

# Gesture Classification with Machine Learning using Kinect Sensor Data

Sambit Bhattacharya, Bogdan Czejdo and Nicolas Perez

Department of Mathematics and Computer Science

Fayetteville State University

Fayetteville, NC, USA

{sbhattac, bczejdo, nperez3}@uncfsu.edu

**Abstract**— We present approaches for gesture classification and gesture segmentation by using machine learning on the Kinect sensor's data stream. Our work involved three phases. Firstly we developed gesture classification from a known vocabulary of gestures in an edited data stream. Secondly we extended those techniques to detect and classify a gesture in an unedited stream which also captures random movements. Thirdly, we apply rules to filter out movements that were not intentional gestures and yet resembled certain gestures in our vocabulary.

**Keywords**- human gesture; machine learning; gesture classification; gesture recognition; gesture segmentation; Kinect sensor

## I. INTRODUCTION

Recognition and classification of human gestures is an active area of research since it is an important component of human-machine interaction. Most existing approaches use low level image processing operators to extract salient features from the moving human body in video data, which are then used to train machine learning algorithms [1,2]. Such feature extraction techniques pose a significant challenge and work under certain restrictive assumptions such as lighting conditions and even the hardware design of the video capturing device [1]. Gesture segmentation also remains a problem with multiple approaches described in the literature [2]. The Microsoft Kinect offers a hardware and software platform which solves some key computer vision problems with high levels of accuracy and well understood error characteristics and limitations [3]. In this paper we present our approach in solving the problem of gesture recognition by using machine learning on the Kinect's data stream of locations of joints on a human body, in other words skeletal data. Our approach consists of three parts. Firstly we develop gesture classification from a known vocabulary of gestures in an edited data stream. Secondly we extend those techniques to detect and classify a deliberate gesture in an unedited stream which also captures random movements. Thirdly, we apply rules to filter out movements that were not intentional gestures and yet resembled certain gestures in our vocabulary; in other words we developed techniques for both the gesture segmentation and classification problem in an unedited stream.

This work was partially supported by National Science Foundation (award id: 0959958) and an internal grant award of Fayetteville State University's Integrated STEM Academic Success (ISAS) program.

## II. THE KINECT JOINT DATA STREAM AND GESTURE VOCABULARY

### A. The Kinect Sensor and its Data Streams

Here we describe in brief the main functions of Kinect sensor. We next describe how the Kinect data stream is read and how we modify the data organization for processing purposes.

The Kinect sensor can measure depth data (estimates of distance from Kinect to pixels in scenery) and it can identify and locate human skeleton. This skeleton is identified through positions of skeletal joints computed by Kinect from depth data. The Kinect processes and sends the skeletal data along with the depth data at the rate of 30 frames a second. It provides data streams for skeletal joints of up to two humans.

The precision of Kinect measurements and distortions was extensively studied in [3]. This experimental study reports that the random error of depth measurements increases quadratically with increasing distance from the sensor and the depth resolution also decreases quadratically with increasing distance from the sensor. It also measured the error bounds and resolutions limits and recommends that for mapping applications the data should be acquired within 1–3 m distance to the sensor. Our experiments were performed with sufficient variations within these limitations.

The Kinect sensor sends skeletal data to the computer as a temporal sequence of X, Y, Z coordinates of all 20 tracked joints. The data is grouped into time-stamped frames and the  $i^{\text{th}}$  frame can be represented as  $C_i^{1x} C_i^{1y} C_i^{1z} \dots C_i^{20x} C_i^{20y} C_i^{20z}$ , where  $C$  is a single coordinate value and the super-script has joint number followed by axis letter. When recording a data stream that starts at the  $0^{\text{th}}$  frame for which  $n+1$  frames have been seen so far, we organize the values as  $C_0^{1x} C_1^{1x} C_2^{1x} \dots C_n^{1x}$   $C_0^{1y} C_1^{1y} C_2^{1y} \dots C_n^{1y}$   $C_0^{1z} C_1^{1z} C_2^{1z} \dots C_n^{1z}$   $C_0^{2x} C_1^{2x} C_2^{2x} \dots C_n^{2x}$   $C_0^{2y} C_1^{2y} C_2^{2y} \dots C_n^{2y}$   $C_0^{2z} C_1^{2z} C_2^{2z} \dots C_n^{2z}$   $C_0^{3x} C_1^{3x} C_2^{3x} \dots C_n^{3x}$   $C_0^{3y} C_1^{3y} C_2^{3y} \dots C_n^{3y}$   $C_0^{3z} C_1^{3z} C_2^{3z} \dots C_n^{3z}$   $C_0^{4x} C_1^{4x} C_2^{4x} \dots C_n^{4x}$   $C_0^{4y} C_1^{4y} C_2^{4y} \dots C_n^{4y}$   $C_0^{4z} C_1^{4z} C_2^{4z} \dots C_n^{4z}$   $C_0^{5x} C_1^{5x} C_2^{5x} \dots C_n^{5x}$   $C_0^{5y} C_1^{5y} C_2^{5y} \dots C_n^{5y}$   $C_0^{5z} C_1^{5z} C_2^{5z} \dots C_n^{5z}$   $C_0^{6x} C_1^{6x} C_2^{6x} \dots C_n^{6x}$   $C_0^{6y} C_1^{6y} C_2^{6y} \dots C_n^{6y}$   $C_0^{6z} C_1^{6z} C_2^{6z} \dots C_n^{6z}$   $C_0^{7x} C_1^{7x} C_2^{7x} \dots C_n^{7x}$   $C_0^{7y} C_1^{7y} C_2^{7y} \dots C_n^{7y}$   $C_0^{7z} C_1^{7z} C_2^{7z} \dots C_n^{7z}$   $C_0^{8x} C_1^{8x} C_2^{8x} \dots C_n^{8x}$   $C_0^{8y} C_1^{8y} C_2^{8y} \dots C_n^{8y}$   $C_0^{8z} C_1^{8z} C_2^{8z} \dots C_n^{8z}$   $C_0^{9x} C_1^{9x} C_2^{9x} \dots C_n^{9x}$   $C_0^{9y} C_1^{9y} C_2^{9y} \dots C_n^{9y}$   $C_0^{9z} C_1^{9z} C_2^{9z} \dots C_n^{9z}$   $C_0^{10x} C_1^{10x} C_2^{10x} \dots C_n^{10x}$   $C_0^{10y} C_1^{10y} C_2^{10y} \dots C_n^{10y}$   $C_0^{10z} C_1^{10z} C_2^{10z} \dots C_n^{10z}$   $C_0^{11x} C_1^{11x} C_2^{11x} \dots C_n^{11x}$   $C_0^{11y} C_1^{11y} C_2^{11y} \dots C_n^{11y}$   $C_0^{11z} C_1^{11z} C_2^{11z} \dots C_n^{11z}$   $C_0^{12x} C_1^{12x} C_2^{12x} \dots C_n^{12x}$   $C_0^{12y} C_1^{12y} C_2^{12y} \dots C_n^{12y}$   $C_0^{12z} C_1^{12z} C_2^{12z} \dots C_n^{12z}$   $C_0^{13x} C_1^{13x} C_2^{13x} \dots C_n^{13x}$   $C_0^{13y} C_1^{13y} C_2^{13y} \dots C_n^{13y}$   $C_0^{13z} C_1^{13z} C_2^{13z} \dots C_n^{13z}$   $C_0^{14x} C_1^{14x} C_2^{14x} \dots C_n^{14x}$   $C_0^{14y} C_1^{14y} C_2^{14y} \dots C_n^{14y}$   $C_0^{14z} C_1^{14z} C_2^{14z} \dots C_n^{14z}$   $C_0^{15x} C_1^{15x} C_2^{15x} \dots C_n^{15x}$   $C_0^{15y} C_1^{15y} C_2^{15y} \dots C_n^{15y}$   $C_0^{15z} C_1^{15z} C_2^{15z} \dots C_n^{15z}$   $C_0^{16x} C_1^{16x} C_2^{16x} \dots C_n^{16x}$   $C_0^{16y} C_1^{16y} C_2^{16y} \dots C_n^{16y}$   $C_0^{16z} C_1^{16z} C_2^{16z} \dots C_n^{16z}$   $C_0^{17x} C_1^{17x} C_2^{17x} \dots C_n^{17x}$   $C_0^{17y} C_1^{17y} C_2^{17y} \dots C_n^{17y}$   $C_0^{17z} C_1^{17z} C_2^{17z} \dots C_n^{17z}$   $C_0^{18x} C_1^{18x} C_2^{18x} \dots C_n^{18x}$   $C_0^{18y} C_1^{18y} C_2^{18y} \dots C_n^{18y}$   $C_0^{18z} C_1^{18z} C_2^{18z} \dots C_n^{18z}$   $C_0^{19x} C_1^{19x} C_2^{19x} \dots C_n^{19x}$   $C_0^{19y} C_1^{19y} C_2^{19y} \dots C_n^{19y}$   $C_0^{19z} C_1^{19z} C_2^{19z} \dots C_n^{19z}$   $C_0^{20x} C_1^{20x} C_2^{20x} \dots C_n^{20x}$   $C_0^{20y} C_1^{20y} C_2^{20y} \dots C_n^{20y}$   $C_0^{20z} C_1^{20z} C_2^{20z} \dots C_n^{20z}$  where the subscript  $0, n$  means the sequence starting at 0 and ending at  $n$ . We additionally perform a data scaling step as recommended in [ref] for each  $C$  along the X, Y and Z directions for every joint as  $C_{\text{scaled}} = (C_{\text{original}} - \min) / (\max - \min)$  where  $\max$  and  $\min$  are the maximum and minimum values of that particular feature. For running the machine learning programs described next, we consider this sequence as a single feature vector where