifying data having similar features. Therefore, we can classify images which have similar characteristics such as color and shape by using clustering. This time, we use k-means clustering to classify the images into two clusters and constructed data sets. The images for data sets are collected by entering a general name into the search word in the image search of google. We prepare 400 images for training and test. There are animals in 200 images among 400 images. Other than that, they are images of the background. Furthermore, we convert 3D array of input image data is converted to a 2D array data.

The flow of learning by k-means clustering and the objective function are described as follows:

$$f(\{C_k\}) = \sum_{k=1}^k \sum_{x_i \in C_k} (\bar{x}_k - x_i)^2. \quad (1)$$

After making the data sets by k-means clustering, CNN learns them. We define as the learning steps are 300. The learning rate of CNN is 0.00001. When it learns by CNN, the input images (432 × 324 pixels) is resized to 28 × 28 pixels.

3. Simulation results

We classified each of the images of the animal and the images of the background into 2 classes using k-means clustering. Images of the animal are divided into 134 and 66 images. Images of the background are divided into 109 and 91 images. One with a large number of clustered images is referred to as cluster A, and the other with a small number as cluster B. We construct the following data set to compare the effects of clustering.

[Data set1] : Only before clustering (400 images)
[Data set2] : Randomly extracted (200 images)
[Data set3] : Cluster A (243 images)
[Data set4] : Cluster B (157 images)
[Data set5] : Cluster A 50% + Cluster B 50% (200 images)

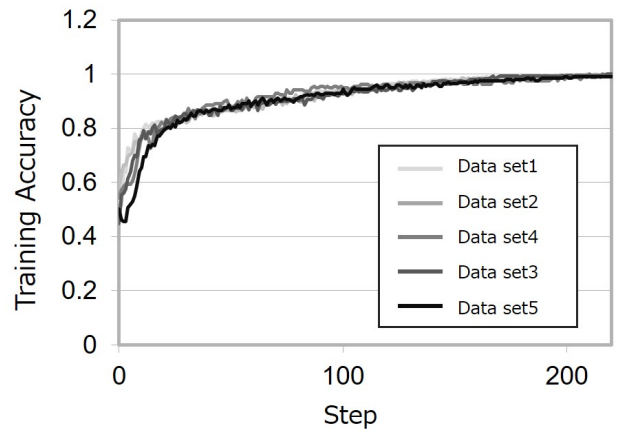


Figure 1: Training accuracy of each data sets.

In Fig. 1, we show the training accuracy and step of each data set on CNN. Each data sets have different number of images. We found that the influence of training accuracy was small even if there were few images for learning.

Table 1: Average of training and test accuracies.

Data sets	1	2	3	4	5
Training accuracy [%]	100	100	100	100	100
Test accuracy [%]	96	88	96.5	66	93.5

From Table 1, the test accuracy of data set3 is only a little bit better than the data set1. When we compared data set2 and data set5, it increased the test accuracy of data set5 with the same number of images.

4. Conclusion

We used k-means clustering to construct data sets that achieve high accuracy with few images. We classified images into two classes for each feature by k-means clustering. From these images, we constructed 5 data sets and trained using CNN. We compared with training accuracy and investigated the effect of clustered data sets. In simulation results, the proposed method obtains a little bit better test accuracy than the conventional method.

In future works, we will investigate the effective data set for learning by increasing the number of clusters.