

Pen Spinning Hand Movement Analysis Using MediaPipe Hands

Tung-Lin “Dove” Wu¹ Taishi “Wabi” Senda²

¹Taiwanese Pen Spinning Forum (TWPS)

²JapEn Board (JEB)

Emails: howard0100000@gmail.com, wabi.penspinning@gmail.com

Abstract

We challenged to get data about hand movement in pen spinning using MediaPipe Hands [1] and OpenCV [2]. The purpose is to create a system that can be used to objectively evaluate the performance of pen spinning competitions. Evaluation of execution, smoothness, and control in competitions are quite difficult and often with subjectivity. Therefore, we aimed to fully automate the process by using objective numerical values for evaluation.

Uncertainty still exists in MediaPipe’s skeletal recognition, and it tends to be more difficult to recognize in brightly colored backgrounds. However, we could improve the recognition accuracy by changing the saturation and brightness in the program. Furthermore, automatic detection and adjustment of brightness is now possible.

As the next step to systematize the evaluation of pen spinning using objective numerical values, we adopted “hand movements”. We were able to visualize the ups and downs of the hand movements by calculating the standard deviation and L2 norm of the hand’s coordinates in each frame. The results of hand movements are quite accurate, and we feel that it is a big step toward our goal. In the future, we would like to make great efforts to fully automate the grading of pen spinning.

Keywords: L2 norm, Standard deviation, OpenCV, Computer vision

1 Introduction

Execution and control play a big role in pen spinning, and it is often an important thing to examine whether a pen spinner could control the pen well or not by analyzing hand movement of the pen spinning video. However, pen spinning community currently does not have an objective way to tell how big the spinner's hand movement is. We want to resolve this problem by evaluating the hand movement of pen spinning videos without any subjectivity and bias, so we need to get the coordinate of the hand of each frame in the video and analyze this data. The analyzing results should not be affected by the camera angle and the distance between camera and hand.

We first used MediaPipe Hands [1] made by Google to get the coordinates of the hand of the video. Details of the hand landmarks are in figure 1. The origin of the vector space is the top left corner of the frame (see figure 2). The coordinate of each landmark is the pixel it located. Thus the coordinates of two same video would not be the same if their resolutions are different. If we zoom in/out the video or rotate the video, the coordinate of the data would also change. In section 3, we would show how to resolve the scaling problem by normalizing the data, and how to deal with the rotation problems. Then in section 4, we designed two main methods to calculate the final result: L2 norm and standard deviation.

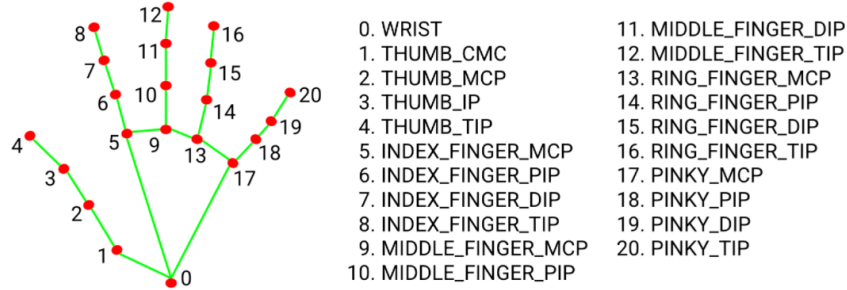


Figure 1: Hand landmarks of MediaPipe Hands

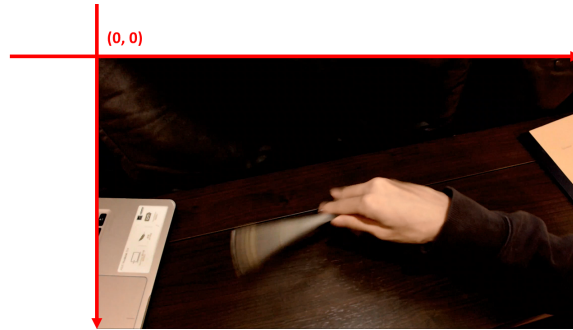


Figure 2: Coordinate system [3]

2 Data Preprocessing

Cutting out the parts other than pen spinning and adjusting the brightness and saturation of the video are necessary before analyzing the video. The data would have a lot of noise if we did not consider these problems.

The reason why we need to get rid of the parts that are not associative with pen spinning first is because that our research is based on skeletal estimation. If we involved coordinates of non-spinning parts, the results would be inaccurate.

The next step is to change the brightness and saturation by using OpenCV [2]. In figure 3a and 3b, skeletal estimation is very unstable when the background of the pen spinning video is bright. The blue dots in the top left corner of figure 3a are all coordinates that are not belong to any part of the hand. This situation make the graph in figure 3b become really noisy, hence loses its reliability. In order to improve this situation, we made it possible to adjust the brightness and saturation automatically in the code. Figure 3c and 3d are the result after the adjustment. We successfully removed most of the noise in figure 3c, thus figure 3d could reflect the real result of the original video.

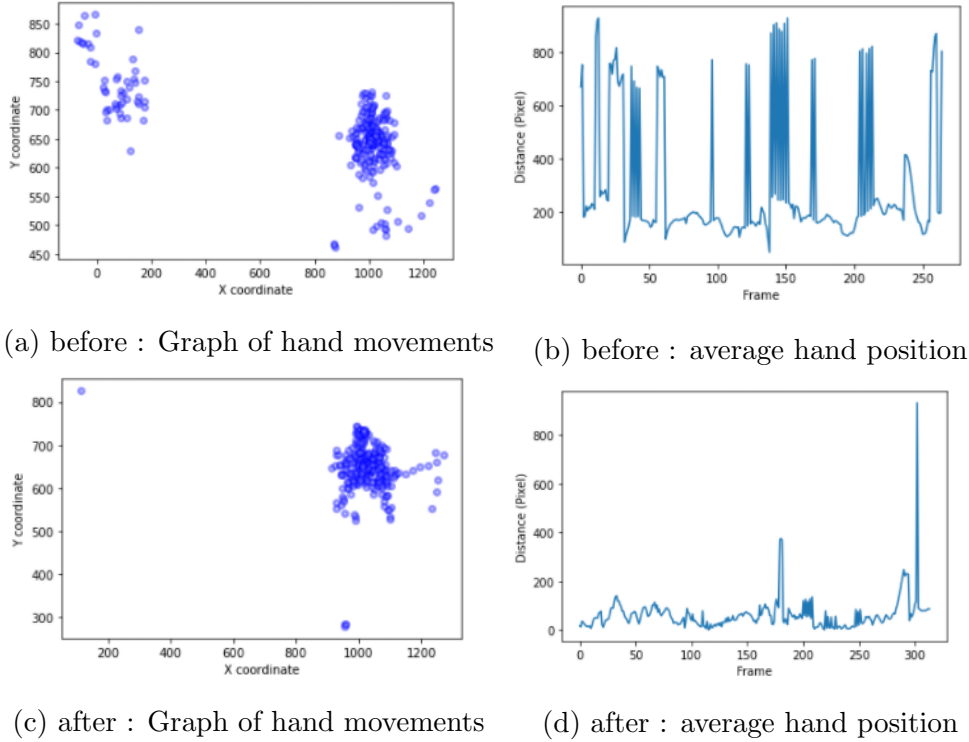


Figure 3: fukrou for JapEn 15th [4]

3 Scaling and Rotation problems

The same video should always have the same result even if there exists some differences like scaling and rotation. We divided the scaling problem into two parts: resolution and zoom in/out. Since the coordinate of the hand is depend on which pixel the hand is, the coordinate would be different if we did not consider the resolution and zoom in/out problems. For example, figure 4a, 4b, and 4c are exactly the same video, but one is 720p, the other is 1080p, and the other is zoom out a little bit. In order to make these three videos have the same results, we decided to calculate the hand size first and use the hand size to normalize the result. The formula of how we get the hand size is shown below,

$$s = \frac{1}{n} \sum_{i=1}^n \sqrt{(x_{i,0} - x_{i,5})^2 + (y_{i,0} - y_{i,5})^2} \quad (1)$$

where $x_{i,j}, y_{i,j}$ ($i = 1, 2, \dots, n$) are the coordinates of ith frame of landmark j (see figure 1).

The rotation problem is really important because every spinner have their unique angle setup. Our method need to make sure that the result of the original video (figure 4b) should be the same with the rotated video (figure 4d) if they are exactly the same combo. We will discuss the details of how we resolved this situation in section 4.

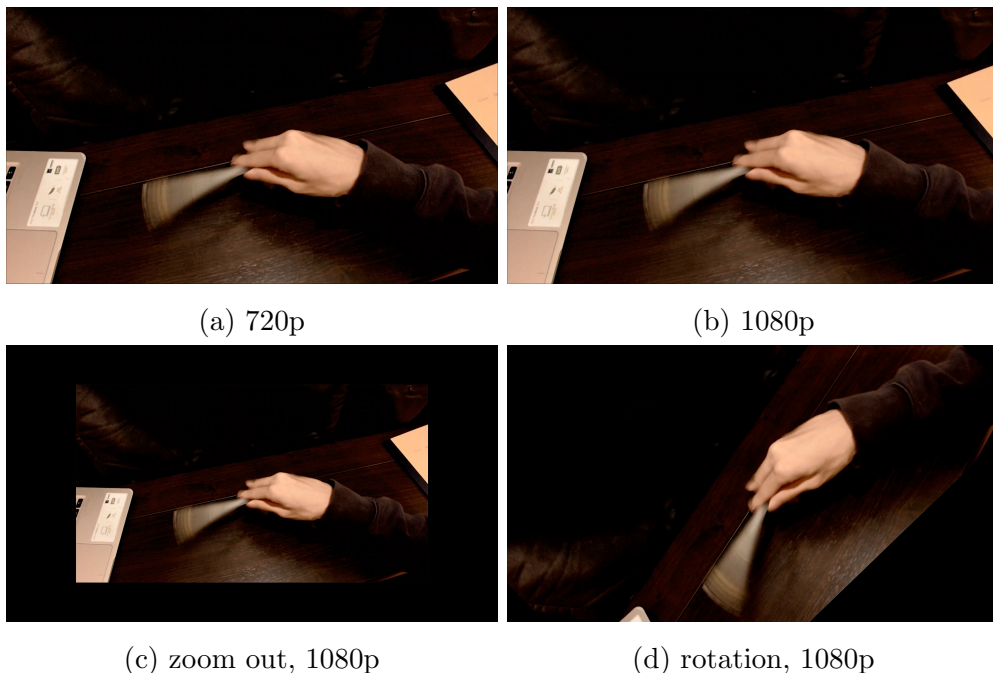


Figure 4: Dove for PSO20 R1 [3]

Other factors like reflection, shifting, and FPS do not affect our calculations. Thus these would not be the issues while computing the results.

4 Methods

4.1 L2 Norm

L2 norm is a method to calculate the distance. Since any rotation of the data does not change the distance between each coordinate, there is no need to worry about it. In this method, we used landmark 0's coordinate (see figure 1) of each frame of the video to be our data. The order of calculations we did is as follows. First, we averaged the coordinates of the data to get the mean position of the hand (see equation 2),

$$\begin{aligned}\bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i \\ \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i\end{aligned}\tag{2}$$

where x_i, y_i ($i = 1, 2, \dots, n$) are the coordinates of i th frame of landmark 0. Then, we derived the average distance between (\bar{x}, \bar{y}) and the coordinate of landmark 0 for each frame. In order to make the derived values correspond to the scaling of the video, we need to normalized the data by the hand size s of the video. The details about how to get s is in equation 1. Hence we resolved the problem by dividing the data by s . The whole process is in equation 3,

$$\begin{aligned}d &= \frac{1}{n} \sum_{i=1}^n \sqrt{\left(\frac{1}{s}x_i - \frac{1}{s}\bar{x}\right)^2 + \left(\frac{1}{s}y_i - \frac{1}{s}\bar{y}\right)^2} \\ &= \frac{1}{n \cdot s} \sum_{i=1}^n \sqrt{(x_i - \bar{x})^2 + (y_i - \bar{y})^2}\end{aligned}\tag{3}$$

where s is the hand size in equation 1, x_i, y_i ($i = 1, 2, \dots, n$) are the coordinates of i th frame of landmark 0.

This method is trying to measure the distance of hand movement in each frame. Furthermore, normalizing the results ensures that every pen spinning video is in the same standard. By comparing the results with other video, we could know that whether certain moves are large or small.

4.2 Standard Deviation

Calculating the standard deviation is always a good method to know the unstable rate of the data. However, it only consider one dimension each time. Since we used landmark 0's coordinate (see figure 1) of each frame of the video to be our data, we need to calculate the standard deviation, σ_x, σ_y , of x axis and y axis respectively and designed a function to combine those two values. The most important thing for this function is that we need to make sure the results of figure 4a, 4b, 4c, 4d should be the same. For the scaling problem

(figure 4a, 4b, 4c), we only need to divide the original coordinates by the hand size. The formulas are written below,

$$\begin{aligned}\sigma_x &= \sqrt{\frac{\sum_{i=1}^n (\frac{1}{s}x_i - \frac{1}{s}\bar{x})^2}{n}} = \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}} \\ \sigma_y &= \sqrt{\frac{\sum_{i=1}^n (\frac{1}{s}y_i - \frac{1}{s}\bar{y})^2}{n}} = \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n}}\end{aligned}\tag{4}$$

where s is the hand size in equation 1, x_i, y_i ($i = 1, 2, \dots, n$) are the coordinates of landmark 0 of i th frame, \bar{x}, \bar{y} are the mean of all x and y of landmark 0 in equation 2. By doing so, the data is normalized by the hand size s . Hence the scaling problem is resolved.

The rotation problem is more challenging while using standard deviation as the method. The data rotates about any coordinate can be treated as rotate about the origin with some shift in x and y axis. However, shifting does not change the result of σ_x and σ_y , thus we only need to consider the situation that rotate about the origin.

The rotation matrix is shown below,

$$\begin{bmatrix} x'_i \\ y'_i \end{bmatrix} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x_i \\ y_i \end{bmatrix}\tag{5}$$

where x'_i, y'_i ($i = 1, 2, \dots, n$) are the coordinates x_i, y_i of landmark 0 of i th frame rotated by angle θ . After the rotation, the new standard deviation $\sigma_{x'}, \sigma_{y'}$ become

$$\begin{aligned}\sigma_{x'} &= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n (x'_i - \bar{x}')^2}{n}} \\ &= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n ((x_i \cos \theta - y_i \sin \theta) - (\bar{x} \cos \theta - \bar{y} \sin \theta))^2}{n}} \\ &= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n ((x_i - \bar{x}) \cos \theta - (y_i - \bar{y}) \sin \theta)^2}{n}} \\ \sigma_{y'} &= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n (y'_i - \bar{y}')^2}{n}} \\ &= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n ((x_i \sin \theta + y_i \cos \theta) - (\bar{x} \sin \theta + \bar{y} \cos \theta))^2}{n}} \\ &= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n ((x_i - \bar{x}) \sin \theta + (y_i - \bar{y}) \cos \theta)^2}{n}},\end{aligned}\tag{6}$$

where s is the hand size in equation 1, \bar{x}' and \bar{y}' are the mean of all x' and y' of landmark 0 in equation 2, . The goal is trying to find a function $f(\sigma_{x'}, \sigma_{y'})$ such that

$$f(\sigma_{x'}, \sigma_{y'}) = f(\sigma_x, \sigma_y)\tag{7}$$

After some calculations, we found out that equation 7 holds if $f(\sigma_x, \sigma_y) = \sqrt{\sigma_x^2 + \sigma_y^2}$. The

proof is shown below. Let $A = (x_i - \bar{x})$, $B = (y_i - \bar{y})$ in equation 6.

$$\begin{aligned}
\sqrt{\sigma_{x'}^2 + \sigma_{y'}^2} &= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n ((A \cos \theta - B \sin \theta)^2 + (A \sin \theta + B \cos \theta)^2)}{n}} \\
&= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n ((A^2(\cos^2 \theta + \sin^2 \theta) + B^2(\cos^2 \theta + \sin^2 \theta))}{n}} \\
&= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n (A^2 + B^2)}{n}} \\
&= \frac{1}{s} \sqrt{\frac{\sum_{i=1}^n ((x_i - \bar{x})^2 + (y_i - \bar{y})^2)}{n}} \\
&= \sqrt{\sigma_x^2 + \sigma_y^2}
\end{aligned} \tag{8}$$

Hence the results of figure 4b, 4d would be the same by using equation 8.

5 Results

The biggest difference between L2 norm (5a) and standard deviation (5b) is that L2 norm does not consider the variation of the data. Thus using standard deviation as the metric function is more suitable in this case, which could allow us to measure the dispersion of a set of values. However, the results of these two method are really similar to each other. In figure 5a, the y axis is the normalized hand movement distance of each frame. Figure 5b shows the standard deviation in each 15 frames. Both graphs indicate the hand movement is larger at the middle and the ending section of the video.

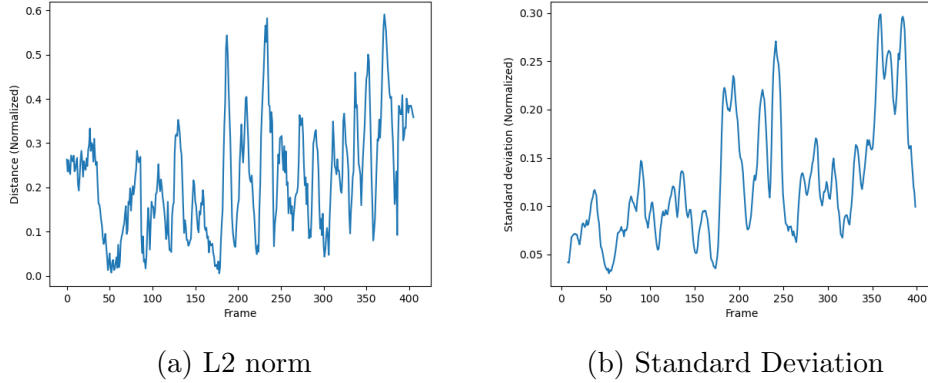


Figure 5: eban for P.S.D.C Final [5]

Table 1 is the result of figure 4a, 4b, 4c, 4d. The values in each row are nearly the same. The differences are come from the uncertainty of MediaPipe Hands, since some of the frames are hard to detect the hand position. However, the error is small enough so that we could ignore it without any problem.

	720p	1080p	1080p zoom out	1080p rotation
L2 Norm	0.3732	0.3663	0.3814	0.3606
Standard Deviation	0.4497	0.4410	0.4479	0.4246

Table 1: Dove for PSO20 R1 [3]

For table 2, we took fukrou for JapEn 15th as an example. The result after adjusting the brightness and saturation is much more reasonable than doing nothing. We also analyzed plenty of pen spinning videos in table 3.

	Original video	Adjust Brightness and Saturation
L2 Norm	1.0060	0.3848
Standard Deviation	1.2947	0.5908

Table 2: fukrou for JapEn 15th [4]

	Eban for P.S.D.C Final [5]	iroziro for JapEn 15th [6]	Wabi for TsumRab- bit [7]	i.suk for Simple 4th [8]
L2 Norm	0.2183	0.3079	0.4931	0.8280
Standard Deviation	0.2521	0.3502	0.5854	0.9786

Table 3: Testing results

6 Conclusion

We have successfully used MediaPipe Hands [1] to generate data on hand movements in pen spinning. This allows us to visualize the ups and downs of the hand movements about the pen spinning performance.

This research can help pen spinning competition’s judges to judge the control of the combo objectively. Everyone who want to know how much the hand movement is in his combo in order to improve his spinning skill can also use this method. As this research continues to develop, the computer would be able to evaluate pen spinning videos individually. We can expect a fair judging system in the future since subjectivity and bias are no longer involved in a program.

We only collected data of hand position. In order to automatically evaluate the execution, control, smoothness, and presentation of pen spinning, we need to get more information on various data such as circular orbit and speed change of the pen. We leave it as a future work to finish these analyses. However, making the program able to calculate the unstable rate of a pen spinning video is definitely a big first step for this community.

Furthermore, we can not only apply this research to pen spinning, but also to dance and finger skating competitions where the stability of the performance is important once we have the data of the coordinates. Certainly it would take a tremendous amount of time to achieve this goal, but we are confident that it is completely feasible. We are all looking forward to applying cutting edge technique to pen spinning and any other field one day.

References

- [1] Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Andrei Tkachenka, George Sung, Chuo-Ling Chang, and Matthias Grundmann. Mediapipe hands: On-device real-time hand tracking. *arXiv preprint arXiv:2006.10214*, 2020.
- [2] Gary Bradski and Adrian Kaehler. *Learning OpenCV: Computer vision with the OpenCV library*. ” O’Reilly Media, Inc.”, 2008.
- [3] psdove. [PSO20] Aestheticism R1 Dove. <https://youtu.be/y3ENwIk68j8>, 2020.
- [4] fukrou ps. JapEn15th _ fukrou. <https://youtu.be/kI5oiAP7VaU>, 2019.
- [5] eban. P.S.D.C vs malimo. <https://youtu.be/ir6yHVJ2cec>, 2014.
- [6] duriziro. japan 15th. <https://youtu.be/ovS4Dgcrba8>, 2019.
- [7] Wabi PS. for tsumrabbitt. https://youtu.be/FnjH4HEq_R8, 2020.
- [8] iSuKps. for Simple 4th. <https://youtu.be/oWZxJvJm5Jc>, 2020.