

2020年度

筑波大学情報学群情報科学類

卒業研究論文

題目

Analysis of Sports Movements Based on
Time-Series Pose Data

(時系列姿勢データに基づくスポーツ動作の解析)

主専攻 ソフトウェアサイエンス主専攻

著者 スコット アトム

指導教員 福井 和広

Abstract

A long standing goal in the field of classification is to develop effective techniques which are applicable to real life problems. With recent advancements in computer vision, significant strides towards this goal has been made due to the increased ability to capture richer and more fine-grained human motion data. For example, it is now possible to capture three dimensional human pose informatino at a frame rate of 120 frames per second with little error using motion capture technology. Such data can be used to improve classification in human action as we demonstrate in domain of injury prevention.

In parallel, exponential growth in computer processing power has led to an incredible improvement in machine learning methods. For example there have been multiple breakthroughs in object detection, voice recognition and language modeling etc. Using pose estimators based on deep learning, we can extract three dimension human pose information from video without the use of motion capture. Thus allowing a lightweight and easy-to-use application of motion analysis in enviroments where the setup cost of motion capture makes the use of it unrealistic.

Despite this progress, there are relatively few studies which effectively couple thse advancementswith human motion analysis. Therefore in this thesis, we introduce two studies that demonstrate the use of motion analysis with machine learning.

In the first study, we search for factors that differentiate the motion between athletes who are at risk of injury using motion capture and ground reaction force data of a commonly used sensorimotor performance indicator known as the single-leg drop jump.

In the second study, we propose an analysis framework for extracting important information from video and numerical data in fencing matches to assist experts in their analysis work.

We show that our approaches can provide useful information to assist those working in the domain of human motion analysis.

Contents

1 Analysis of Single-Legged Jumping Motion Based on 3D Time-Series Pose Data	1
1.1 Sensorimotor Performance Indictors	2
1.1.1 Ground Reaction Forces	3
1.1.2 Pose Characteristics	3
1.2 Subspace Based Classification	4
1.2.1 Shape Subspace	4
1.2.2 Subspace Method	5
1.2.3 Mutual Subspace Method	6
1.2.4 Grassmann Discriminant Analysis	7
Linear Discriminant Analysis	7
Grassmann Manifold and Grassmann Kernel	7
1.3 Proposed Method	9
1.4 Experiments	10
1.5 Dataset	11
1.6 Evaluation Metrics	14
1.6.1 Accuracy	14
1.6.2 F1 Score	14
1.6.3 Area Under the Curve	14
1.7 Experiment 1.	15
1.7.1 Experiment setup	15
1.7.2 Results	15
1.7.3	15
1.8 Discussion	16
2 Development of an Analysis Framework for Fencing Based on 2D/3D Time-Series Pose Data	17
2.1 Analytical Frameworks in Fencing	18
2.2 Proposed Method	20
2.2.1 Generation of Panoramic Images	21
2.2.2 Extraction of Human Regions	21
2.2.3 Inpainting Human Regions	21
2.2.4 Selecting Keyframe Images	22
2.2.5 Connecting Keyframe Images	22

2.2.6	extraction of skeletal and positional information	24
2.2.7	Detect and track the corresponding player	24
2.2.8	extraction of skeletal and location information	24
2.2.9	Visualization and analysis of skeletal and positional information	26
2.3	Experiments	27
2.3.1	Dataset	27
2.4	Results and Discussion	28
2.4.1	Generation of panoramic images	28
2.4.2	clustering, visualization	28
2.5	conclusion	30
3	Conclusion	31
Appendices		32
Acknowledgements		34
References		35

List of Figures

1.1	Single leg drop jump[1].	2
1.2	Two jump sequences sampled at 1fps. The successful jump is colored in green and the unsuccessful jump is colored in red.	13
2.1	Images of study conducted in [2]	19
2.2	Pipeline of Proposed Method	20
2.3	process until the person area is inpainted	21
2.4	SIFT feature matching with RANSAC method	22
2.5	23
2.6	Pose Estimation with Video Pose	24
2.7	The process of extracting the skeletal and location information of the corresponding player from an unattended panorama.	25
2.8	Clustering results. Up to the second principal component is illustrated.	28
2.9	An example of a cluster posture frame. The label is the same as in Figure 2.8 . . .	29

List of Tables

1.1	List of binary labels	11
1.2	List of body markers	12

Chapter 1

Analysis of Single-Legged Jumping Motion Based on 3D Time-Series Pose Data

There are a wide range of benefits can be received by progressing research in human motion analysis. One example that could significantly benefit many people is the prevention of injury in sport. For professional athletes, sustaining an injury directly leads to an economic loss. For elite young athletes a single injury has the potential to squander hopes and dreams of playing as a professional. For societies with healthcare, injuries are a significant economic cost to society, with studies showing the mean medical cost of a high school varsity athlete was \$709 per injury, \$2223 per injury in human capital costs, and \$10432 per injury in comprehensive costs [3].

Although there are many studies on the topic of injury prevention, these are usually from either a rigorously medical perspective [4] or with a focus on machine learning with empirical (big) data [5]. This is a sizeable gap in which we try to bridge.

In this study, we aim to find factors that differentiate the motion between athletes who are at risk of injury. In order to do so, we develop models and techniques that allow us to analyse granular motion data of a well studied sensori-motor control task, the single-leg drop jump (SDJ).

We demonstrate that the methods introduced can be used to detect anomalous motions. Moreover, we show that by interpreting the decision mechanism that led to the results can provide valuable practical information.

1.1 Sensorimotor Performance Indictors

"Sensorimotor" indicates the involvement of both sensory and motor functions. The ability to execute sensorimotor tasks well is crucial for performance in athletic activity. Moreover, there is an abundance of research suggesting the connection between sensorimotor performance and injury. Therefore, great effort is directed to understanding these skills. In this section we will briefly review the literature intersecting sensorimotor performance and injury. In particular we will also introduce a method that involves an action known as the single leg drop jump (SLDJ) which is frequently used to measure sensorimotor performance.

SLDJ is a unilateral horizontal drop jump, in other words the subject hops off an elevated platform of usually 30cm. It has shown to be more reliable, in means of reproducibility, than other sensorimotor performance tasks and is known to show similarities to the bilateral drop jump which is used in many areas including athlete assessment, performance monitoring, talent identification and rehabilitation [6].

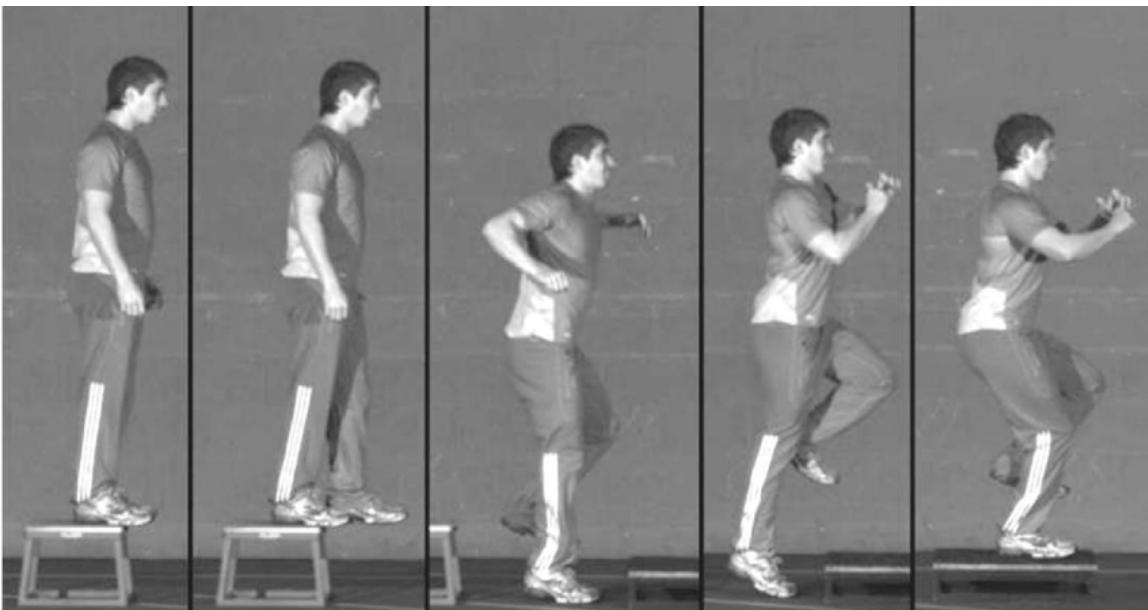


Figure 1.1: Single leg drop jump[1].

1.1.1 Ground Reaction Forces

There are largely two methods of measuring a SLDJ. The first is by measuring the ground reactions forces (GRFs) using a special measurement device known as a force plate. Several indicators such as Time to stabilization (TTS) [7], dynamic postural stability index (DPSI) [8] center of pressure (COP) [7] are then calculated and are assessed. A number of studies show that TTS, DPSI, COP can be used to differente participants between chronic ankle instability (CAI), functional ankle instability (FAI), anterior cruciate ligament etc [9]. Yet, to date, the interrelations among TTS, DPSI, and other indicators are largely unknown [8].

It is important to note that there are multiple variations to calculate the above methods due to the possibility of different sample rate, filter settings or trial length. Such variances can cause a difference on outcome values that may lead to contradictory results [10]. This makes it difficult to determine whether one can correclty infer probability of injury using GRF measurements.

1.1.2 Pose Characteristics

The second method of measurement of SLDJ is the measurement of pose characteristics using either manual notation or automated motion capture. In most cases captured data is then analysed manually, without the use of machine learning techniques. It has been shown that women have a greater valgus knee angle at time of initial contact then men performing a SLDJ [11]. However, due to the high-dimensional and multivariate nature of this data, there only a handful of research focusing on the measurement of pose characteristics. In this paper we aim to overcome this difficulty using subspace based classifiers and introduce a method that can be used to conduct classification using motion capture data.

1.2 Subspace Based Classification

Subspace analysis is a term used to describe a general framework used in computer vision, that is used for the comparison and classification of subspaces. In this section, we introduce typical approaches in subspace analysis, in particular the classification of subspaces, and describe how it can be extended to analyse shapes such as human poses.

1.2.1 Shape Subspace

Given a matrix $X \in \mathbb{R}^{n \times 3}$ where each row represents a three-dimensional coordinate of some shape, in this case a point on the human body, we can generate a shape matrix by translating the coordinate system so the each coordinate is relative to center, i.e the root, of the shape[12].

The columns that a shape matrix spans is called a shape subspace. A shape subspace is invariant under affine projection. It is also invariant to changes of coordinates caused by camera rotation and object motions. This characteristic of the shape subspace has made it effective in various tasks, such as motion segmentation and sequential factorization. Since shape subspace is invariant to affine projection, a subspace similarity metric, such as canonical angles, can effectively represent the geometrical relation between two shape subspaces.

The following methods (subspace method, mutual space method, Grassmann Discriminant Analysis etc.) are used to classify subspaces. Thus they can be applied to shape subspaces without any modification.

1.2.2 Subspace Method

The Subspace method assumes an input vector x and k -class subspaces. Each class subspace approximates a data distribution for a single class. This approximation is obtained by applying Principle Component Analysis (PCA) to each class. The similarity S of the input vector x to the i^{th} class subspace \mathcal{P} is defined based on either:

- The length of the projection of x to \mathcal{P} [13].
- The minimum angle between x and \mathcal{P} [14].

The length of an input vector x is often normalized to 1. In this case these two criteria are identical. Since they are the same, from here on we think of the angle-based similarity S defined by the following equation:

$$S = \cos^2 \theta = \sum_{i=1}^k \frac{(x^\top \cdot \Phi_i)^2}{\|x\|^2}$$

Φ_i is the i^{th} orthogonal normal basis vector of the class subspace \mathcal{P} , which are obtained from applying the PCA to a set of patterns of the class. In more rigorous terms, these orthonormal basis vectors can be obtained as the eigenvectors of the correlation matrix

$$\sum_{i=1}^l x^{(i)} x^{(i)\top} \text{ where } x^{(i)} \in \mathbb{X}$$

of the class (\mathbb{X} is the training dataset).

Learning Phase

1. Generate k class subspaces from each class by using PCA.

Recognition Phase

1. Calculate S between x and each subspace \mathcal{P}, \mathcal{Q} etc.
2. Classify the x into the class of the subspace where S was calculated to be the highest.

1.2.3 Mutual Subspace Method

The Mutual Subspace Method (MSM) is an extension of the Subspace Method (SM), where instead of having an input vector x , we use an input subspace \mathcal{P} . MSM is commonly used for image set classification [15].

The Subspace method assumes an input subspace and k class subspaces. Let us define the input subspace to be a d_p -dimensional subspace \mathcal{P} and the class subspaces to be d_q -dimensional subspaces $\{\mathcal{P}, \mathcal{Q}, \mathcal{R}, \dots\}$.

The similarity S between, for example, \mathcal{P} and \mathcal{Q} was originally defined as the minimum canonical angle θ_1 . Canonical angles [16] are uniquely defined as:

$$\cos^2 \theta_i = \max_{\substack{u_i \perp u_j (=1, \dots, i-1) \\ v_i \perp v_j (=1, \dots, i-1)}} \frac{|(u_1, v_i)|^2}{\|u_i\|^2 \|v_i\|^2}$$

Where $u_i \in \mathcal{P}, v_i \in \mathcal{Q}, \|u_i\| \neq 0, \|v_i\| \neq 0$.

We can also include the remaining canonical angles when calculating the similarity.

$$\tilde{S} = \frac{1}{t} \sum_{i=1}^t \cos \theta_i^2$$

This value \tilde{S} reflects the structural similarity between two subspaces. It is also defined on the t smallest canonical angles. For practical applications, the canonical angles between subspaces \mathcal{P}_1 and \mathcal{Q} are obtained by calculating the singular values $\{\lambda_j\}_{j=1}^m$ of the correlation matrix between their basis matrices $P, Q \in \mathbb{R}^{d \times m}$, i.e. solving the SVD of $VV^\top = P^\top Q$. This corresponds to finding the rotation of each basis that is closest to the opposing subspace. From solving this problem, the canonical angles can be obtained by $\theta_j = \cos^{-1}(\lambda_j)$, where $j = 1, \dots, m$.

Learning Phase

1. Generate k class subspaces from each class by using PCA.

Recognition Phase

1. Calculate \tilde{S} between the input subspace and each dictionary subspace.
2. Classify the input subspace into the class where the dictionary subspace was calculated to be the highest.

1.2.4 Grassmann Discriminant Analysis

Next, we introduce here the concept Grassman Discriminant Analysis(GDA)[17] which is a discriminatory mechanism to separate classes of signals. First we will, briefly explain its predecessor, Linear Discriminant Analysis (LDA) [18]. LDA is a dimensionality reduction technique that is well known and has been successfully used for classification. In essence, GDA is conducted as kernel LDA but with the Grassmann kernels.

Linear Discriminant Analysis

The goal of Linear Discriminant Analysis (LDA) is to project a dataset onto a lower-dimensional space that has good class-separability. Due to the curse of dimensionality, dimensionality reduction is an import step that is necessary in order avoid overfitting. It also has the added benefit of reducing computational costs.

Learning Phase

1. Compute d -dimensional mean vectors for the different classes from the dataset.
2. Compute the inbetween-class and within-class scatter matrices.
3. Computer the eigenvectors and corresponding eigenvalues for the scatter matrices.
4. Sort the eigenvectors by decreasing eigenvalues and choose k eigenvectors with the largest eigenvalues to form a $n \times k$ dimensional matrix W

Recognition Phase

1. Project the samples onto the new subsapce, $YY=XX^TWW$ (where XX is a $n \times d$ -dimensional matrix representing the n samples, and yy are the transformed $n \times k$ -dimensional samples in the new subspace).
2. Conduct classification using a classifier of choice (a simple K-Neighbors classifier is commonly used).

Grassmann Manifold and Grassmann Kernel

The Grassmann manifold $\mathcal{G}(m, d)$ is defined as the set of m -dimensional linear subspaces of \mathbb{R}^d . A Grassmann manifold can be embedded in a reproducing kernel Hilbert space by the use of a Grassmann kernel. In this case, the most popular kernel is the projection kernel k_p , which can be defined as $k_p(\mathcal{Y}_1, \mathcal{Y}_2) = \sum_{j=1}^m \cos^2 \theta_j$, which is homologous to the subspace similarity. We can measure the distance between two points on a Grassmann manifold by using this projection kernel [?], and a subspace \mathcal{Y} can be represented by a vector with regards to a reference subspace dictionary $\{\mathcal{Y}_q\}_{q=1}^N$ as $y = k_p(\mathcal{Y}, \mathcal{Y}_q) = [k_p(\mathcal{Y}, \mathcal{Y}_1), k_p(\mathcal{Y}, \mathcal{Y}_2), \dots, k_p(\mathcal{Y}, \mathcal{Y}_N)] \in \mathbb{R}^N$.

$$\begin{aligned}
Ra(\alpha) &= \frac{\alpha^\top S_b \alpha}{\alpha^\top S_w \alpha} = \\
&= \frac{\alpha^\top K(V - e_N e_N^\top / N) K \alpha}{\alpha^\top (K(I_N - V)K + \sigma^2 I_N) \alpha} = \\
&= \frac{\alpha^\top b \alpha}{\alpha^\top (w + \sigma^2 I_N) \alpha}, \tag{1.1}
\end{aligned}$$

where K is the kernel matrix, e_N is a vector of ones that has length N , V is a block-diagonal matrix whose c -th block is the matrix $e_{N_c} e_{N_c}^\top / N_c$, and $b = K(V - e_N e_N^\top / N)K$. For example, the kernel matrix, K , is calculated as the similarity matrix between subspaces Y_q and Y_w . The term $\sigma^2 I_N$ is used for regularizing the covariance matrix $w = K(I_N - V)K$. It is composed of the covariance shrinkage factor $\sigma^2 > 0$, and the identity matrix I_N of size N . The set of optimal vectors α are computed from the eigenvectors of $(w + \sigma^2 I_N)^{-1} b$.

We apply the GDA algorithm to the reference subspaces Y_i^c to generate reference vectors y_i^c . When given an unknown bioacoustic signal $x_{in}(t)$, we compute its SSA subspace \mathcal{Y}_{in} and map it onto the manifold to generate a vector y_{in} ; then we predict its corresponding bioacoustic class (e.g. species) based on the nearest reference vector (1-NN).

1.3 Proposed Method

1.4 Experiments

In this chapter we conduct two experiments to demonstrate that our proposed method outperforms other traditional methods that have been introduced in previous chapters. The implementation is available at the website xxxx.

1.5 Dataset

We collected data of 524 instances of a single-leg drop jump landing from 144 college football (soccer) players from the University of Tsukuba.

The dataset consists of the following attributes:

- 3 dimensional coordinates of 29 body markers capture by multiple infrared cameras. A full list of the where the body markers were placed is available in Table 1.2.
- An integer (1-9) indicating the classification of ankle instability, labelled by an expert (label description: 1=Healthy, 2=Structural instability, 3=Subjective instability, 4=Sprained more than 3 times, 5=2 and 3, 6=2 and 4, 7=3 and 4, 8=2 and 3 and 4, 9=healthy but with a history of one or two sprains).
- Binary labels (0 or 1) indicating classes. For example, whether the single-leg drop jump was successful (0) or not (1). After consultation with an expert, we defined a successful jump as a jump where the subject was able to keep their balance on a single leg for more than 5 seconds after landing. A full list of the binary labels are available in Table 1.1.
- Force plate data.
- Weight (kg) for each of the individual 144 athletes.

Label	Explanation
Success	0 if the subject can balance on one foot for 5 seconds after SDJ, otherwise 1
Jump Leg	0 if subject jumped with right foot, otherwise 1
Healthy	0 if all of the below are 0, otherwise 1
Prior Injury	0 if subject has no previous injuries, otherwise 1
Structural Instability	0 if expert determines that the joint is loose due to ligament damage after palpation, otherwise 1
Subjective Instability	0 if the subject is determined to have instability in the ankle joint by a questionnaire, otherwise 1
Prone	1 if the subject has 1 or more sprains and none of the above three labels are 1, otherwise 0

Table 1.1: List of binary labels

Index	Name
0	Cervical vertebra 7
1	Calcaneal tuberosity
2	Fibula Head
3	Greater trochanter
4	Suprasternal notch
5	Anterior superior iliac spine (L)
6	lateral epicondyle
7	medial malleolus
8	Posterior superior iliac spine (L)
9	medial condyle
10	medial epicondyle
11	lateral malleolus
12	Base of first metatarsal bone
13	Base of second metatarsal bone
14	Base of fifth metatarsal bone
15	Head of first metatarsal bone
16	Head of second metatarsal bone
17	Head of fifth metatarsal bone
18	Fibular trochlea of calcaneus
19	First distal phalanges
20	Anterior superior iliac spine (R)
21	Posterior superior iliac spine (L)
22	Acromion (L)
23	Acromion (R)
24	Sustentaculum tali
25	Thoracic vertebrae 8
26	Tibial tuberosity
27	Scaphoid bone
28	Xiphoid process

Table 1.2: List of body markers

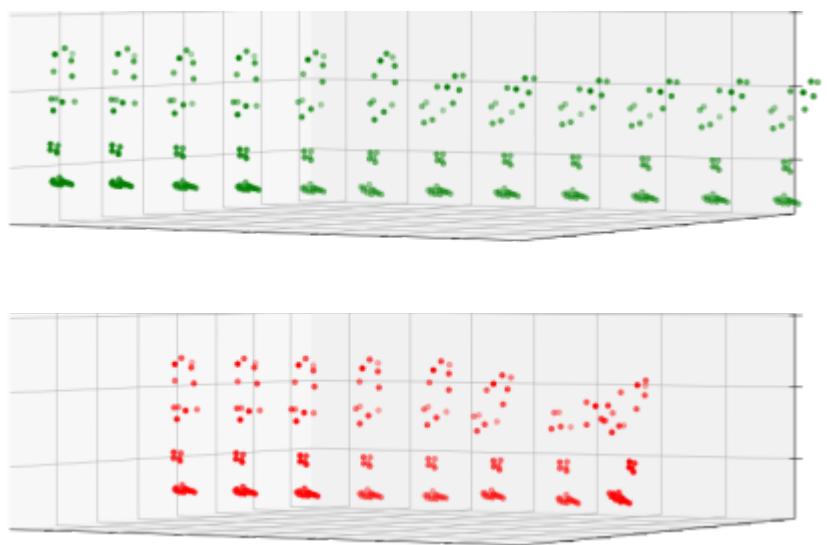


Figure 1.2: Two jump sequences sampled at 1fps. The successful jump is colored in green and the unsuccessful jump is colored in red.

1.6 Evaluation Metrics

To evaluate the performance of different classification methods, we use a 5-fold cross-validation strategy to compute the classification accuracy, F1 score and the area under the receiver operating characteristic curve (AUC).

Although accuracy is a sufficient metric for performance in datasets with symmetric classes distribution and class value, it is an inadequate metric for datasets with class imbalance and datasets where values of false positive and false negatives are differ.

Therefore, we included the two additional metrics explained above. These metrics are frequently used in medical cases such as cancer detection from images [19][20].

1.6.1 Accuracy

Accuracy is computed as the fraction of correct predictions.

$$\text{accuracy}(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} 1(\hat{y}_i = y_i)$$

1.6.2 F1 Score

The F1 score can be interpreted as an average of the classifier's ability to:

- Not label as negative samples as a positive (i.e. Precision).
- Find all positive samples (i.e. Recall).

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

1.6.3 Area Under the Curve

The receiver operating characteristic (ROC) plots the true positive rate (TPR) to the false positive rate (FPR) of a binary classifier with varying discrimination thresholds.

The area under the ROC curve (AUC) summarizes the information of the ROC in one number, between 0 and 1. AUC = 0.5 indicates an uninformative classifier (a classifier with uniform prediction for any sample) or random classifier if the classes are symmetric. Therefore no realistic classifiers should have an AUC < 0.5. AUC = 1 indicates a perfect classifier.

The algorithm to compute AUC with an in-depth explanation can be found in [21].

1.7 Experiment 1.

In this experiment we evaluate x classifiers accuracy on multiple binary classification problems regarding our single-leg drop jump dataset.

1.7.1 Experiment setup

We experiment with the following setup.

- (A) Pose data up until landing.
- (B) All pose data.
- (C) Pose data + Force Plate data.

1.7.2 Results

1.7.3

1.8 Discussion

Chapter 2

Development of an Analysis Framework for Fencing Based on 2D/3D Time-Series Pose Data

In many competitive sports, the analysis of video footage and numerical data can help gain an advantage in athletic performance. However, analysis requires not only specialized knowledge but also a great deal of time and effort. Therefore, it is very important to streamline and automate the process.

In recent years, advances in computer vision and machine learning have reduced the burden of these tasks in various fields. In this study, using such technology, we propose an analysis framework for extracting importance information from video and numerical data to assist experts in their analysis work.

Specifically, we will

1. Create a panoramic image of the entire court from fencing footage captured by a Handycam
2. Extract skeletal information of each athlete and their position on the court
3. Examine methods for analyzing and visualizing matches and athletes using the skeletal and positional information.

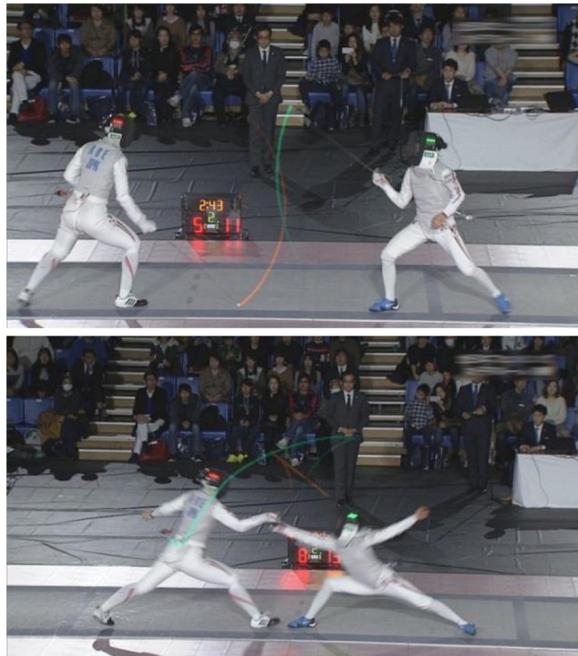
2.1 Analytical Frameworks in Fencing

There are many studies that focus on applying computer vision to fencing. For example, Takahashi[2] et al. proposed a robust detection and tracking method to visualize the trajectory of a fencing sword for live TV programs. Since live TV programs require not only accuracy but also real-time computation, they successfully track fast-moving fencing swords by using supervised machine learning to detect the swords and using a particle filter to predict their positions in the next frame.

Athow[22] et al. proposed a computer vision-based scoring system to assist referees in foil events in fencing. They used color blob detection to track fencers who wore color patches on their hands. In addition, there have been cases of computer vision and motion analysis: Makawaski et al. [23] created a new local trace image to represent fencing motions, and showed that the dynamics of the motion is useful for analyzing similar motion patterns. They showed that motion dynamics can be useful for analyzing similar motion patterns. In the same study, he also created a dataset for analyzing footwork, contributing to the development of the field.

However, there is still no research that proposes a comprehensive framework for extracting posture and positional information from video in fencing and performing tactical analysis to gain an advantage in competition.

In this study, we propose an analysis framework that uses computer vision and machine learning to extract posture and position information to assist experts in their analysis tasks. is introduced



(a) Broadcast images of sword trajectories at All Japan Fencing Championship 2017



(b) Photos of the RGB/IR camera



(c) Photos of equipment in operations area

Figure 2.1: Images of study conducted in [2]

2.2 Proposed Method

The framework proposed in this study can be divided into three main stages. In the first stage, a panoramic image is generated from a video footage. In the second stage, the player's posture coordinates and absolute coordinates on the court are estimated from the panoramic image. In the last stage, we analyze and visualize the obtained skeletal and positional information of the players using dimensionality reduction and clustering techniques. The details of each element are described below.

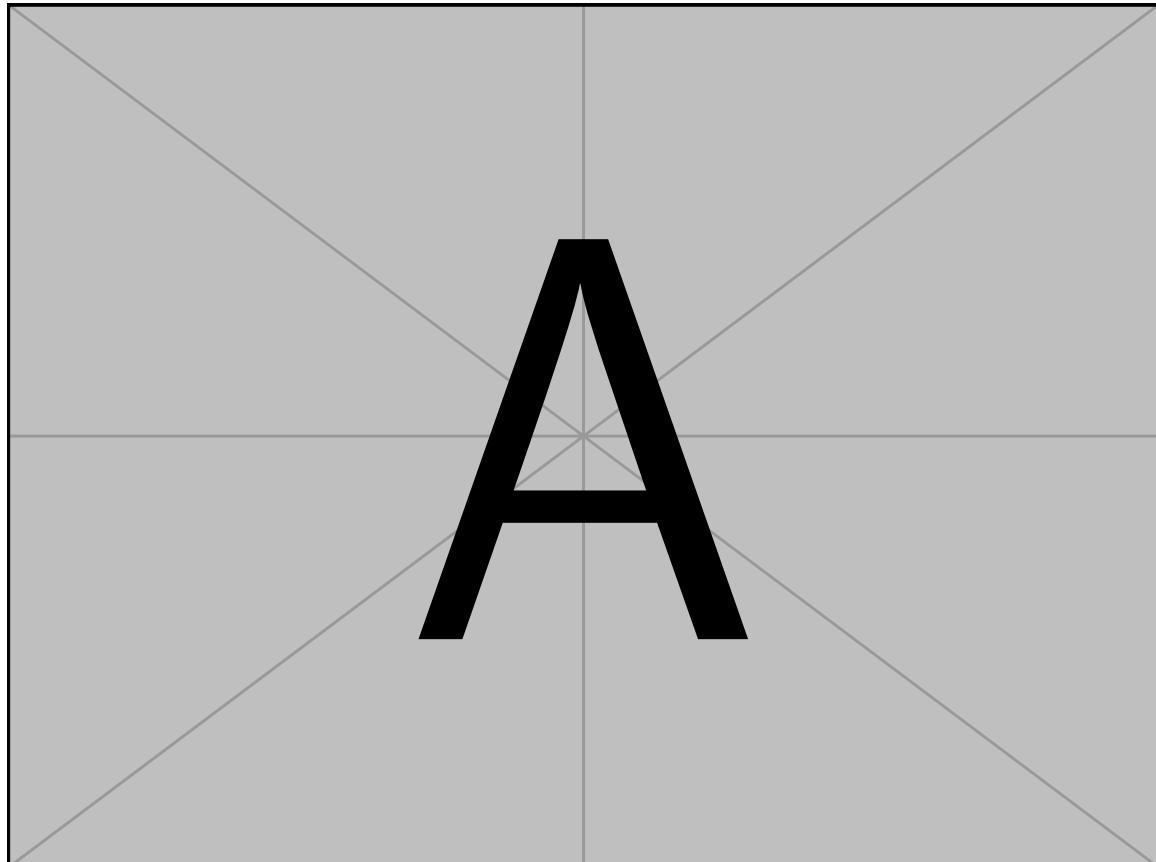


Figure 2.2: Pipeline of Proposed Method

2.2.1 Generation of Panoramic Images

Since it is strategically useful to know where the players are on the court, it is useful to generate panoramic images showing the entire court. In addition, the position of the players on the pitch can be estimated from the panoramic image by using person detection, so that the positional information of the players can be obtained not only visually but also numerically.

2.2.2 Extraction of Human Regions

Usually, panoramic background images are generated using images of only stationary objects, so to use conventional methods, moving objects, i.e., humans, must be eliminated in each frame, and the eliminated areas must be complemented in some way. In this study, this is done in two steps. For the extraction of human regions, we use segmentation with HR-NET[24], and label each pixel of each frame as human region (=1) or not human region (=0). As a post-processing step, we perform smoothing operations in the vertical, horizontal, and temporal directions to mask the human region more reliably in each frame. 2.3(b) shows the masking of the person area after post-processing.

2.2.3 Inpainting Human Regions

The masked human region is inpainted by calculating the optical flow in the front-back direction based on the Flow-edge Guided Video Completion[?] method. The result is an unattended video where the person is seamlessly removed from the video.

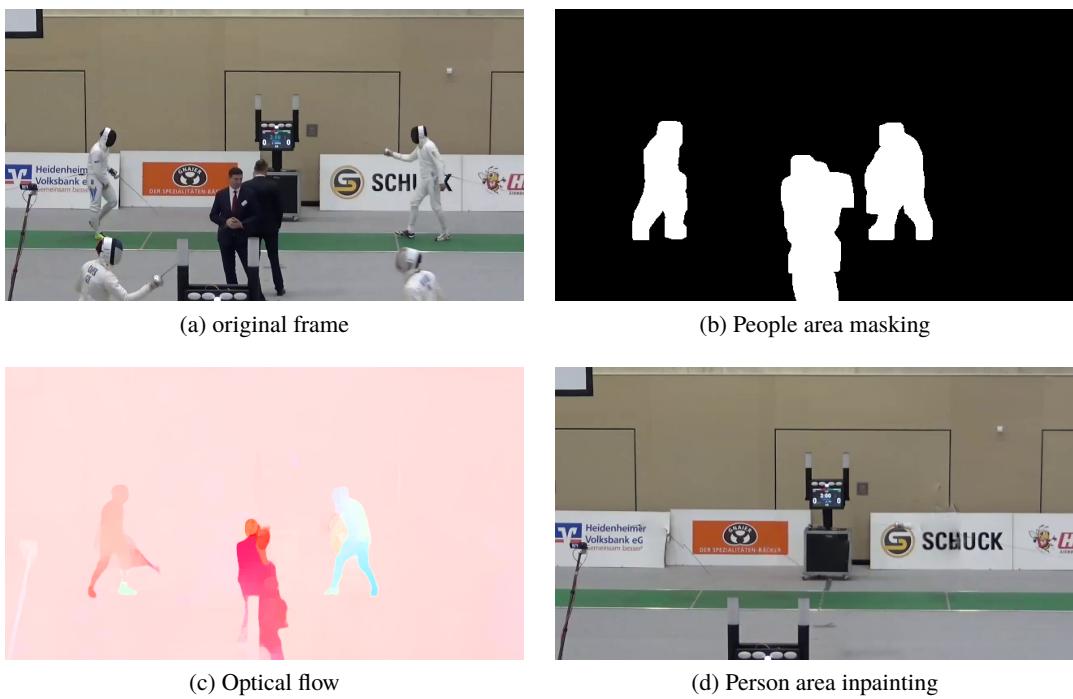


Figure 2.3: process until the person area is inpainted



Figure 2.4: SIFT feature matching with RANSAC method

2.2.4 Selecting Keyframe Images

The unattended image generated above does not contain any human, so there are almost no moving objects, which is suitable for panoramic images. However, the time complexity of the conventional Image Stitching method [25] is $O(4N^3)$ when N is the number of frames. It would take an enormous amount of time. Therefore, we selected keyframes to create a high-quality panoramic video that captures the entire court even with a small number of frames.

First, the estimated projection matrices of the first frame of the video and the other frames are flattened and used as feature vectors. To estimate the projection matrix, we compute the SIFT feature [26] in the image and use keypoint matching with the RANSAC method [27]. Figure 2.4 shows the result of keypoint matching between two frames as an example.

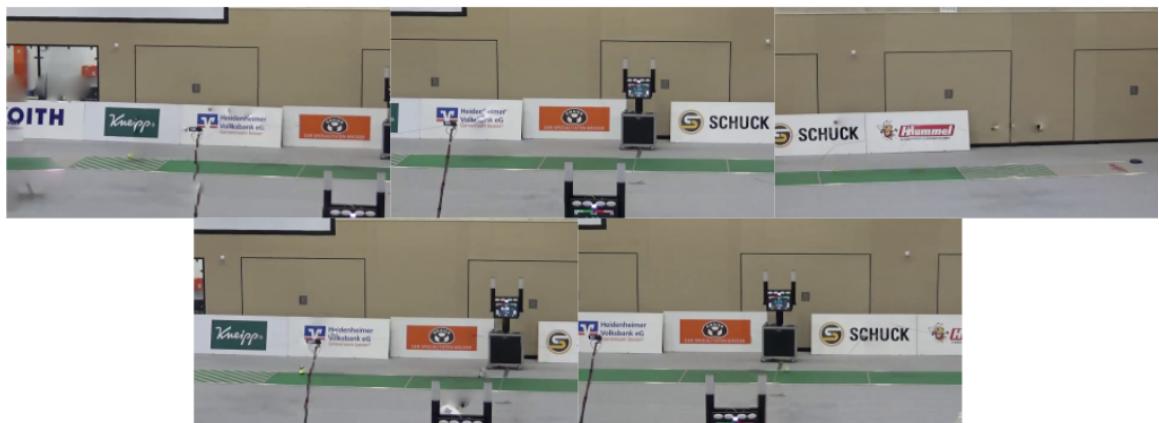
The feature vectors are reduced to two dimensions by Principal Component Analysis (PCA)[28] and T-Sributed Stochastic Neighbor Embeddings (T-SNE)

2.2.5 Connecting Keyframe Images

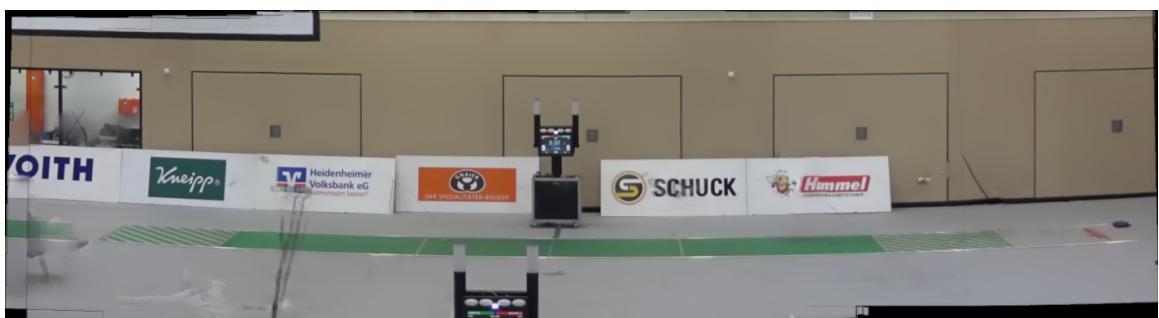
Generate a single unattended background panorama by stitching together keyframe images using the images selected in ???. For stitching, we use the Image Stitcher Class in OpenCV[29], which has a full pipeline of image registration (feature matching, wave correction, etc.) and image composition (warping, blending, etc.).

To remove noise from unattended background panoramas, the process of ?? is performed several times to generate several slightly different unattended background panoramas and apply a pixel-wise median filter. Figure 2.7(a) shows the result.

To create panoramic images from the unattended background panoramas, we compute the projection matrix by matching the SIFT features with the RANSAC method on the unattended background panoramas and each frame. The player can be overlaid on the unattended panorama by applying a matrix transformation using the projection matrix calculated above to the bounding box of the player detected by the method described in ???. Figure 2.7(c) shows how this process is applied.



(a) Example of selected keyframe images



(b) Unattended panoramic image

Figure 2.5

2.2.6 extraction of skeletal and positional information

By applying the existing object detection and posture estimation techniques to the panoramic video, we can easily obtain the skeletal and positional information of the player. The details of the method are described below.

2.2.7 Detect and track the corresponding player

To detect and track the players, we use FairMOT[30]. FairMOT solves the problems of traditional multi-object tracking (MOT) methods by using a single-shot deep neural network without anchors. FairMOT solves the problems of conventional multi-object tracking (MOT) methods by using a single-shot deep neural network without anchors, and is SoTA in terms of both tracking accuracy and inference speed.

In addition to the player in question, there are referees and players from other games in the match video, so it is necessary to sort them out. To do so, we first manually annotate the unattended panoramic video with the areas of the court where the players in the game should be. If the detected person is out of the area, the person is removed as noise. If the person is not out of the area, the person is used as the corresponding player in the overlay created through the process of 2.2.5. In Figure 2.7(b), the players that are out of the region are shown as "DELETE" and the players that are not out of the region are shown as "KEEP". When only one player was detected for the overlay, linear completion was used. When three or more players were detected, players were selected using bipartite matching [31], a Hungarian algorithm that treats skeletal information as a graph.

2.2.8 extraction of skeletal and location information

PifPaf[32] is used to obtain skeletal information for the bounding boxes of the relevant players detected by ???. The skeletal information obtained using this pose estimator is shown in Figure 2.7(d).

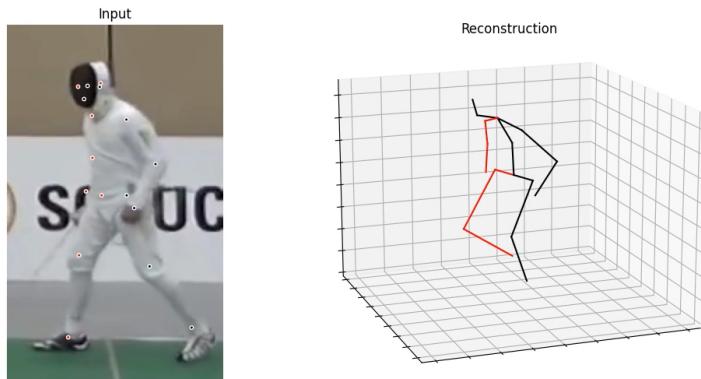
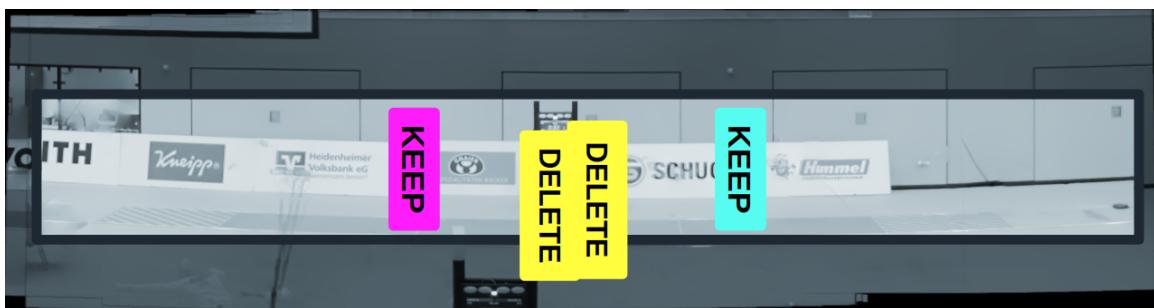


Figure 2.6: Pose Estimation with Video Pose



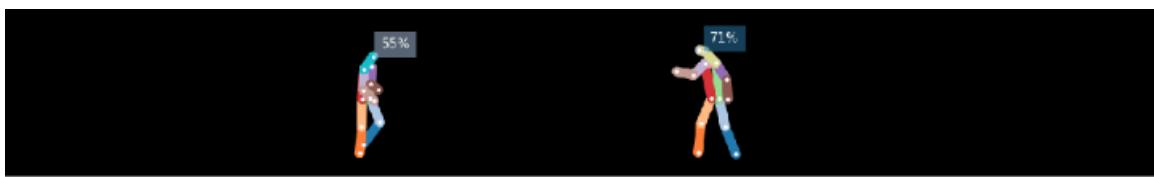
(a) unattended panorama image



(b) Area of the court in the unattended panoramic image



(c) Overlay of the corresponding player in the unattended panorama



(d) Visualization of skeletal information of the corresponding player (posture estimation using PifPaf)

Figure 2.7: The process of extracting the skeletal and location information of the corresponding player from an unattended panorama.

2.2.9 Visualization and analysis of skeletal and positional information

. From the posture and positional information, clustering of action sequences and correlation with scoring situations can be analyzed to understand frequent plays in fencing and how players changed their playing style depending on the situation.

As an example application of the proposed framework, we perform clustering of skeletal information. The skeletal information is dimensionally compressed using PCA[28] and clustering is performed using the mixed Gaussian model[33]. The mixture Gaussian model can represent a Gaussian distribution with different parameters for each cluster. For the number of clusters, we use the minimum value of the Bayesian information criterion (BIC) [34]. The BIC is a measure of the negative log-likelihood with a constraint that the number of parameters should not be too large.

2.3 Experiments

2.3.1 Dataset

The dataset was created from 100 match videos of fencing epee events provided by the Japan Fencing Association. The match videos were shot with a single Handycam and lasted about 15 minutes at 60 frames per second.

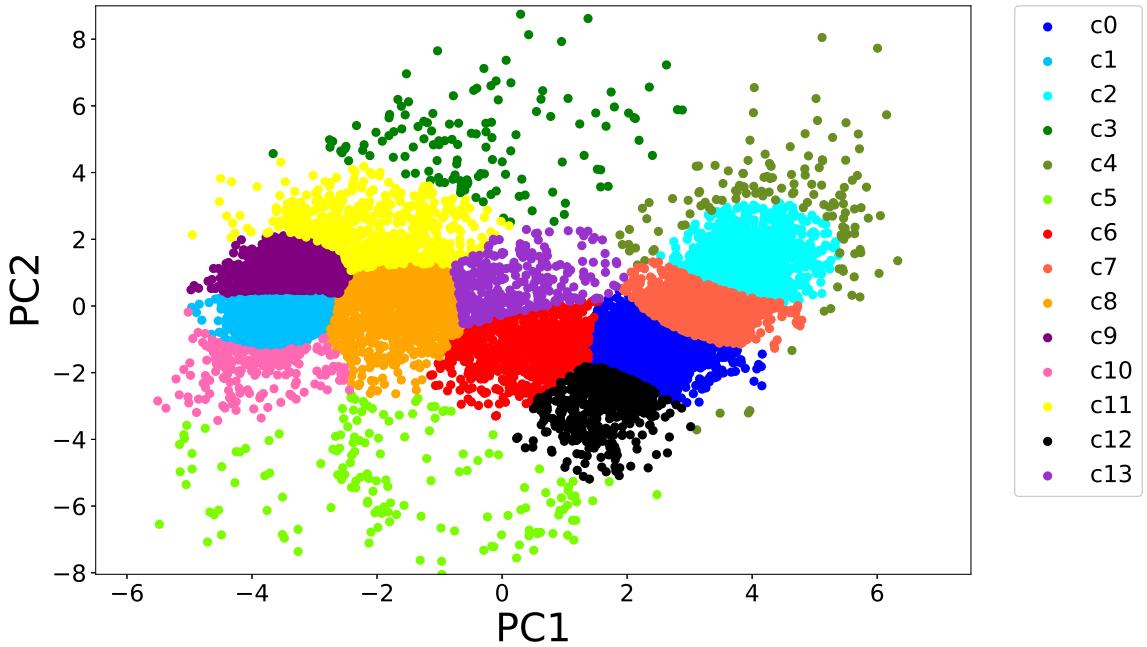


Figure 2.8: Clustering results. Up to the second principal component is illustrated.

2.4 Results and Discussion

2.4.1 Generation of panoramic images

We were able to generate panoramic images and create video mosaics more quickly than conventional methods by devising our own method for obtaining posture and position information. For analysis using posture and position information, we will conduct more detailed analysis using the position and skeletal data of the players obtained by the above method.

2.4.2 clustering, visualization

As an example application of the skeletal information obtained by the proposed framework, clustering was performed. The results of clustering up to the second principal component of the skeletal information with a mixed Gaussian model are shown in Figure 2.8. The legend indicates each cluster. Some of the representative posture frames for each cluster are shown in Figure 2.9. A characteristic posture was observed in each cluster. The most common posture observed in cluster c2 was the most neutral posture, which can be shifted to any posture. Cluster c4, on the other hand, had a lower posture than the posture frames of the other clusters, with the epee sword sticking out in front of the body. The posture of cluster c4 was lower than that of the other clusters, and it was considered to be an attacking posture, based on the characteristics of fencing. Thus, clustering using the proposed framework may make it easier to classify the plays. In addition, the posture frames of clusters c8 and c13 are not directly related to the play. Clustering can also remove these data, which would otherwise have to be done manually.

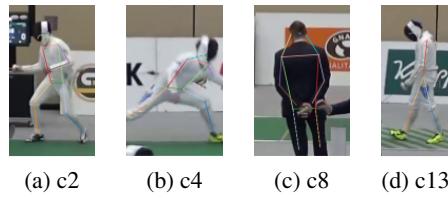


Figure 2.9: An example of a cluster posture frame. The label is the same as in Figure 2.8

2.5 conclusion

In this study, we proposed an analytical framework for extracting video and numeric values to assist experts in their analysis tasks using computer vision and machine learning in the epee fencing event. As a future prospect, we will conduct a more rigorous evaluation of the proposed method and try to extend the framework to other events.

Chapter 3

Conclusion

Appendices

Therefore we find the frame where the average foot position on the y-axis is lowest, with the condition that the frame must be in within 39 frames after the highest point.

Acknowledgements

There are many people I must thank for contributing to my extended bachelor research experience. First I must thank my adviser Kazuhiro Fukui who, through many discussions has helped me finish this thesis. I would like to thank my close collaborators and lab mates from the computer vision lab who I've had the distinct pleasure of working with and learning from.

I am thankful to the institutions and people who have made this research possible by permitting the use of necessary data. The single-leg jump data was provided by Shun Kunugi. The fencing data was provided by the Research Center for Medical and Health Data Science, The Institute of Statistical Mathematics, Research Organization of Information and Systems, and the Japan Fencing Association.

References

- [1] James Wild, Neil Bezodis, Richard Blagrove, and Ian Bezodis. A biomechanical comparison of accelerative and maximum velocity sprinting: Specific strength training considerations. *Professional Strength and Conditioning*, 01 2011.
- [2] Masaki Takahashi, Shinsuke Yokozawa, Hideki Mitsumine, Tetsuya Itsuki, Masato Naoe, and Satoshi Funaki. Real-time visualization of sword trajectories in fencing matches. *Multimedia Tools and Applications*, 79(35):26411–26425, 2020.
- [3] Sarah B Knowles, Stephen W Marshall, Tyler Miller, R Spicer, J Michael Bowling, D Loomis, RW Millikan, Jingzhen Yang, and FO Mueller. Cost of injuries from a prospective cohort study of north carolina high school athletes. *Injury prevention*, 13(6):416–421, 2007.
- [4] Gregory D Myer, Kevin R Ford, and Timothy E Hewett. Rationale and clinical techniques for anterior cruciate ligament injury prevention among female athletes. *Journal of athletic training*, 39(4):352, 2004.
- [5] J Alderson. A markerless motion capture technique for sport performance analysis and injury prevention: Toward a ‘big data’ , machine learning future. *Journal of Science and Medicine in Sport*, 19:e79, 2015.
- [6] Markus Stålbom, David Jonsson Holm, John B. Cronin, and Justin W.L. Keogh. Reliability of kinematics and kinetics associated with Horizontal Single leg drop jump assessment. A brief report. *Journal of Sports Science and Medicine*, 6(2):261–264, 6 2007.
- [7] Duncan P. Fransz, Arnold Huurnink, Idsart Kingma, and Jaap H. van Dieën. How does postural stability following a single leg drop jump landing task relate to postural stability during a single leg stance balance task? *Journal of Biomechanics*, 47(12):3248–3253, 9 2014.
- [8] Arnold Huurnink, Duncan P. Fransz, Idsart Kingma, Vosse A. de Boode, and Jaap H.van Dieën. The assessment of single-leg drop jump landing performance by means of ground reaction forces: A methodological study. *Gait and Posture*, 73:80–85, 9 2019.
- [9] ERIK A. WIKSTROM, MARK D. TILLMAN, and PAUL A. BORSA. Detection of Dynamic Stability Deficits in Subjects with Functional Ankle Instability. *Medicine & Science in Sports & Exercise*, 37(2):169–175, 2 2005.

- [10] Duncan P. Fransz, Arnold Huurnink, Vosse A. De Boode, Idsart Kingma, and Jaap H. Van Dieën. Time to stabilization in single leg drop jump landings: An examination of calculation methods and assessment of differences in sample rate, filter settings and trial length on outcome values. *Gait and Posture*, 41(1):63–69, 2015.
- [11] Kyla A. Russell, Riann M. Palmieri, Steven M. Zinder, and Christopher D. Ingersoll. Sex differences in valgus knee angle during a single-leg drop jump. *Journal of Athletic Training*, 41(2):166–171, 4 2006.
- [12] Yosuke Igarashi and Kazuhiro Fukui. 3D object recognition based on canonical angles between shape subspaces. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 6495 LNCS(PART 4):580–591, 2011.
- [13] WATANABE and S. Evaluation and Selection of Variables in Pattern Recognition. *Computer and Information Science II*, pages 91–122, 1967.
- [14] Taizo Iijima, Hiroshi Genchi, and Kenichi Mori. THEORY OF CHARACTER RECOGNITION BY PATTERN MATCHING METHOD. pages 50–56, 1973.
- [15] Akinari Sakai, Naoya Sogi, and Kazuhiro Fukui. Gait Recognition Based on Constrained Mutual Subspace Method with CNN Features. In *2019 16th International Conference on Machine Vision Applications (MVA)*, pages 1–6. IEEE, 2019.
- [16] Françoise Chatelin. *Eigenvalues of Matrices*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2012.
- [17] Jihun Hamm and Daniel D Lee. Grassmann discriminant analysis: A unifying view on subspace-based learning. In *Proceedings of the 25th International Conference on Machine Learning*, pages 376–383, 2008.
- [18] Keinosuke Fukunaga. *Introduction to Statistical Pattern Recognition (2nd Ed.)*. Academic Press Professional, Inc., USA, 1990.
- [19] Korsuk Sirinukunwattana, Shan E.Ahmed Raza, Yee Wah Tsang, David R.J. Snead, Ian A. Cree, and Nasir M. Rajpoot. Locality Sensitive Deep Learning for Detection and Classification of Nuclei in Routine Colon Cancer Histology Images. *IEEE Transactions on Medical Imaging*, 35(5):1196–1206, 5 2016.
- [20] Dan C. Cireşan, Alessandro Giusti, Luca M. Gambardella, and Jürgen Schmidhuber. Mitosis detection in breast cancer histology images with deep neural networks. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 8150 LNCS, pages 411–418, 2013.
- [21] Tom Fawcett. An introduction to ROC analysis. *Pattern Recognition Letters*, 2006.
- [22] Stephanie Athow and Jeff McGough. Using computer vision to assist the scoring of modern fencing.

- [23] Filip Malawski and Bogdan Kwolek. Recognition of action dynamics in fencing using multi-modal cues. *Image and Vision Computing*, 75:1 – 10, 2018.
- [24] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019.
- [25] Matthew Brown and David G Lowe. Automatic panoramic image stitching using invariant features. *International journal of computer vision*, 74(1):59–73, 2007.
- [26] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.
- [27] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
- [28] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987.
- [29] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.
- [30] Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. Fairmot: On the fairness of detection and re-identification in multiple object tracking. *arXiv preprint arXiv:2004.01888*, 2020.
- [31] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- [32] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Pifpaf: Composite fields for human pose estimation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2019.
- [33] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [34] Richard H Jones. Bayesian information criterion for longitudinal and clustered data. *Statistics in medicine*, 30(25):3050–3056, 2011.