

ResNet50 전이학습을 활용한 손동작 인식 기반 가위바위보 게임 구현

박창준*, 김창기⁰, 손성규*, 이경진*, 유희경*, 곽정환(교신저자)*

*한국교통대학교 AI로봇공학과,

⁰한국교통대학교 AI로봇공학과

e-mail: {1726018, 1715086, 1915153, 1926034, jgwak}@ut.ac.kr*, 1715071@ut.ac.kr⁰

Implementation of hand motion recognition-based rock-paper-scissors game using ResNet50 transfer learning

Changjoon Park*, Changki Kim⁰, Seongkyu Son*, Kyoungjin Lee*,

Heekyung Yoo*, Jeonghwan Gwak(Corresponding Author)*

*Dept. of AI Robotics, Korea National University of Transportation,

⁰Dept. of AI Robotics, Korea National University of Transportation

● 요약 ●

GUI(Graphical User Interface)를 대신하는 차세대 인터페이스로서 NUI(Natural User Interface)에 기대가 모이는 것은 자연스러운 흐름이다. 본 연구는 NUI의 손가락 관절을 포함한 손동작 전체를 인식시키기 위해 웹캠과 카메라를 활용하여 다양한 배경과 각도의 손동작 데이터를 수집한다. 수집된 데이터는 전처리를 거쳐 데이터셋을 구축하며, ResNet50 모델을 활용하여 전이학습한 합성곱 신경망(Convolutional Neural Network) 알고리즘 분류기를 설계한다. 구축한 데이터셋을 입력시켜 분류학습 및 예측을 진행하며, 실시간 영상에서 인식되는 손동작을 설계한 모델에 입력시켜 나온 결과를 통해 가위바위보 게임을 구현한다.

키워드: 합성곱 신경망(CNN), ResNet50, 전이학습, 손동작 인식, 가위바위보 게임

I. Introduction

최근 PC와 VR, AR 등과 융합된 NUI(Natural User Interface)[1]는 눈동자 인식, 손동작 인식 등 동작 인식이 가능하다. 인간 손에는 손가락 관절이 존재하므로 손동작은 눈동자로 표현할 수 있는 동작보다 수가 많으며 효율적이다. 따라서 손동작 인식은 마우스나 키보드 같은 부가적인 장치 없이 사용자와 컴퓨터 간의 상호작용을 원활하게 할 수 있는 효과적인 방법이다.

본 논문에서는 손가락 관절을 포함한 손 전체 동작을 인식시키기 위하여 시각적 이미지를 분석하는 데 쓰이는 인공 신경망 중 하나인 합성곱 신경망(Convolutional Neural Network)[2] 알고리즘을 사용한다. 이는 다차원 배열 데이터 처리가 가능하며 일반적인 신경망과는 다르게 이미지에서 특징을 추출해 처리하므로 컬러 이미지 처리에 특화되어 있다. 또한 부족한 데이터셋을 보완하기 위해 전이학습을 진행하고 방대한 이미지를 통해 학습한 ResNet50 모델[3]을 활용하여 합성곱 신경망을 설계한다. 이러한 합성곱 신경망은 다음과 같은 과정으로 진행된다.

첫째, 웹캠 및 카메라를 통해 여러 각도와 배경으로 영상을 수집하고 프레임 간격으로 이미지화하여 라벨을 붙여줌으로써 데이터셋을 구축한다. 둘째, 구축된 데이터셋은 ResNet50 모델을 사용하여 전이학습한 합성곱 신경망을 활용하여 분류학습을 진행한다. 셋째, 웹캠을 통해 실시간으로 입력받은 손동작 이미지를 분류된 클래스에 맞게 예측이 진행되는지를 평가한다. 마지막으로 설계된 모델에 손동작을 인식하여 나온 예측 결과들을 통해 가위바위보 게임[4, 5]을 구현한다.

II. Preliminaries

2.1 CNN (Convolutional Neural Network)

CNN[3]은 Deep Learning 알고리즘 중 하나로 이미지 처리에 뛰어난 성능을 보이는 신경망이다. CNN의 기본 구조는 Convolution Layer와 Pooling Layer를 반복적으로 적용하여 여러 겹 쌓는 특징 추출 부분, Fully Connected Layer를 구성하고 출력층에 Softmax를

적용한 분류 부분으로 구성된다.

학습 과정은 다음과 같이 진행하였다. 먼저 필터를 입력 이미지에 교차 적용하여 곱 연산을 하고 활성화 함수를 거친 후 결과 값을 추출한다. 결과 값은 1차원 배열 데이터로 변환되어 Fully Connected Layer를 거쳐 분류된다. 본 논문에서는 Input_image로 여러 배경에서 촬영한 가위바위보 영상을 224*224 이미지로 데이터 전처리하여 사용하며, 흔들림이 있는 Noise Data는 제거한다. 또한 총 4개의 분류를 얻기 위해 Softmax 출력값을 조정한다.

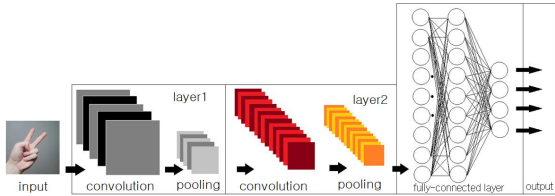


Fig. 1. Convolution Neural Network의 기본적인 구조



Fig. 2. ResNet50기반 전이학습

2.2 ResNet50기반 전이학습

관련 작업에 대해 이미 훈련된 모델의 일부를 가져와 새 모델에서 재사용함으로써 손쉽게 데이터를 구축하는 기술이다. 주로 학습데이터가 부족한 분야의 모델 구축을 위해 사용한다. ResNet50은 Image Net 데이터베이스의 백 만개가 넘는 영상에 대해 훈련한 50개의 층을 갖는 CNN 모델이다. 본 논문에서는 ResNet50의 가중치를 불러와 층을 추가하고 가위바위보 분류기를 학습한다. 이번 연구에서 구축한 입력 데이터의 수가 가위바위보 행동별 3천 장으로 학습을 원활히 진행하기 힘들다고 판단하였다. 따라서 ResNet50을 활용한 전이학습을 통해서 입력 데이터 부족으로 인한 모델 성능 문제를 해소한다.

III. The Proposed Scheme

본 연구에서는 가위바위보 손동작의 분류를 위해 데이터를 수집 및 전처리하는 과정을 거친다. 이때 생성된 데이터셋을 합성곱 신경망에 입력하여 손동작을 분류할 수 있도록 학습시키고 실제 데이터를 분류한다. 또한 가위바위보 알고리즘을 구성하고 게임을 직관적으로 진행할 수 있도록 python의 tkinter 모듈을 이용하여 GUI를 구성한다. 연구는 1080p camera, CPU-i9 10900k, GPU-RTX3090, RAM-16GB, cuda11.4, cudnn11.4, python3.9.5와 python의 tkinter, tensorflow-gpu2.5, OpenCV4.5.4 환경에서 진행한다.

1. 데이터수집 및 전처리

가위바위보에 대한 손동작을 합성곱 신경망 모델에 학습시키기 위한 데이터 수집 및 전처리하는 단계이며 데이터의 수집은 카메라(스마트폰 및 웹캠)를 통해 진행된다.

가위바위보의 데이터셋 종류는 가위 1(엄지와 검지), 가위 2(검지와 중지), 바위, 보 총 4개로 구성된다. 해당 4개의 손동작을 카메라 앵글 중앙에 각각 비추어 영상을 촬영 및 저장하고, 위 영상 데이터를 OpenCV를 이용하여 10프레임당 1장의 이미지 데이터로 변환 및 데이터셋화 한다. 해당 데이터셋을 합성곱 신경망 모델에 입력하기 위해 영상의 중앙에서부터 정사각형 모양으로 이미지를 자르고 벡터의 차원을 224*224*3의 크기로 변환하여 저장한다.

남은 전처리는 ResNet50 전이학습 모델에서 지원하는 preprocess_input 함수를 이용한다. 가위 1, 가위 2, 바위, 보의 라벨당 데이터 개수는 (표 1)과 같이 각 라벨마다 3천 개로 총 1만 2천 개의 데이터로 구성되며, 전체 데이터의 80%는 특징 데이터, 나머지 20%는 검증 데이터로 나누어 모델 학습에 사용한다. 손동작 영상은 합성곱 신경망 모델의 편향 학습을 막기 위해 다양한 손의 각도와 배경에서 (그림 3)과 같이 촬영한다.

Table 1.

	가위1	가위2	바위	보	전체
특징	2,400	2,400	2,400	2,400	9,600
검증	600	600	600	600	2,400
전체	3,000	3,000	3,000	3,000	12,000



Fig. 3.

2. 합성곱 신경망(CNN) 모델 제작 및 학습

현 단계에서는 앞서 구축한 손동작 이미지 데이터셋을 분류하기 위해 합성곱 신경망 모델을 생성하고 모델 학습을 진행한다. 해당 모델은 ResNet50의 예측률이 뛰어난 가중치를 이용해 합성곱 신경망 분류기를 전이 학습시켜 예측률을 높인다.

먼저 tensor flow에서 제공하는 ResNet50 모델을 불러와 weights='imagenet' 및 include_top=false 옵션을 적용하여 ResNet50의 가중치인 imagenet을 학습에 사용하고 말단의 Fully-connected layer를 제거한다. 이후 앞서 정의한 손동작 데이터셋을 구별할 수 있도록 4개의 class를 출력하는 합성곱 신경망 분류기를 만든다. 분류기는 Flatten layer/Fully-connected layer (activation='relu')/output layer (activation='softmax') 순서로 구성된다. 합성곱 신경망 분류기는 (그림 4)와 같이 Fully-connected

있는 알고리즘을 간단하게 구현하였으며, 상대방의 가위바위보 라벨 값을 랜덤하게 적용된다.

자신이 카메라로 촬영하고 있는 손동작 데이터를 직관성 있게 확인할 수 있도록 OpenCV를 이용하여 카메라에 비추고 있는 손의 모습이 보일 수 있게 영상 처리를 하고, 합성곱 신경망 분류기가 예측한 라벨을 영상 위에 출력한다. 또한 스트리밍 데이터를 예측할 때 손동작을 잘못 인식할 수 있으므로 게임 시작 버튼이 작동한 후 손동작 데이터를 3초간 10개의 이미지를 처리 및 예측한다. 이후 예측한 확률값을 모두 더하여 확률이 제일 높은 라벨을 최종 예측 라벨로 선정한다. 선정된 라벨값과 랜덤하게 정해진 상대방 라벨 정보를 가위바위보 승패 알고리즘에 입력한 후 출력값을 중앙에 띄운다. (그림 11)은 가위바위보 게임을 실행한 초기 화면이다.



Fig. 11. 가위바위보 게임 실행 영상

IV. Conclusions

개발한 게임을 테스트하기 위해 배경과 각도를 달리하여 게임을 진행한다. 배경은 흰 배경과 모니터 상에서 띄운 풍경 사진을 배경으로 변화를 주고 각도는 정면, 윗면, 아랫면, 손바닥, 손등, 대각선 총 6개의 각도로 구분하여 테스트한다. 웹캠에 가위바위보 행동을 인식하고 하단부에 start 버튼을 누르면 3초 뒤에 컴퓨터가 무작위로 낸 가위바위보와 비교하여 화면에 결과를 표시한다. 테스트를 계속한 결과, 가위바위보 중 각도에 구애받지 않는 바위의 인식이 가장 높았고 가위는 손바닥 각도일 때 가장 인식이 높았지만, 화면상 앞에 위치한 손가락에 의해 뒤의 손가락이 가려지는 아랫면, 윗면 각도에서는 가위와 보를 명확히 구분하지 못하는 결과를 확인했다. 보도 마찬가지로 정면, 아랫면, 윗면에서 인식을 저하를 확인했고 손가락 사이의 거리를 크게 할수록 인식이 높아지는 것을 알 수 있다.



Fig. 12-1. 하얀 배경에서 연구한 결과-가위



Fig. 12-2. 하얀 배경에서 연구한 결과-바위

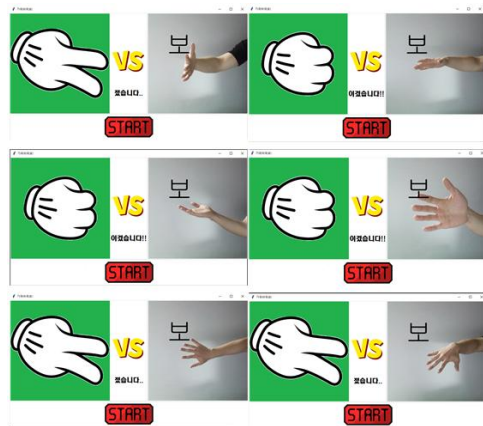


Fig. 12-3. 하얀 배경에서 연구한 결과-보



Fig. 13-1. 배경 변화를 주며 연구한 결과가위



Fig. 13-2. 배경 변화를 주며 연구한 결과바위



Fig. 13-3. 배경 변화를 주며 연구한 결과보

더불어 모델 학습 부분에서 가위의 데이터셋을 총 2개(가위 1, 가위 2)로 구분했으나, (그림 14)를 통해 학습시키지 않은 중지와 약지, 약지와 소지 형태의 가위도 인식하는 모습을 확인할 수 있었다.



Fig. 14-1. 중지와 약지



Fig. 14-2. 약지와 소지

본 논문에서는 합성곱 신경망 모델을 이용하여 손동작을 인식하는 알고리즘 개발 및 NUI를 적용한 가위바위보 게임 시스템을 제안한다. 가위 1, 가위 2, 바위, 보를 각 3천 장, 총 1만 2천 장으로 구성된 데이터셋과 ResNet50의 가중치를 활용하여 합성곱 신경망 분류기의 학습을 진행한다.

해당 분류기 학습결과, 가위바위보를 분류하는 합성곱 신경망 모델은 약 96%의 정확률을 얻었다. 또한 바위의 인식률이 가장 높았고 가위의 인식률이 가장 낮음을 확인한다. 즉 가위와 보, 정면에서 바라보는 가위나 바위 등 손가락 관절이 겹치는 부분을 미세하게 인식하지 못하는 모습을 보인다. 해당 부분에 있어 인식률을 높이는 방향으로 학습모델을 개선하기 위해서는 기존보다 더 많은 양의 이미지 데이터가 필요하므로 더욱 다양한 각도와 배경에서의 이미지 데이터를 수집해야 한다. 이러한 방법을 통해 더 나은 모델로 개선하면 미세한 손동작도 인식할 수 있으므로 손동작으로 제어하는 부분에 있어 낮은 오류를 보여줄 것이고, 기존 터치 중심의 GUI보다 더욱 큰 편리함을 가져와 준다.

또한 VR, AR 등에서도 손동작으로 여러 콘텐츠를 이용하는 부분에 있어서 높은 인식률로 부드럽고 세세한 손동작들을 수행할 수 있으므로 더 높은 질과 다양한 콘텐츠를 구현해낼 것으로 기대된다.

ACKNOWLEDGEMENT

This results was supported by "Regional Innovation Strategy (RIS)" through the National Research Foundation of Korea(NRF) funded by the Ministry of Education(MOE) (2021RIS-001).

REFERENCES

- [1] Hee-Sun Choi, "A Study on NUI with Fingertip Detection and Tracking", Chung-Ang University Seoul Campus Academic Information Center, pp.12~24, 2012
- [2] Hyun-soo Lee, "A Structure of Convolutional Neural Networks for Image Contents Search" Chung-Ang University Seoul Campus Academic Information Center, pp.1~14, 2018
- [3] Sung-Wook Park, Do-Yeon Kim, "Comparison of Image Classification Performance in Convolutional Neural Network according to Transfer Learning", Journal of Korea Multimedia Society Vol. 21, No. 12, December 2018, pp.1388-1393, 2018
- [4] Jeong-min Seo, Byeong-ju Kim, Hyeong-man Moon , Chang-sun Park, Jung-hwan Hwang, "Paper-scissors-stone Game System Used by Hand Image Recognition Technique" Korea Multimedia Society, pp.510~512 ,2010
- [5] Yeon-Su Jang, Da-Ye Kim, Dong-Jin Park, YunSung Han, Soobin Jeon, Dongmahn Seo, "A Rock-paper-scissors Game Using Hand Image Recognition Technology based on Artificial Neural Network", Korea Information Processing Society, pp.659-662, 2020