

# *Research on false eyewitness Detection algorithm of Asian giant hornet Image based on support Vector Machine*

**Jiaxu Li**

*College of Physics and Optoelectronic Engineering, Shenzhen University Shenzhen City, Guangdong Province 518000*

**Keywords:** Asian giant hornet, sample imbalance, support vector machine, data enhancement

**Abstract:** Recently, *Vespa mandarinia* have been found frequently in Washington state, causing severe potential damage to the local ecosystem. To establish effective pest management programs, Washington State collected reports of people witnessing these wasps. However, due to the existence of a large number of wrong reports and the limited resources of government agencies, it is not possible to conduct field visits to all reports. In recognition of report images, since the image data set is extremely unbalanced, we use data enhancement, cutting and rotating the positive images to increase the number of positive images. In addition, due to the fact that the traditional neural network is easy to overfit in this kind of imbalanced data set, we use the One-Class SVM model to transform the classification problem into the outlier test problem and the experimental results show that the accuracy of our algorithm is 98%.

## **1. Introduction**

An invasive species is an introduced organism that negatively alters its new environment. [1] Although their spread can have beneficial aspects, invasive species adversely affect the invaded habitats and bioregions, causing ecological, environmental, and economic damage. [2]

The *Vespa mandarinia*, also known as the Asian giant hornet, is indeed an invasive species to the Washington State for the reason that it is a predator to the native European honeybees. The hornet is originated in the temperate and tropical eastern Asia, including parts of Japan, China, India and Sri Lanka. Therefore, it is unlikely to occur outside of Washington state and Vancouver Island. Although the hornets are voracious predators of other insects that are considered agricultural pests, small amount of them are capable of devastating a whole colony of European honeybees in a short time leading to the imbalance of ecosystem. Such pest must be eradicated before it spread too far.

In order to prevent the hornets from spreading and maintain the balance of the ecology system in Washington State, accurate identification is conducted by the governors. Civilians are encouraged to report the sightings of the hornets with digital evidence like images and videos.

There are a lot of misjudgments in the information uploaded by eyewitnesses. Therefore, in this paper, we establish a general prediction model, which only uses the possibility of misclassification of the given data file output, so as to automatically judge the authenticity of the reported sightings.

## 2. Data Preprocessing

The given image data set uploaded by people in Washington State has the characteristics of imbalance between positive and negative examples and scarcity, which hinders model training using machine learning methods. We take the following steps to sanitize the data set:

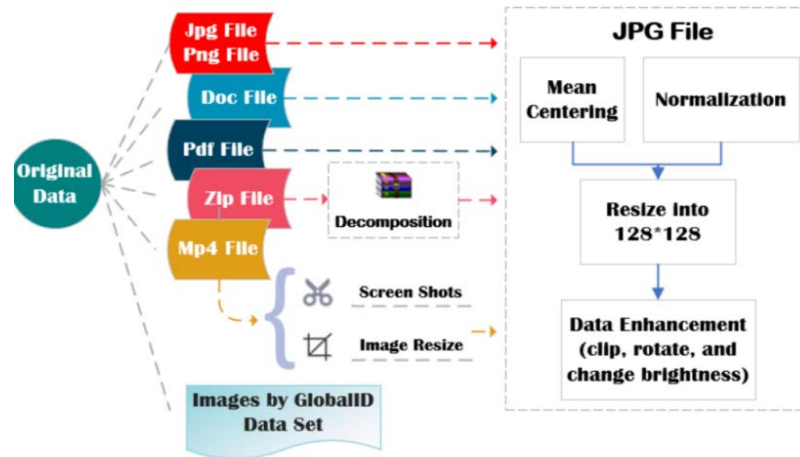


Figure 1: Data Preprocessing Flow Chart

**Step 1** For files in doc, pdf, and compressed package formats in the data set, extract the pictures in the files and save them in jpg format. Three screenshots, which are relatively clear and sized wasp pictures, are manually taken for video format files.

**Step 2** Every picture is preprocessing using mean centering and normalization method. In order to facilitate further processing, pictures are sized into 128\*128.

**Step 3** Using the method of data enhancement, we clip, rotate, and change the brightness and contrast of the "positive" and "unverified" images to make all kinds of data as close as possible to keep the data set balanced.

## 3. Support Vector Domain Description (SVDD)

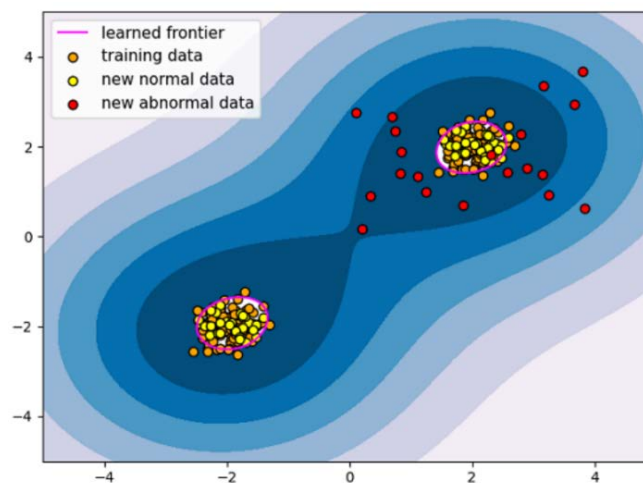


Figure 2: The framework of Support Vector Domain Description (SVDD)

Considering the current situation of extremely unbalanced data sets, we can look at the problem from a completely different perspective. We can transform the original problem into a One-Class Learning or Novelty Detection issue. Such methods only model one of those categories rather than paying attention to extracting the difference between classes.

We propose the classical One-Class SVM [3], whose essential algorithm is Support Vector Domain Description (SVDD), for solving the One Class Learning or Novelty Detection problems. In this way, our method can greatly eliminate the impact of data imbalance.

### 3.1 Introduction to SVDD

The basic idea of SVDD is to construct a minimum hypersphere (Hypersphere refers to the sphere in more than three-dimensional space. The corresponding two-dimensional space is the surface, and the three-dimensional space is the sphere.) Since there is only one class, the minimum hypersphere should consist of all the data that belong to the Negative ID. When classifying a new data, we only need to confirm that whether it's in the hypersphere (As shown in figure2).

The optimization objective of SVDD is to find a minimum sphere with center  $a$  and radius  $R$ :

$$F(R, a, \xi_i) = R^2 + C \sum_i \xi_i \quad (1)$$

Let the sphere satisfy:

$$(x_i - a)^T(x_i - a) \leq R^2 + \xi_i \quad \forall_i, \xi_i \geq 0 \quad (2)$$

Where  $\xi_i$  is a slack variable, whose function is to prevent the model from being destroyed by individual extreme data points. When most of the data points are concentrated in a small area along with a few data points which is far away from the former points, the whole hypersphere will become huge because of the noise data, which makes the model sensitive to outliers and finally lead to overfitting. Hence, the model needs to tolerate some data points that do not satisfy the constraints, give them some flexibility, and at the same time, we should cover every data point in the training set to meet the constraints so that we can solve them using Lagrange Multiplier later. Mention that the size of the slack variable is related to each data point.

Now we apply the Lagrange Multiplier:

$$L(R, a, \alpha_i, \xi_i) = R^2 + C \sum_i \xi_i - \sum_i \alpha_i \{R^2 + \xi_i - (x_i^2 - 2ax_i + a^2)\} - \sum_i \gamma_i \xi_i \quad (3)$$

Mention that  $\alpha_i \geq 0$  and  $\gamma_i \geq 0$ , the following formula is obtained by deriving the parameter and making the derivative equal to 0:

$$\begin{aligned} \sum_i \alpha_i &= 1, a = \frac{\sum_i \alpha_i x_i}{\sum_i \alpha_i} = \sum_i \alpha_i x_i \\ C - \alpha_i - \gamma_i &= 0 \quad \forall_i \end{aligned} \quad (4)$$

Formula 4 is replaced by Lagrangian function:

$$L = \sum_i \alpha_i (x_i \cdot x_i) - \sum_{i,j} \alpha_i \alpha_j (x_i \cdot x_j) \quad (5)$$

Note that  $0 \leq \alpha_i \leq C$ ,  $\sum_i \alpha_i = 1$ , among which  $0 \leq \alpha_i \leq C$ ,  $\sum_i \alpha_i = 1$  is derived from  $\alpha_i \geq 0$ ,  $\gamma_i \geq 0$  and  $C - \alpha_i - \gamma_i = 0$ . The vector inner product in formula 1 can also be solved by kernel

function like SVM:

$$L = \sum_i \alpha_i K(x_i, x_i) - \sum_{i,j} \alpha_i \alpha_j K(x_i, x_j) \quad (6)$$

Formula 5 can be solved like SVM. After the training, judge whether a new data point  $Z$  is of this class, then it depends on whether the data point is in the trained hypersphere. If it is in the hypersphere, that is  $(z - a)^T(z - a) \leq R^2$ , it is determined to belong to this class, and the center of the hypersphere is represented by the support vector. Then the judgment condition for judging whether the new data belongs to this class is as follows:

$$(z \cdot z) - 2 \sum_i \alpha_i (z \cdot x_i) + \sum_{i,j} \alpha_i \alpha_j (x_i \cdot x_j) \leq R^2 \quad (7)$$

If kernel function is applied, then formula 3 can be rewritten as follow:

$$K(z, z) - 2 \sum_i \alpha_i K(z, x_i) + \sum_{i,j} \alpha_i \alpha_j K(x_i, x_j) \leq R^2 \quad (8)$$

### 3.2 Experiments based on SVDD

In our experiments, due to the extreme imbalance between the Negative ID and Positive ID, we only take the photos which belong to the Negative ID as the training set. Correspondingly, any other data (like Positive ID or Unverified) will be taken as an outlier.

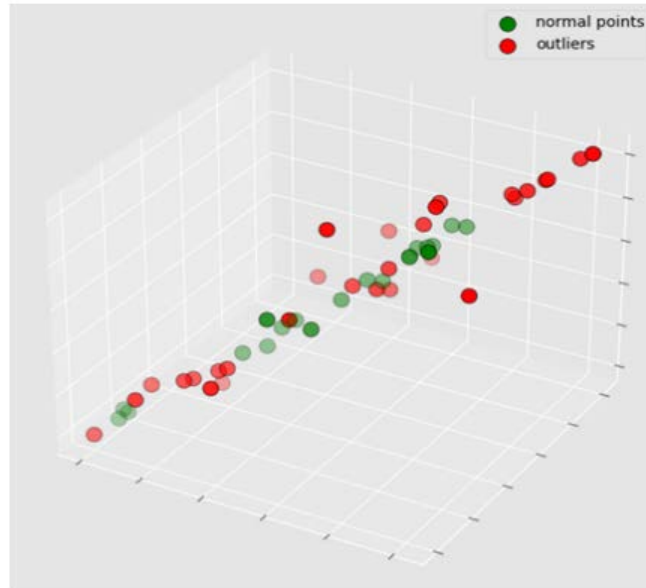


Figure 3: The effect chart after PCA dimension reduction

In our experiments, due to the extreme imbalance between the Negative ID and Positive ID, we only take the photos which belong to the Negative ID as the training set. Correspondingly, any other data (like Positive ID or Unverified) will be taken as an outlier [4].

$$k(x, x') = e^{-\frac{\|x-x'\|^2}{2\sigma^2}} \quad (9)$$

### 3.3 Model Results on the Image Classification Task

We develop a test set consists of 20 Negative ID photos, 10 Positive ID photos and 20 Unverified ID photos. The final accuracy will be calculated by the following formula:

$$\text{Accuracy} = \frac{N + O_1 + O_2}{50}$$

Where N stands for Num of Negative ID,  $O_1$  stands for Num of Positive ID,  $O_2$  stands for Num of Unverified ID.

Repeating the test procedure for five times, we get the final accuracy equal to 98%, which means the imbalanced problem of the data set is addressed successfully.

According to the SVDD method, data points are conducted in a high dimensional space and can not be visualized directly, we apply the Principle Component Analysis (PCA) for reducing the demension. Figure 3 shows that the normal points are approximately concentrated in the middle of the plot while the outliers are separated on the edge of the picture.

### 4. Conclusion

In this paper, in order to solve the problem that the image error rate provided by Asian giant hornet eyewitnesses is too high, we introduce the support vector machine model to classify the wrong recognition images. The experimental results show that the algorithm proposed in this paper still achieves the accuracy of 98% in the case of uneven data.

### References

- [1] Mark A Davis and Ken Thompson. *Eight ways to be a colonizer; two ways to be an invader: a proposed nomenclature scheme for invasion ecology*. *Bulletin of the ecological society of America*, 81(3):226–230, 2000.
- [2] J. Ehrenfeld. *Ecosystem consequences of biological invasions*. *Annual Review of Ecology, Evolution, and Systematics*, 41:59–80, 2010.
- [3] Yunqiang Chen, Xiang Sean Zhou, and Thomas S Huang. *One-class svm for learning in image retrieval*. In *Proceedings 2001 International Conference on Image Processing (Cat.No. 01CH37205)*, volume 1, pages 34–37. IEEE, 2001.
- [4] Florent Perronnin, Jorge Sánchez, and Thomas Mensink. *Improving the fisher kernel for large-scale image classification*. In *European conference on computer vision*, pages 143–156. Springer, 2010.