



# Design and Implementation of a Body Fat Classification Model using Human Body Size Data

Taejun Lee<sup>1</sup>, Hakseong Kim<sup>1</sup>, and Hoekyung Jung<sup>1\*</sup>, *Member, KIICE*

<sup>1</sup>Department of Computer Engineering, PaiChai University, Daejeon 155-40, Republic of Korea

## Abstract

Recently, as various examples of machine learning have been applied in the healthcare field, deep learning technology has been applied to various tasks, such as electrocardiogram examination and body composition analysis using wearable devices such as smart watches. To utilize deep learning, securing data is the most important procedure, where human intervention, such as data classification, is required. In this study, we propose a model that uses a clustering algorithm, namely, the K-means clustering, to label body fat according to gender and age considering body size aspects, such as chest circumference and waist circumference, and classifies body fat into five groups from high risk to low risk using a convolutional neural network (CNN). As a result of model validation, accuracy, precision, and recall results of more than 95% were obtained. Thus, rational decision making can be made in the field of healthcare or obesity analysis using the proposed method.

**Index Terms:** Classification, CNN, fat rate percentage, human body size, K-means clustering

## I. INTRODUCTION

Deep learning is being applied in various fields, including the healthcare field. Particularly, as the personal smartphone and wearable device markets grow, body composition analysis, electrocardiogram testing, and diet management are provided in an easily accessible manner to improve eating habits.

Securing data is important when applying deep learning. In the past, only medical records, medical insurance information, and data measured by users were used; however, information collected through various channels, such as the public, social media services, and genetic data, have recently been used [1].

It can be said that the pre-processing process of the collected information is the most important step. The pre-processing requires human intervention, such as removing missing values or outliers and clustering, which is in the domain of unsupervised learning [2-6].

In this study, using only easy-to-measure aspects, such as waist circumference, reported in the Size Korea public dataset provided by the National Institute of Technology, the body fat percentage data were clustered according to gender and age to create a target setpoint. Accordingly, we propose a model that undergoes a pre-processing step generate the desired setpoint.

It is believed that a model using both supervised and unsupervised learning can be used in obesity analysis or body shape management service cases.

## II. RELATED WORKS

In this section, we examine human body size data provided by the National Institute of Technology and Standards and related research.

Received 30 October 2021, Revised 06 December 2021, Accepted 09 December 2021

\*Corresponding Author Hoekyung Jung (E-mail: [hkjung@pcu.ac.kr](mailto:hkjung@pcu.ac.kr), Tel:+82-42-520-5640)

Department of Computer Engineering, PaiChai University, Daejeon, 35345 Republic of Korea

Open Access <https://doi.org/10.56977/jicce.2023.21.2.110>

print ISSN: 2234-8255 online ISSN: 2234-8883

© This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

Copyright © The Korea Institute of Information and Communication Engineering

### A. Human Body Size Data

Since 1978, the National Institute of Technology and Standards governed by the Ministry of Trade, Industry, and Energy has been promoting the human body dimensional survey dissemination project to measure and standardize the human body information of Koreans and produce products that are convenient for Koreans to use with ergonomic designs [7].

The latest data are provided in the 7th human body size dataset released in December 2015. The total number of people measured was 6,413, which were categorized based on age and gender as reported in Table 1.

**Table 1.** People groups included in the 7th Korean human body size survey project

	Man	Woman	Total
16-19	998	926	1,924
20-29	869	668	1,573
30-39	655	675	1,330
40-49	311	359	670
50-59	221	358	579
60-69	145	228	373
<b>Total</b>	<b>3,199</b>	<b>3,214</b>	<b>6,413</b>

### B. Related Research Conducted using the Human Body Size Data

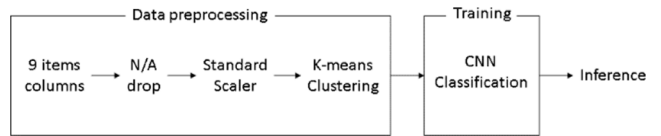
There are various studies that are in line with the human body size survey project. Specifically, there is a study that divides the lower body types of men in their 30s to design appropriate pant patterns [8]. In addition, another study divided the upper body type of male college students in their 20s into four clusters by extracting factors through principal component and K-means cluster analyses [9]. The first study was significant in that it suggested changes in clothing design or dimensions by statistically analyzing changes in human body dimensions.

## III. DATA PREPROCESSING AND MODEL DESIGN

In this section, we discuss the overall structure of the proposed model, pre-processing of the training data, and design.

### A. Overall Structure

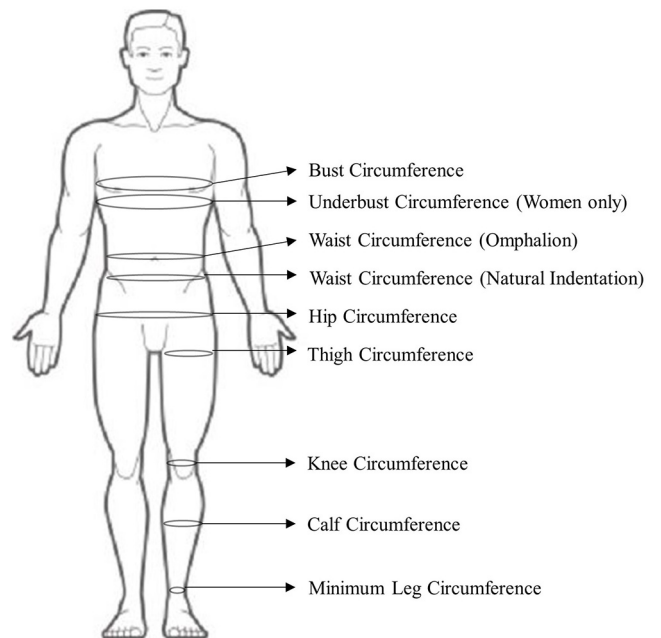
Fig. 1 shows the data pre-processing and model training processes for the human body size data. The data pre-processing procedure is discussed in Section III-B, and the convolutional neural network (CNN) classification process is discussed in Section III-C.



**Fig. 1.** Overall structure.

### B. Data Item Selection and Pre-processing

Several items of human body size data were originally considered in the model. However, only nine of them were used, namely, the bust circumference, under bust circumference, waist circumference (natural indentation), waist circumference (omphalion), hip circumference, thigh circumference, knee circumference, calf circumference, and minimum leg circumference. The names set by the National Institute of Technology and Standards were used as is, and there were no data regarding the under bust circumference for men. Fig. 2 depicts the nine human body aspects using a model photo.



**Fig. 2.** Nine human body sizes selected for this study.

Nine items were cut in the column direction from the existing human body size data, and rows with missing values were removed. Afterwards, the data were scaled so that the mean was 0 and variance was 1, followed by standardization.

Figs. 3 and 4 show box plots of the data distribution before and after standardization for men in their 20s. The nine items mentioned above were sequentially expressed from left to right.

After the above process, clustering was performed according to body fat percentage using the K-means algorithm. Fig.

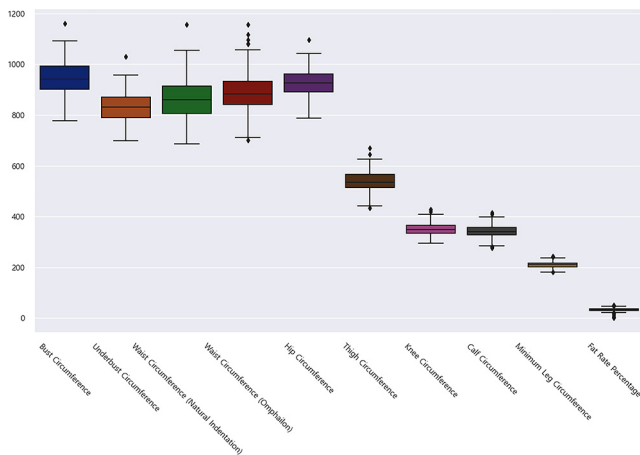


Fig. 3. Box plot for men in their 20s before standardization.

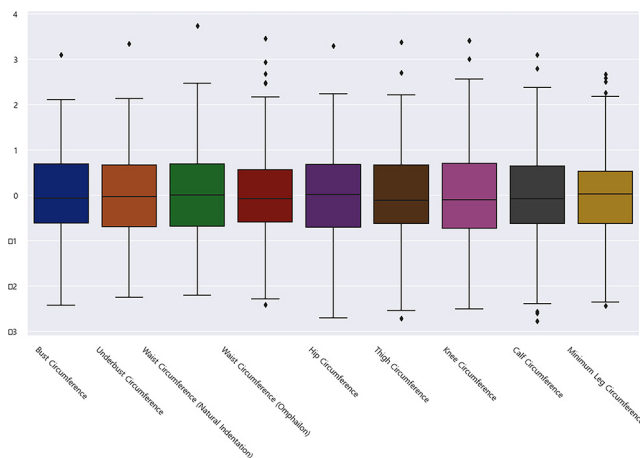


Fig. 4. Box plot for men in their 20s after standardization.

```

Input: Sample point  $X$ , Number of clusters  $k$ , Maximum of iteration  $i$ 
Output: Centroid  $C$ 
    Initialization : At  $X$ , randomly select  $k$  centroids
                    as the initial cluster centroid  $\mu^{(j)}, j \in \{1, \dots, k\}$ 
                     $C_1, \dots, C_k \leftarrow \emptyset$ 
1: for  $a = 1$  to  $i$  do
2:   for each  $x \in X$  do
3:      $j = \operatorname{argmin}_{j=1, \dots, k} \|x - \mu_j\|^2$ 
4:      $C_j \leftarrow C_j \cup x$ 
5:   end
6:   for  $j = 1$  to  $k$  do
7:     if  $C_j \neq \emptyset$  then  $\mu^{(j)} = \frac{1}{|C_j|} \sum_{x \in C_j} x$ 
8:   end
9: end
    
```

Fig. 5. K-means clustering algorithm.

5 shows the pseudocode of the K-means algorithm [6].

The elbow method was used to determine the value of  $k$ . Fig. 6 shows the results when only men in their 20s were considered; similar results were shown for other age groups and genders.

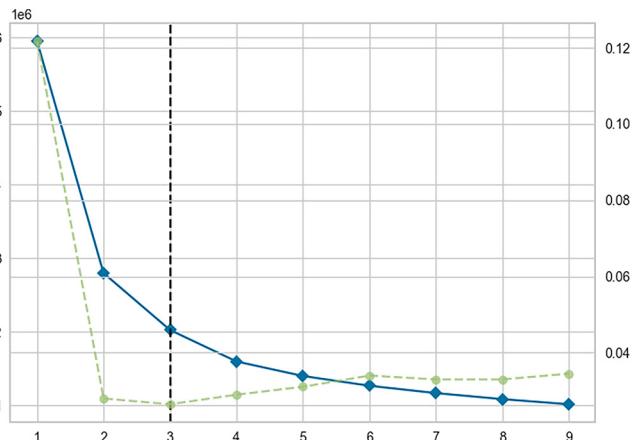


Fig. 6. Value of K when the elbow method is applied to men in 20s

Table 2. Clustering results

	Cluster number	Men's fat rate	Women's fat rate
20	Lowest	13.7741935483	24.3294478527
	Low	16.1882575757	27.9261682242
	Average	19.6666666666	31.6592039800
	High	23.8366197183	34.8061538461
	Highest	29.2052631578	42.7833333333
30	Lowest	14.9828124999	24.7364197530
	Low	18.3919540229	28.4421874999
	Average	20.2247311827	31.4851190476
	High	23.0694267515	35.4985714285
	Highest	27.6900000000	40.4210526315
40	Lowest	16.5818181818	26.3729411764
	Low	19.3210526315	29.4008771929
	Average	20.9579439252	32.4777777777
	High	24.5845070422	36.7283018867
	Highest	28.1173913043	42.7222222222
50	Lowest	<b>17.275</b>	28.2071428571
	Low	<b>18.6491525423</b>	<b>31.0783132530</b>
	Average	21.3803571428	<b>32.9857142857</b>
	High	<b>24.2875000000</b>	35.22125
	Highest	<b>25.0644444444</b>	37.6268292682
60	Lowest	16.3800000000	29.0653061224
	Low	21.2761904761	<b>33.1511627906</b>
	Average	23.2384615384	<b>33.7983870967</b>
	High	26.1259259259	<b>34.7949999999</b>
	Highest	28.0785714285	39.6866666666

However, if the group is divided into three groups, such as a cluster representing below average body fat percentage, average cluster, and cluster representing above average, results with a large amount of data in one cluster inevitably occur, making detailed classification difficult.

In this study, we did not use the optimal value but divided

the values into five categories from the lowest to the highest case; the corresponding results are reported in Table 2.

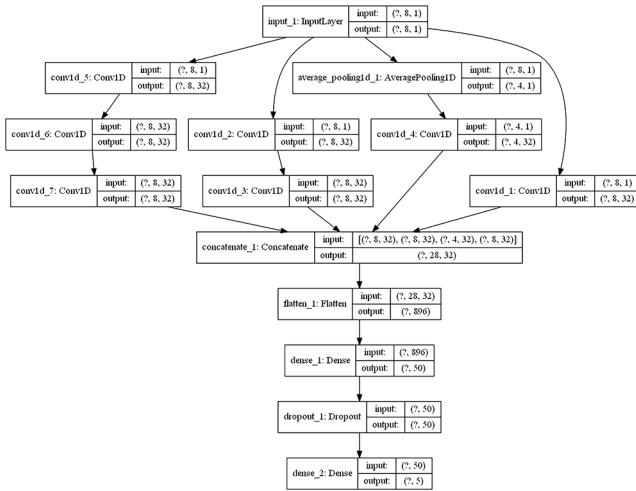
From the cluster results reported in Table 2, observe that when the average body fat percentage in each cluster is calculated, men and women in their 20s to 40s are evenly distributed.

However, after 50 years of age, the difference in the average body fat percentage interval corresponding to each group was not significant. The values represented in bold in Table 2 indicate cases where there was a difference of 1 or less. These results are because of the fact that the number of samples from the 50s onwards was much smaller than the number of people represented in the 7th body size dataset reported in Table 1.

### C. CNN Classification Model Design

A CNN classification model was used to transform the Inception-v2 module developed by Google to fit the dimensions of the human body size dataset [10]. The reason for using this model is that it suggests extracting features from different hierarchies from different stems and combining them again, which is widely used.

Fig. 7 shows the model modified to fit the dimensions of the human body data to the Inception-v2 module.



**Fig. 7.** CNN classification model proposed in this study using the inception-v2 module.

As shown in Fig. 7, the Inception module can be imported and used directly from the Keras library. This is because the input data dimensions must match the preprocessed data dimensions.

For the input layer (?, 8, 1), eight items of data should be given as input; in the case of men, this means eight items of

the human body size without the data for under bust circumference. For women (?, 9, 1) all nine items of human body size are used.

The first flow from the left is a section in which the features are extracted across three convolutional layers by extending them to 32 dimensions. The second flow is an interval in which features are extracted across two convolutional layers by extending them to 32 dimensions. The third flow is an average pooling layer to calculate the statistics and reduce the number of features by half to 32. This is a section that increases in dimensions. In the fourth flow, features are extracted by extending them to 32 dimensions using only one convolutional layer.

Because all four flows mentioned above have the same output dimension of 32, they are combined (concatenation layer) into one dimension (flatten layer) and then go through the dense layer using only 50 nodes. The final output dimensions became (?, 5) to classify the final five items using the dropout layer, which was used by leaving only 0.7% of the nodes. All convolutional layers used a rectified linear unit (ReLU).

Because this was a multi-classification task, the SoftMax function was used as the final activation function. Assuming that is the number of neurons in the output layer, is the input signal from the previous layer, and is the  $t$ th output, the SoftMax function can be defined as stated in Eq. (1).

$$y_k = \frac{\exp(a_k)}{\sum_{i=1}^n \exp(a_i)} \quad (1)$$

### D. Setting Model Hyperparameters

For the CNN classification model, the categorical cross-entropy error was used as the loss function. Assuming that the number of data dimensions is  $k$ , value estimated by the neural network is  $y_k$ , and target and correct answer is  $t_k$ , the cross-entropy function can be defined as expressed in Eq. (2).

$$CEE = -\sum_k t_k \log y_k \quad (2)$$

Adam was used as the optimizer where the number of training iterations (epochs) was set to 100 for both training and validation. The batch size was set to 1000 for both training and validation.

For the model evaluation, the training and validation dataset proportion was 80 to 20, and a random shuffle was used. The Stratify technique was used to prevent the class ratio from being biased toward either training or validation datasets. Table 3 lists the hardware specifications.

**Table 3.** Hardware specifications

<b>CPU</b>	AMD Ryzen 7 2700x (Core x8, 3.70 Ghz)
<b>RAM</b>	16 GB
<b>GPU</b>	Nvidia Geforce GTX 1660ti
<b>Storage</b>	SSD 500 GB / HDD 2TB
<b>OS</b>	Windows 10
<b>Virtual Environment</b>	Anaconda 3

## IV. MODEL TRAINING RESULTS

In this section, we discuss the performance indicators, considerations, and inference results of the final training results of the model.

### A. Performance Indicators

Accuracy, precision, recall, and final loss were used as performance indicators. The accuracy, precision, and recall are defined as expressed in Equations 3 and 4. The loss values were previously mentioned in Section III-D.

**Table 4.** Final learning results

		Accuracy	Precision	Recall	Loss	Times [s]
<b>20</b>	Men	0.8745	0.8962	0.8600	0.2943	5.52
	Women	0.9015	0.9249	0.8864	0.2359	3.29
<b>30</b>	Men	0.9173	0.9324	0.9019	0.2217	3.31
	Women	0.9259	0.9451	0.9241	0.1955	3.13
<b>40</b>	Men	0.8992	0.9053	0.8871	0.2450	3.28
	Women	0.8819	0.9164	0.8750	0.2681	3.14
<b>50</b>	Men	0.8864	0.9048	0.8636	0.2709	3.65
	Women	0.8951	0.9197	0.8811	0.2711	3.26
<b>60</b>	Men	0.9204	0.9174	0.8850	0.2676	3.33
	Women	0.9344	0.9543	0.9126	0.2138	3.38

**Table 5.** Final verification results

		Accuracy	Precision	Recall	Loss
<b>20</b>	Men	0.9770	0.9770	0.9770	0.1079
	Women	0.9548	0.9548	0.9548	0.1408
<b>30</b>	Men	0.9694	0.9694	0.9694	0.1018
	Women	0.9555	0.9555	0.9555	0.1014
<b>40</b>	Men	0.9354	0.9354	0.9354	0.1675
	Women	0.9166	0.9166	0.9166	0.1553
<b>50</b>	Men	0.9772	0.9772	0.9772	0.1405
	Women	0.9861	0.9861	0.9861	0.0893
<b>60</b>	Men	1.0	1.0	1.0	0.1245
	Women	0.9347	0.9318	0.8913	0.1489

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN} \tag{3}$$

$$Precision = \frac{TP}{TP + FP} \quad Recall = \frac{TP}{TP + FN} \tag{4}$$

## B. Learning Results

Table 4 reports the final learning results for both men and women in their 20s to 60s according to the performance indicators. The last column represents the learning time.

Table 5 summarizes the verification results based on the performance indicators.

## C. Discussion

As a result of the proposed model training, when averaged by performance index, observe that the accuracy is 0.90366, precision is 0.92165, recall is 0.88768, loss value is 0.249, and learning time is 3.529 seconds.

When learning was verified correctly, the average values were 0.96067, 0.96038, 0.95633, and 0.12779 for accuracy, precision, recall, and loss, respectively.

This result indicates that clustering was performed correctly when labeling was performed with the clustering algorithm using the nine items of the human body size data. Subsequently, the average value of the corresponding cluster was calculated.

However, a correct classification cannot be achieved because the human body size data, which are the learning data, do not properly cluster the average body fat percentage in the groups classifying participants in their 50s to 60s compared to those classifying the participants in their 20s to 40s.

## D. Inference Results

The cluster model files corresponding to the previously learned age and gender were saved as Pickle files, and the CNN classification model files were saved as H5 files so that inferences could be made separately. The average body fat percentage of each group was divided based on age and gender and saved in a CSV file.

The environment in which the user inputs data and inference result can be checked was implemented based on the Flask Web framework. The processing procedure comprised only the loading and execution of the aforementioned files. Fig. 8 shows the form used by the user to input the data, and Fig. 9 shows the corresponding results.

Body shape analysis form (test)

age:

gender:

bust circumference:

underbust circumference (women only):

waist circumference (natural indentation):

waist circumference (omphailon, belly-based):

hip circumference:

thigh circumference:

knee circumference:

calf circumference:

minimum leg circumference:

Fig. 8. Human body size input form using flask.

Cluster results according to your gender and age: Lowest body fat percentage  
(Average body fat percentage for the same gender and age as you:  
13.774193548387096)

Fig. 9. Result for the test represented in Fig. 8.

## V. CONCLUSIONS

In this study, easy-to-measure data comprising nine items, such as the chest and waist measurements, were utilized through human body size dataset, which is a public dataset provided by the National Institute of Technology and Standards. Based on age and gender, clustering was performed using the K-means algorithm, and the average body fat percentage of each cluster was calculated to obtain meaningful clustering results.

Based on these results, a model was proposed to use the Inception-v2 module to obtain five classification groups from the highest body fat percentage to the lowest body fat percentage by modifying the Inception-v2 module according to the data shape. As a result of the verification, accuracy, precision, and recall of more than 95% were obtained.

There was no significant difference between the results of the clusters in the 50s and 60s groups compared with the results of the clusters between those in their 20s and 40s.

Future research should supplement the current results and use a better clustering algorithm or CNN architecture. In addition, it is believed that rational decision-making will be possible in the field of healthcare or obesity analysis using the device through 3D scanning for circumference measurement.

## ACKNOWLEDGMENTS

This research was supported by a 2023 Pai Chai University research grant.

## REFERENCES

- [1] Smart & Company Ltd, Convergence of artificial intelligence and healthcare, growing the technology commercialization market [Internet], Available: <http://m.elec4.co.kr/article/articleView.asp?idx=26993>.
- [2] J. Kim, B. Kang, and H. Jung, "Determination of coagulant input rate in water purification plant using K-means algorithm and GBR algorithm," *Journal of the Korea Institute of Information and Communication Engineering*, vol. 25, no. 6, pp. 792-798, Jun. 2021. DOI: 10.6109/jkiice.2021.25.6.792.
- [3] P.-H. Huynh, V. H. Nguyen, and T.-N. Do, "Enhancing gene expression classification of support vector machines with generative adversarial networks," *Journal of Information and Communication Convergence Engineering*, vol. 17, no. 1, pp. 14-20, Mar. 2019. DOI: 10.6109/jiice.2019.17.1.14.
- [4] R. Chefira and S. Rakrak, "A knowledge extraction pipeline between supervised and unsupervised machine learning using gaussian mixture models for anomaly detection," *Journal of Computing Science and Engineering*, vol. 15, no. 1, pp. 1-17, Mar. 2021. DOI: 10.5626/JCSE.2021.15.1.1.
- [5] J. Li, G. Huang, and Y. Zhou, "A sentiment classification approach of sentences clustering in webcast barrages," *Journal of Information Processing Systems*, vol. 16, no. 3, pp. 718-732, Jun. 2020. DOI: 10.3745/JIPS.04.0174.
- [6] J. Kim, E. Park, K. Han, J. Lee, and H. J. Lee, "A two-stage learning method of cnn and k-means rgb cluster for sentiment classification of images," *Journal of Intelligence and Information Systems*, vol. 27, no. 3, pp. 139-156, 2021. DOI: 10.13088/jiis.2021.27.3.139.
- [7] Korea Agency for Technology and Standards, Size korea (korean human body size measurements) [Internet], Available: <https://size.korea.kr/>.
- [8] E.-K. Kim and Y. R. Nam, "Analysis of lower body shape of men in their 30s for pants pattern designs – focus on changes in human dimensions and body type classification," *Journal of the Korea Fashion & Costume Design Association*, vol. 23, no. 2, pp. 133-146, 2021. DOI: 10.30751/kfcd.2021.23.2.133.
- [9] C. S. Joung, "A study on the types of upper body shape of male university students," *Journal of The Korean Society Design Culture*, vol. 25, no. 1, pp. 453-463, Mar. 2019. DOI: 10.18208/ksdc.2019.25.1.453.
- [10] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, USA, pp. 2818-2826, 2016. DOI: 10.1109/CVPR.2016.308.



**Taejun Lee**

He received his bachelor's degree and master's degree in computer science from Pai Chai University in 2020 and 2023, respectively. Since 2023, he has been pursuing a Ph.D. at Pai Chai University. His current research interests are deep learning, machine learning, big data, and computer vision.



**Hakseong Kim**

He received his bachelor's degree in electronics from Jeon Buk University in 1998. Then, he received his master's degree in information and communication engineering from Chung Nam University in 2010. His current research interests are deep learning, machine learning, big data, AI, and instrumentation & control.



**Hoekyung Jung**

He received his master's degree in 1987 and Ph.D. in 1993 from the Department of Computer Engineering of Kwangwoon University, Korea. From 1994 to 1995, he worked for ETRI as a researcher. Since 1994, he has worked in the Department of Computer Engineering at Pai Chai University, where he now works as a professor. His current research interests include multimedia document architecture modeling, information processing, embedded system, machine learning, big data, and IoT.